



**ESPE**

**UNIVERSIDAD DE LAS FUERZAS ARMADAS**  
**INNOVACIÓN PARA LA EXCELENCIA**

**DEPARTAMENTO DE ELÉCTRICA, ELECTRÓNICA Y  
TELECOMUNICACIONES**

**CARRERA DE INGENIERÍA EN ELECTRÓNICA, AUTOMATIZACIÓN Y  
CONTROL**

**TRABAJO DE TITULACIÓN, PREVIO A LA OBTENCIÓN DEL TÍTULO  
DE INGENIERA EN ELECTRÓNICA, AUTOMATIZACIÓN Y CONTROL**

**TEMA: “DESARROLLO DE MODELOS MULTIVARIANTES PARA LA  
DISCRIMINACIÓN Y CUANTIFICACIÓN DE SUSTANCIAS  
EXPLOSIVAS”**

**AUTORA: DE LA CRUZ MOSQUERA, DINA KAROLAY**

**DIRECTORA: ING. GUAMÁN NOVILLO, ANA VERÓNICA PHD.**

**SANGOLQUÍ**

**2020**



**ESPE**  
UNIVERSIDAD DE LAS FUERZAS ARMADAS  
INNOVACIÓN PARA LA EXCELENCIA

**DEPARTAMENTO DE ELÉCTRICA, ELECTRÓNICA Y  
TELECOMUNICACIONES  
CARRERA DE INGENIERÍA EN ELECTRÓNICA, AUTOMATIZACIÓN Y  
CONTROL**

**CERTIFICACIÓN**

Certifico que el trabajo de titulación, “**DESARROLLO DE MODELOS MULTIVARIANTES PARA LA DISCRIMINACIÓN Y CUANTIFICACIÓN DE SUSTANCIAS EXPLOSIVAS**” fue realizado por la señorita **DE LA CRUZ MOSQUERA, DINA KAROLAY** el mismo que ha sido revisado en su totalidad, analizado por la herramienta de verificación de similitud de contenido; por lo tanto, cumple con los requisitos teóricos, científicos, técnicos, metodológicos y legales establecidos por la Universidad de Fuerzas Armadas ESPE, razón por la cual me permito acreditar y autorizar para que lo sustente públicamente.

Sangolqui, enero de 2020

  
.....  
**ING. ANA VERÓNICA GUAMÁN NOVILLO Ph.D.**  
C. C: 1103996946



**DEPARTAMENTO DE ELÉCTRICA, ELECTRÓNICA Y  
TELECOMUNICACIONES  
CARRERA DE INGENIERÍA EN ELECTRÓNICA, AUTOMATIZACIÓN Y  
CONTROL**

**AUTORÍA DE RESPONSABILIDAD**

Yo, **DE LA CRUZ MOSQUERA, DINA KAROLAY**, declaro que el contenido, ideas y criterios del trabajo de titulación: **“DESARROLLO DE MODELOS MULTIVARIANTES PARA LA DISCRIMINACIÓN Y CUANTIFICACIÓN DE SUSTANCIAS EXPLOSIVAS”** es de mi autoría y responsabilidad, cumpliendo con los requisitos teóricos, científicos, técnicos, metodológicos y legales establecidos por la Universidad de Fuerzas Armadas ESPE, respetando los derechos intelectuales de terceros y referenciando las citas bibliográficas. Consecuentemente el contenido de la investigación mencionada es veraz.

Sangolquí, enero de 2020



.....  
**DINA KAROLAY DE LA CRUZ MOSQUERA**

C. C: 1724472293



# ESPE

UNIVERSIDAD DE LAS FUERZAS ARMADAS  
INNOVACIÓN PARA LA EXCELENCIA

**DEPARTAMENTO DE ELÉCTRICA, ELECTRÓNICA Y  
TELECOMUNICACIONES  
CARRERA DE INGENIERÍA EN ELECTRÓNICA, AUTOMATIZACIÓN Y  
CONTROL**

**AUTORIZACIÓN**

Yo, **DE LA CRUZ MOSQUERA, DINA KAROLAY** autorizo a la Universidad de las Fuerzas Armadas ESPE publicar el trabajo de titulación: **“DESARROLLO DE MODELOS MULTIVARIANTES PARA LA DISCRIMINACIÓN Y CUANTIFICACIÓN DE SUSTANCIAS EXPLOSIVAS”** en el Repositorio Institucional, cuyo contenido, ideas y criterios son de mi responsabilidad.

Sangolqui, enero de 2020

.....  
**DINA KAROLAY DE LA CRUZ MOSQUERA**

C. C: 1724472293

# DEDICATORIA

*Dedicado a mi familia,  
gracias a su amor y apoyo pude cumplir  
uno de los sueños más importantes en mi vida*

# AGRADECIMIENTO

A la Dra. Ana Guamán, por su guía en la elaboración de este trabajo de titulación, sin sus sugerencias, consejos y toda la orientación que me brindo, el desarrollo de este proyecto no habría sido posible.

A mis abuelitas Dina y Graciela, quienes me enseñaron que una mujer es capaz de llegar tan lejos como se lo proponga si lo hace de forma honesta y con el corazón, en ustedes encontré un gran ejemplo de mujeres fuertes e independientes.

A mis padres, Eddy y Ximena, por todos los sacrificios que han hecho por mí para que pueda cumplir esta meta. Gracias a todo el amor y apoyo que me brindaron durante mi carrera universitaria he podido salir adelante ante toda adversidad.

A mis hermanos Emily y Sebastian, los dos regalos más lindos que me pudieron dar mis padres. Gracias por su amor incondicional y la alegría que me han brindado durante esta etapa y toda mi vida.

Finalmente, agradezco a David, a quien tuve la suerte de conocer durante esta etapa universitaria y se convirtió en mi mejor amigo. Gracias por todo el amor y apoyo que me has dado a lo largo de este camino, pero sobre todo gracias por creer en mí.

Karolay De La Cruz Mosquera

2019

# ÍNDICE DE CONTENIDOS

<b>CERTIFICACIÓN</b> .....	i
<b>AUTORÍA DE RESPONSABILIDAD</b> .....	ii
<b>AUTORIZACIÓN</b> .....	iii
<b>DEDICATORIA</b> .....	iv
<b>AGRADECIMIENTO</b> .....	v
<b>ÍNDICE DE CONTENIDOS</b> .....	vi
<b>ÍNDICE DE TABLAS</b> .....	xiii
<b>ÍNDICE DE FIGURAS</b> .....	xv
<b>RESUMEN</b> .....	xxii
<b>ABSTRACT</b> .....	xxiii
<b>CAPÍTULO 1: Introducción</b> .....	1
1.1. Antecedentes .....	1
1.2. Justificación e Importancia .....	5
1.3. Objetivos .....	7
1.4. Descripción del Proyecto .....	8
<b>CAPÍTULO 2: Marco Conceptual</b> .....	11
2.1. Prototipo Nariz Electrónica.....	11

	vii
2.2. Machine Learning .....	14
2.2.1. Deep Learning .....	17
2.3. Proceso para la Elaboración de Modelos de Machine Learning y Deep Learning .....	18
2.3.1. Preprocesamiento.....	18
2.3.2. Ingeniería de Características.....	19
2.3.3. Aprendizaje.....	19
2.3.4. Evaluación .....	25
2.4. Desarrollo de Interfaz Gráfica de Usuario .....	25
<b>CAPÍTULO 3: Metodología Experimental .....</b>	<b>28</b>
3.1 Descripción de los Experimentos.....	28
3.2. Distribución de los Datos en Conjuntos de Entrenamiento y Prueba .....	31
3.3. Preprocesamiento de las Señales de Entrenamiento .....	34
3.3.1. Implementación de Filtro.....	34
3.3.2. Corrección de Línea Base.....	36
3.3.3. Alineamiento de Picos .....	39
3.3.4. Concatenación de Sensores .....	40
3.3.5. Detección y Eliminación de Outliers.....	41
<b>CAPÍTULO 4: Análisis Discriminante y Regresión de Mínimos Cuadrados Parciales</b> <b>(PLS-DA y PLS-R) .....</b>	<b>49</b>
4.1. Conceptos Básicos .....	49

	viii
4.1.1. Análisis Discriminante Lineal (LDA) .....	49
4.1.2. Regresión Lineal.....	50
4.1.3. Método de Mínimos Cuadrados .....	52
4.1.4. Regresión de Mínimos Cuadrados Parciales (PLS-R).....	53
4.1.5. Análisis Discriminante de Mínimos Cuadrados Parciales PLS-DA.....	54
4.1.6. Construcción de Componentes PLS .....	55
4.1.7. Construcción del Modelo de Predicción.....	56
4.2. Generación de Modelos PLS-DA y PLS-R .....	57
4.2.1. Escalamiento.....	57
4.2.2. Balanceo de Clases .....	60
4.2.3. Validación Cruzada .....	60
4.3. Resultados del Desempeño de los Modelos PLS-DA y PLS-R.....	61
4.3.1. Modelo PLS-DA.....	61
4.3.2. Modelo PLS-R.....	65
<b>CAPÍTULO 5: Regresión Logística .....</b>	<b>72</b>
5.1. Conceptos Básicos .....	72
5.1.1. Función Logística Sigmoide.....	72
5.1.2. Regularización L2.....	73
5.1.3. Regresión Logística .....	74
5.1.4. Construcción de Componentes PLS .....	74

	ix
5.1.5. Construcción de Modelo de Regresión Logística.....	76
5.2. Generación de Modelos de Regresión Logística .....	77
5.2.1. Escalamiento.....	78
5.2.2. Balanceo de Clases .....	78
5.2.3. Validación Cruzada .....	78
5.3. Resultados del Desempeño de los Modelos de Regresión Logística.....	79
<b>CAPÍTULO 6: Redes Neuronales Artificiales.....</b>	<b>85</b>
6.1. Conceptos Básicos .....	85
6.1.2. Construcción de Componentes PLS .....	87
6.1.3. Construcción del Modelo MLP .....	88
6.2. Generación de Modelos de Redes Neuronales MLP .....	95
6.2.1. Escalamiento.....	96
6.2.2. Balanceo de Clases .....	96
6.2.3. Validación Cruzada .....	96
6.3. Resultados del Desempeño de los Modelos de Redes Neuronales MLP.....	97
6.3.1. Modelos de Clasificación .....	97
6.3.2. Modelos de Regresión .....	102
<b>CAPÍTULO 7: Deep Learning .....</b>	<b>109</b>
7.1. Conceptos Básicos .....	109
7.1.1. Red Neuronal Recurrente (RNN) .....	109

	x
7.1.2. Red Neuronal de Memoria Larga a Corto Plazo (LSTM) .....	111
7.1.3. Construcción del Modelo de Deep Learning LSTM .....	113
7.2. Generación de Modelos de Deep Learning LSTM .....	115
7.2.1. Escalamiento .....	116
7.2.2. Balanceo de Clases .....	116
7.2.3. Validación Cruzada .....	116
7.3. Resultados del Desempeño de los Modelos de Deep Learning LSTM .....	117
7.3.1. Modelos de Clasificación .....	117
7.3.2. Modelos de Regresión .....	121
<b>CAPÍTULO 8: Análisis de Resultados y Pruebas</b> .....	<b>128</b>
8.1. Análisis Comparativo de los Modelos de Clasificación .....	128
8.1.1. Modelos de Tarea 1: Clasificación de Alcohol/Sustancias Explosivas en Estado Puro y Mezclas .....	129
8.1.2. Modelos de Tarea 2: Clasificación de Alcohol/Sustancias Explosivas en Estado Puro .....	131
8.1.3. Modelos de Tarea 3: Clasificación de Alcohol/Pólvora en Estado Puro .....	133
8.1.4. Modelos de Tarea 4: Clasificación de Alcohol/TNT en Estado Puro .....	135
8.1.5. Modelos de Tarea 5: Clasificación de Pólvora en Estado Puro/TNT en Estado Puro .....	137

8.1.6.	Modelos de Tarea 6: Clasificación de Alcohol/Pólvora en Estado Puro y Mezcla de Pólvora .....	139
8.1.7.	Modelos de Tarea 7: Clasificación de Alcohol/TNT en Estado Puro y Mezcla de TNT.....	141
8.1.8.	Modelos de Tarea 8: Clasificación de Alcohol/Pólvora en Estado Puro/TNT en Estado Puro .....	143
8.1.9.	Tiempo de Entrenamiento de Modelos de Clasificación.....	146
8.2.	Análisis Comparativo de los Modelos de Cuantificación.....	146
8.2.1.	Modelos de Tarea 1: Regresión de Alcohol/Sustancias Explosivas en Estado Puro y Mezclas .....	147
8.2.2.	Modelos de Tarea 2: Regresión de Alcohol/Sustancias Explosivas en Estado Puro.....	148
8.2.3.	Modelos de Tarea 3: Regresión de Alcohol/Pólvora en Estado Puro.....	149
8.2.4.	Modelos de Tarea 4: Regresión de Alcohol/TNT en Estado Puro .....	151
8.2.5.	Modelos de Tarea 5: Regresión de Pólvora en Estado Puro/TNT en Estado Puro.....	152
8.2.6.	Modelos de Tarea 6: Regresión de Alcohol/Pólvora en Estado Puro y Mezcla de Pólvora .....	154
8.2.7.	Modelos de Tarea 7: Regresión de Alcohol/TNT en Estado Puro y Mezcla de TNT .....	155

8.2.8. Modelos de Tarea 8: Regresión de Alcohol/Pólvora en Estado Puro/TNT en Estado Puro .....	156
8.2.9. Tiempo de Entrenamiento de los Modelos de Regresión .....	157
8.3. Exploración del Desempeño de la Interfaz Gráfica de Usuario.....	159
8.4. Discusión de Resultados .....	161
<b>CAPÍTULO 9 Conclusiones y Recomendaciones.....</b>	<b>163</b>
9.1. Conclusiones.....	163
9.2. Recomendaciones .....	168
<b>REFERENCIAS .....</b>	<b>169</b>

# ÍNDICE DE TABLAS

<b>Tabla 1</b> <i>Sensores empleados en el prototipo e-nose</i> .....	12
<b>Tabla 2</b> <i>Clasificación de algoritmos de machine learning empleados</i> .....	16
<b>Tabla 3</b> <i>Métricas de desempeño para tareas de clasificación</i> .....	22
<b>Tabla 4</b> <i>Métricas de desempeño del proyecto</i> .....	25
<b>Tabla 5</b> <i>Base de datos 1- Experimentos clasificados por concentración</i> .....	29
<b>Tabla 6</b> <i>Base de datos 1- Experimentos clasificados por sustancia</i> .....	29
<b>Tabla 7</b> <i>Base de datos 2 - Experimentos clasificados por sustancia</i> .....	30
<b>Tabla 8</b> <i>Base de datos 2 - Experimentos clasificados por concentración</i> .....	30
<b>Tabla 9</b> <i>Selección de las tareas de clasificación y regresión de los modelos multivariantes</i> .....	32
<b>Tabla 10</b> <i>Base de datos 1- Experimentos clasificados por concentración de sustancia</i> .....	33
<b>Tabla 11</b> <i>Numero de experimentos previo y posterior a la eliminación de outliers</i> .....	48
<b>Tabla 12</b> <i>Desempeño del modelo de clasificación de alcohol, pólvora y TNT</i> .....	58
<b>Tabla 13</b> <i>Desempeño del modelo de cuantificación de alcohol, pólvora y TNT</i> .....	60
<b>Tabla 14</b> <i>Base de datos 1- Resultados de los modelos PLS-DA</i> .....	61
<b>Tabla 15</b> <i>Base de datos 2- Resultados de los modelos PLS-DA</i> .....	61
<b>Tabla 16</b> <i>Base de datos 1- Resultados de modelos PLS-R</i> .....	66
<b>Tabla 17</b> <i>Base de datos 2- Resultados de modelos PLS-R</i> .....	66
<b>Tabla 18</b> <i>Base de datos 1- Resultados de modelos de regresión logística</i> .....	79
<b>Tabla 19</b> <i>Base de datos 2- Resultados de modelos de regresión logística</i> .....	79
<b>Tabla 20</b> <i>Base de datos 1- Resultados de modelos de clasificación MLP</i> .....	97
<b>Tabla 21</b> <i>Base de datos 2- Resultados de modelos de clasificación MLP</i> .....	98
<b>Tabla 22</b> <i>Base de datos 1- Resultados de modelos de regresión MLP</i> .....	103

<b>Tabla 23</b> <i>Base de datos 2- Resultados de modelos de regresión MLP.....</i>	103
<b>Tabla 24</b> <i>Base de datos 1- Resultados de los modelos de clasificación de deep learning .....</i>	117
<b>Tabla 25</b> <i>Base de datos 2- Resultados de los modelos de clasificación de deep learning .....</i>	118
<b>Tabla 26</b> <i>Base de datos 1- Resultados de los modelos de regresión de deep learning .....</i>	122
<b>Tabla 27</b> <i>Base de datos 2- Resultados de los modelos de regresión de deep learning .....</i>	122
<b>Tabla 28</b> <i>MSE y R2 de modelos de regresión de alcohol/sustancias explosivas en estado puro y mezclas.....</i>	147
<b>Tabla 29</b> <i>MSE y R2 de modelos de regresión de alcohol/sustancias explosivas en estado puro.....</i>	148
<b>Tabla 30</b> <i>MSE y R2 de modelos de regresión de alcohol/pólvora en estado puro.....</i>	149
<b>Tabla 31</b> <i>MSE y R2 de modelos de regresión de alcohol/pólvora en estado puro.....</i>	151
<b>Tabla 32</b> <i>MSE y R2 de modelos de pólvora en estado puro/TNT en estado puro.....</i>	153
<b>Tabla 33</b> <i>MSE y R2 de modelos de regresión de alcohol/pólvora en estado puro y mezcla de pólvora.....</i>	154
<b>Tabla 34</b> <i>MSE y R2 de modelos de regresión de alcohol/TNT en estado puro y mezcla de TNT.....</i>	155
<b>Tabla 35</b> <i>MSE y R2 de modelos de regresión de alcohol/pólvora en estado puro/TNT en estado puro.....</i>	156

# ÍNDICE DE FIGURAS

<i>Figura 1.</i> Prototipo nariz electrónica parte del proyecto de investigación 2016-PIC-009 .....	12
<i>Figura 2.</i> Diferencia de machine learning y programación clásica .....	15
<i>Figura 3.</i> Flujo de trabajo de un modelo .....	18
<i>Figura 4.</i> Validación cruzada de diez iteraciones k=10 fold .....	20
<i>Figura 5.</i> Matriz de confusión para dos clases .....	21
<i>Figura 6.</i> Curvas ROC para diferentes clasificadores .....	23
<i>Figura 7.</i> Estructura de la interfaz gráfica de usuario .....	26
<i>Figura 8.</i> Representación de los datos recolectados con el prototipo e-nose .....	31
<i>Figura 9.</i> Asignación de etiquetas categóricas y continuas a los experimentos .....	32
<i>Figura 10.</i> Espectro de potencia de los datos de entrenamiento-Base de datos 1 .....	35
<i>Figura 11.</i> Experimento de 2ml de alcohol filtrado - Base de datos 1 .....	36
<i>Figura 12.</i> Identificación de puntos máximos y mínimos de las señales de un experimento .....	38
<i>Figura 13.</i> Señales de un experimento con líneas base corregidas - Base de datos 1 .....	38
<i>Figura 14.</i> Señales correspondientes al sensor 4 alineadas- Base de datos 1 .....	40
<i>Figura 15.</i> Representación del conjunto de datos de entrenamiento .....	40
<i>Figura 16.</i> Señales de los seis sensores de un experimento concatenadas - Base de datos 1 .....	41
<i>Figura 17.</i> Representación geométrica de modelo PCA .....	42
<i>Figura 18.</i> Graficas de sedimentación de los componentes de los modelos PCA - Base de datos 1 .....	45
<i>Figura 19.</i> Detección de outliers en cada uno de los modelos PCA para las observaciones de la base de datos 1 .....	46

<b>Figura 20.</b> Observaciones previo a la eliminación de outliers y después de esta - Base de datos 1.....	47
<b>Figura 21.</b> Ejemplo de análisis discriminante lineal .....	50
<b>Figura 22.</b> Ejemplo de regresión lineal univariante .....	51
<b>Figura 23.</b> Proceso para la construcción de un modelo PLS-R.....	54
<b>Figura 24.</b> Proceso para la construcción de un modelo PLS-DA.....	55
<b>Figura 25.</b> Curvas ROC para la descripción del desempeño de los modelos en datos de entrenamiento y prueba.....	62
<b>Figura 26.</b> Curvas ROC para descripción del desempeño de los modelos para clasificación de sustancias explosivas puras y mezclas.....	63
<b>Figura 27.</b> Curvas ROC para descripción del desempeño de los modelos de clasificación entre dos y tres clases de sustancias.....	64
<b>Figura 28.</b> Curvas ROC para descripción del desempeño de los modelos de clasificación entre pólvora y TNT.....	65
<b>Figura 29.</b> Valores predichos por los modelos de regresión .....	67
<b>Figura 30.</b> Valores predichos por los modelos de regresión .....	68
<b>Figura 31.</b> Valores predichos por los modelos de regresión .....	69
<b>Figura 32.</b> Valores predichos por los modelos de regresión .....	70
<b>Figura 33.</b> Valores predichos por el modelo de regresión de sustancias explosivas puras para la base de datos 2.....	70
<b>Figura 34.</b> Función logística sigmoide .....	73
<b>Figura 35.</b> Proceso para la construcción del modelo de regresión logística .....	74
<b>Figura 36.</b> Grafica de sedimentación para selección de variables latentes .....	75

<b>Figura 37.</b> Pasos para construcción de modelo de Regresión Logística .....	76
<b>Figura 38.</b> Curvas ROC para la descripción del desempeño de los modelos en datos de entrenamiento y prueba.....	80
<b>Figura 39.</b> Curvas ROC para descripción del desempeño de los modelos para clasificación de sustancias explosivas puras y mezclas.....	81
<b>Figura 40.</b> Curvas ROC para descripción del desempeño de los modelos de clasificación entre dos y tres clases de sustancias.....	82
<b>Figura 41.</b> Curvas ROC para descripción del desempeño de los modelos de clasificación entre pólvora y TNT.....	83
<b>Figura 42.</b> Red Perceptrón Multicapa con una capa oculta.....	86
<b>Figura 43.</b> Proceso para la construcción de modelo de red neuronal MLP.....	87
<b>Figura 44.</b> Grafica de sedimentación para selección de variables latentes .....	88
<b>Figura 45.</b> Número de capas ocultas seleccionadas para el modelo MLP .....	89
<b>Figura 46.</b> Funciones de activación para las capas ocultas de una red neuronal MLP .....	91
<b>Figura 47.</b> Función de optimización: Descenso de gradiente .....	94
<b>Figura 48.</b> Función de optimización: Descenso de gradiente estocástico.....	95
<b>Figura 49.</b> Curvas ROC para la descripción del desempeño de los modelos en datos de entrenamiento y prueba.....	99
<b>Figura 50.</b> Curvas ROC para descripción del desempeño de los modelos para clasificación de sustancias explosivas puras y mezclas.....	100
<b>Figura 51.</b> Curvas ROC para descripción del desempeño de los modelos de clasificación entre dos y tres clases de sustancias.....	101

<b>Figura 52.</b> Curvas ROC para descripción del desempeño de los modelos de clasificación entre pólvora y TNT.....	102
<b>Figura 53.</b> Valores predichos por los modelos de regresión .....	104
<b>Figura 54.</b> Valores predichos por los modelos de regresión .....	105
<b>Figura 55.</b> Valores predichos por los modelos de regresión .....	106
<b>Figura 56.</b> Valores predichos por los modelos de regresión .....	107
<b>Figura 57.</b> Valores predichos por el modelo de regresión de sustancias explosivas puras para la base de datos 2.....	108
<b>Figura 58.</b> Arquitectura de una red neuronal clásica y de una red neuronal recurrente.....	110
<b>Figura 59.</b> Estructura de una neurona LSTM.....	112
<b>Figura 60.</b> Proceso para la construcción de modelo de deep learning LSTM.....	112
<b>Figura 61.</b> Arquitectura de la red neuronal de deep learning LSTM .....	113
<b>Figura 62.</b> Curvas ROC para la descripción del desempeño de los modelos en datos de entrenamiento y prueba.....	118
<b>Figura 63.</b> Curvas ROC para descripción del desempeño de los modelos para clasificación de sustancias explosivas puras y mezclas.....	119
<b>Figura 64.</b> Curvas ROC para descripción del desempeño de los modelos de clasificación entre dos y tres clases de sustancias.....	120
<b>Figura 65.</b> Curvas ROC para descripción del desempeño de los modelos de clasificación entre pólvora y TNT.....	121
<b>Figura 66.</b> Valores predichos por los modelos de regresión .....	123
<b>Figura 67.</b> Valores predichos por los modelos de regresión .....	124
<b>Figura 68.</b> Valores predichos por los modelos de regresión .....	125

<b>Figura 69.</b> Valores predichos por los modelos de regresión .....	126
<b>Figura 70.</b> Valores predichos por el modelo de regresión de sustancias explosivas puras para la base de datos 2.....	126
<b>Figura 71.</b> Curva ROC de modelos de clasificación de alcohol/sustancias explosivas en estado puro y mezclas.....	129
<b>Figura 72.</b> Umbral de modelo de clasificación de alcohol/sustancias explosivas en estado puro y mezclas LSTM.....	130
<b>Figura 73.</b> Matriz de confusión del modelo de clasificación de alcohol/sustancias explosivas en estado puro y mezclas.....	130
<b>Figura 74.</b> Curva ROC de modelos de clasificación de alcohol/sustancias explosivas en estado puro.....	131
<b>Figura 75.</b> Matriz de confusión del modelo de clasificación de alcohol/sustancias explosivas en estado puro.....	132
<b>Figura 76.</b> Curva ROC de modelos de clasificación de alcohol/pólvora en estado puro.....	133
<b>Figura 77.</b> Matriz de confusión del modelo de clasificación de alcohol/pólvora en estado puro.....	134
<b>Figura 78.</b> Curva ROC de modelos de clasificación de alcohol/TNT en estado puro .....	135
<b>Figura 79.</b> Matriz de confusión del modelo de clasificación de alcohol/TNT en estado puro ...	136
<b>Figura 80.</b> Curva ROC de modelos de clasificación de alcohol/TNT en estado puro .....	137
<b>Figura 81.</b> Matriz de confusión del modelo de clasificación de pólvora en estado puro /TNT en estado puro.....	139
<b>Figura 82.</b> Curva ROC de modelos de clasificación de Alcohol/Pólvora en estado puro y mezcla de pólvora.....	140

<b>Figura 83.</b> Matriz de confusión del modelo de clasificación de Alcohol/Pólvora en estado puro y mezcla de pólvora.....	141
<b>Figura 84.</b> Curva ROC de modelos de clasificación de Alcohol/TNT en estado puro y mezcla de TNT.....	142
<b>Figura 85.</b> Matriz de confusión del modelo de clasificación de Alcohol/TNT en estado puro y mezcla de TNT.....	143
<b>Figura 86.</b> Curva ROC de modelos de clasificación de Alcohol/Pólvora en estado puro/TNT en estado puro.....	144
<b>Figura 87.</b> Matriz de confusión del modelo de clasificación de Alcohol/Pólvora en estado puro/TNT en estado puro.....	145
<b>Figura 88.</b> Tiempo de entrenamiento de modelos de clasificación.....	146
<b>Figura 89.</b> Valores predichos por el modelo de regresión de alcohol/sustancias explosivas en estado puro y mezclas.....	147
<b>Figura 90.</b> Valores predichos por el modelo de regresión de alcohol/sustancias explosivas en estado puro.....	149
<b>Figura 91.</b> Valores predichos por el modelo de regresión de alcohol/pólvora en estado puro ...	150
<b>Figura 92.</b> Valores predichos por el modelo de regresión de alcohol/TNT en estado puro.....	152
<b>Figura 93.</b> Valores predichos por el modelo de regresión de alcohol/TNT en estado puro.....	153
<b>Figura 94.</b> Valores predichos por el modelo de regresión de alcohol/pólvora en estado puro y mezcla de pólvora.....	154
<b>Figura 95.</b> Valores predichos por el modelo de regresión de alcohol/TNT en estado puro y mezcla de TNT.....	155
<b>Figura 96.</b> Valores predichos por el modelo de regresión de Alcohol/Pólvora /TNT .....	157

	xxi
<b>Figura 97.</b> Tiempo de entrenamiento de modelos de regresión .....	158
<b>Figura 98.</b> Interfaz gráfica de usuario del proyecto de investigación .....	159
<b>Figura 99.</b> Preprocesamiento en la interfaz gráfica de usuario del proyecto de investigación ...	160
<b>Figura 100.</b> Ventana emergente de resultado de entrenamiento de modelo PLS-DA .....	161
<b>Figura 101.</b> Matriz de confusión resultante del entrenamiento de modelo PLS-DA .....	161

## RESUMEN

La detección de olores mediante sistemas de olfato artificial, denominados narices electrónicas, es un tema de investigación actual con aplicaciones a nivel militar como la detección de sustancias explosivas. La necesidad de evitar el tráfico ilegal de este tipo de sustancias debido a temas de seguridad nacional ha motivado la implementación y optimización de un prototipo e-nose para la detección de sustancias explosivas, parte del proyecto de investigación 2016-pic-009. Sin embargo, se ha dado un mayor enfoque a la optimización del hardware y no a los modelos con los cuales el prototipo será capaz de clasificar y cuantificar sustancias explosivas como TNT y pólvora en base doble. Por lo cual, el propósito de este trabajo de investigación es generar y analizar modelos de machine learning mediante las técnicas lineales: mínimos cuadrados parciales y regresión logística, y técnicas no lineales: red neuronal perceptrón multicapa y red neuronal profunda LSTM, integrados en una interfaz gráfica de usuario para el reentrenamiento o prueba de los modelos. Los resultados del proyecto muestran un mejor desempeño en la clasificación de sustancias explosivas con concentraciones entre 3 y 5gr con 1mL de sustancia dopante que con concentraciones entre 0.1 y 3gr con 2mL de sustancia dopante. Además, en cuanto a la cuantificación el  $R^2$  no supero el 0.57 para las condiciones iniciales del prototipo y el 0.22 para las actuales.

### **PALABRAS CLAVE:**

- **MACHINE LEARNING**
- **MÍNIMOS CUADRADOS PARCIALES**
- **REGRESIÓN LOGÍSTICA**
- **RED NEURONAL ARTIFICIAL PERCEPTRÓN MULTICAPA**
- **RED NEURONAL ARTIFICIAL PROFUNDA LSTM**

## ABSTRACT

The detection of odors by artificial smell systems, called electronic noses, is a subject of current research with military-level applications such as the detection of explosive substances. The need to prevent illegal trafficking of these substances due to national security issues has motivated the implementation and optimization of an e-nose prototype for the detection of explosive substances, part of the 2016-pic-009 research project. However, greater focus has been given to hardware optimization and not to the models with which the prototype will be able to classify and quantify explosive substances such as TNT and double-base gunpowder. Therefore, the purpose of this research work is to generate and analyze machine learning models with linear techniques: partial least squares and logistic regression, and nonlinear techniques: multilayer perceptron neural network and LSTM deep neural network, integrated into a graphical user interface for retraining or testing models. The results of the project show a better performance in the classification of explosive substances with concentrations between 3 and 5gr with 1mL of doping substance than with concentrations between 0.1 and 3gr with 2mL of doping substance. In addition, in terms of quantification,  $R^2$  does not exceed 0.57 for the initial conditions of the prototype and 0.22 for the current conditions.

### KEY WORDS:

- **MACHINE LEARNING**
- **PARTIAL LEAST SQUARE**
- **LOGISTIC REGRESSION**
- **MULTILAYER PERCEPTRON ARTIFICIAL NEURONAL NETWORK**
- **LSTM ARTIFICIAL DEEP NEURONAL NETWORK**

# CAPÍTULO 1

## Introducción

### 1.1. Antecedentes

Actualmente, existe un gran interés en la detección de olores mediante instrumentación electrónica. Los sistemas de olfato artificial, denominados narices electrónicas (e-nose), son instrumentos que, mediante una matriz de sensores químicos electrónicos, un circuito de muestreo y un algoritmo clasificador de patrones, son capaces de reconocer olores simples o complejos (J. W. Gardner, 2013; Walt & Sternfeld, 2013). Sin embargo, la lista de técnicas de preprocesamiento y reconocimiento de patrones aplicadas a tecnología e-nose es extensa, por lo cual deben ser seleccionadas en función del tipo de sensores utilizados y la naturaleza del problema (J. W. Gardner, 2013).

Las narices electrónicas son usadas en una amplia variedad de aplicaciones, entre las cuales se encuentran: detección de materiales explosivos, perfumes/aceites esenciales, olores simples como etanol, comida, solventes en polímeros, pinturas, plásticos, entre otros (Gardner, 2013). En el ámbito de la detección de explosivos, han demostrado ser un campo de investigación importante debido a temas de seguridad nacional y aplicaciones militares. Además, es una alternativa al uso de canes entrenados, los cuales son difíciles de entrenar y cuyo cuidado es costoso (Trogler, 2013).

La aplicación de tecnología e-nose para detección de sustancias explosivas es un tema de investigación de la Universidad de las Fuerzas Armadas-ESPE, mediante el proyecto 2016-PIC-009

“Localización de TNT y pólvora de base doble a través de sensado químico en un entorno controlado mediante robótica cooperativa”, estudiantes y docentes de carreras afines a las líneas de investigación han podido realizar contribuciones al proyecto, algunas de las cuales se discuten a continuación. En el proyecto de titulación “Desarrollo de un prototipo electrónico de sensado químico, para la detección de trinitrotolueno (TNT) y pólvora base doble en un ambiente controlado” (López Hernández, 2016), los autores desarrollan un prototipo electrónico para la detección de trinitrotolueno (TNT) y pólvora de base doble (mezcla de nitroglicerina y nitrocelulosa), cuyos resultados presentan un 70% de eficiencia en la detección de las sustancias y un 86,67% de discriminación entre sustancias explosivas y no explosivas, todo esto mediante la implementación del modelo de análisis discriminante de componentes principales (PCDA) con un clasificador de los K vecinos más cercanos (kNN k nearest neighbor). En el proyecto “Optimización e integración de una nariz electrónica autónoma embebida en un sistema robótico para la identificación de sustancias explosivas como TNT y pólvora base doble en ambientes controlados” (Espinosa & Venegas, 2017), se desarrolla un prototipo electrónico basado en sensores químicos de óxido metálico para el reconocimiento y clasificación de sustancias como TNT y pólvora de base doble, para esto se genera un modelo multivariante mediante la combinación del análisis de componentes principales (PCA) y el análisis discriminante lineal (LDA) con un clasificador de los K vecinos más cercanos. Como resultado se obtuvo una tasa de clasificación de 58.88% al discriminar alcohol, TNT y pólvora, y un 80% entre sustancias explosivas y no explosivas. En “Evaluación del desempeño de un UAV, equipado con dispositivos de sensado para la identificación de componentes químicos volátiles empleados en el procesamiento ilícito de sustancias” (Vallejos Brito, 2019), el autor propone el desarrollo de un prototipo de sistema de sensado de sustancias ilícitas cuyo algoritmo de clasificación basado en redes neuronales demostró no ser efectivo en la

discriminación de las sustancias. En el artículo titulado “Electronic nose prototype detection” (López, Triviño, Calderón, Arcentales, & Guamán, 2017), se presenta el proceso de diseño de un prototipo de nariz electrónica para la detección de sustancias explosivas además de la implementación de un modelo PCA para la identificación de sustancias y una configuración experimental para pruebas, el cual permite detectar y diferenciar distintas sustancias de diversas fuentes, una de ellas de la familia de explosivos, con una sensibilidad del 80% y una especificidad del 90% utilizando la validación cruzada Leave-one-out. En “Multivariate Discrimination Model for TNT and Gunpowder Using an Electronic Nose Prototype: A Proof of Concept” (Guaman, Lopez, & Torres-Tello, 2019), los autores presentan una prueba de concepto para discriminar sustancias explosivas, por medio del desarrollo de un modelo de discriminación para la clasificación de TNT y pólvora mediante la combinación de PCA y el análisis discriminante de Fisher (FDA) donde se obtuvo como resultado una precisión del 67% cuando tres sustancias diferentes (dos explosivos y uno no explosivo) fueron discriminadas con un valor de  $p < 0.01$ , y una precisión del 86.6% cuando fueron discriminadas sustancias explosivas y no explosivas. Otros resultados se presentan en “Detection of explosives with laser-induced breakdown spectroscopy” (Wang, Liu, Zhao, Ge, & Huang, 2012), donde uno de los enfoques del trabajo es la creación de un modelo de análisis discriminante de mínimos cuadrados parciales (PLS-DA) para distinguir explosivos de materiales orgánicos como el plástico. El modelo fue probado con 100 muestras de espectros de TNT y de siete tipos de plásticos (700 muestras), cuyo resultado demostró que la discriminación entre TNT y plásticos es posible utilizando PCA y PLS-DA, a pesar de tener una composición atómica similar. En “Simulation of Artificial Noses for the Automated Detection and Classification of Organic Compounds” (B. Lakshmi, K.Parish Venkata Kumar, Dr. K.Nageswara Rao, 2013), se presenta un prototipo de nariz electrónica compuesto por nueve sensores de vapor de peróxido, un

sensor de humedad y uno de temperatura, para la detección de químicos comunes del hogar (isopropanol, acetona, amoníaco, líquido encendedor y vinagre) mediante la construcción de dos redes neuronales artificiales: una red feedforward multicapa estándar entrenada con el algoritmo backpropagation y otra entrenada con el algoritmo difuso ARTmap. Los resultados de este trabajo demuestran que al utilizar la red neuronal entrenada con el algoritmo difuso ARTmap se obtiene un mayor porcentaje de aciertos en la clasificación de químicos que con el algoritmo de Backpropagation.

En los trabajos mencionados se observa que, se han implementado diversos tipos de técnicas para la detección de sustancias explosivas. Sin embargo, en el proyecto del que forma parte la investigación, hasta el momento, no se han desarrollado modelos que permitan realizar tanto la discriminación como la cuantificación de las sustancias, además, es importante señalar que no se ha dado un enfoque al uso de Redes Neuronales Artificiales (ANN) a pesar de que ANN son consideradas uno de los métodos más prometedores para este tipo de aplicación, ya que pueden hacer frente a problemas no lineales y manejar el ruido o la deriva mejor que los enfoques estadísticos convencionales (Fu, Li, Qin, & Freeman, 2007) . Por lo tanto, el presente trabajo de investigación pretende realizar un preprocesamiento de las señales obtenidas con el prototipo e-nose elaborado en (Vallejo & Zurita, 2017) y optimizado en (Jacome, 2019), así como utilizarlas para generar modelos lineales y no lineales que permitirán discriminar y cuantificar sustancias explosivas, además de desarrollar una interfaz que permita integrar los modelos para su utilización con nuevos datos de prueba o reentrenarlos en caso de ser necesario. Todo ello, con el fin de mejorar el procesamiento de las señales provenientes de la nariz electrónica, dado que constituye su etapa final y es aquella que mediante su etapa de preprocesamiento podrá reducir matemáticamente los

efectos no deseados en la respuesta de los sensores, detectar posibles sensores defectuosos, y mediante su etapa de reconocimiento de patrones cuantificar y discriminar las sustancias analizadas.

## **1.2. Justificación e Importancia**

La presente propuesta de proyecto de investigación nace de la necesidad de desarrollar modelos que permitan la clasificación y cuantificación de sustancias explosivas como TNT y pólvora de base doble mediante técnicas lineales y no lineales de aprendizaje automático. Además, es un tema que busca aportar al proyecto de investigación 2016-pic-009 “Localización de TNT y pólvora de base doble a través de sensado químico en un entorno controlado mediante robótica cooperativa”.

En los últimos años en Ecuador la policía ha incorporado a sus operativos de seguridad perros adiestrados para la detección de materiales explosivos. El proceso de adiestramiento dura aproximadamente doce meses con jornadas diarias extensas, tanto para los canes como para sus guías, además, el cuidado de la dieta y salud de los canes implica una alimentación balanceada y visitas regulares al veterinario (Ministerio del Interior, 2019). Todo ello significa costos tanto en tiempo como en dinero que podrían evitarse al utilizar tecnología e-nose.

La detección de explosivos mediante tecnología e-nose implica la aplicación de diferentes técnicas de discriminación y cuantificación, varias de las cuales se discuten a continuación. En el caso del trabajo de investigación titulado “Creación de un modelo de calibración multivariante para determinar el límite de detección de un prototipo nariz electrónica para medición de sustancias explosivas” (Salazar, 2018), cuya base de datos es usada en el presente trabajo de investigación, se realizaron 149 experimentos usando TNT y pólvora base doble en su estado puro y mezclados con pasta dental y jabón, para la generación de modelos de calibración multivariantes usando Regresión

de Componentes Principales (PCR) y Regresión de Mínimos Cuadrados Parciales (PLS), además, se realizó el análisis de componentes principales (PCA) para observar la separación de las sustancias y corregir errores que podrían afectar el desarrollo de los modelos antes mencionados. Los límites de detección del prototipo determinados al utilizar el modelo PLS fueron de: 1.42g para TNT, 1.02g para TNT mezclado, 0.99g para pólvora en estado puro y 0.86g para pólvora mezclada, por otra parte, los coeficientes de determinación del modelo PLS oscilaron en un rango de 0.39 a 0.73 y del modelo PCR en un rango de 0.10 a 0.31 (Salazar, 2018). Algunos de los problemas identificados en este proyecto de titulación son: en la etapa de preprocesamiento se elimina únicamente los valores en los cuales la nariz realiza el proceso de limpieza, sin realizar un análisis previo de que parte de las señales obtenidas en cada uno de los experimentos tiene información útil para la generación de los modelos, por otra parte, los datos predichos por el modelo PCR presentan un desempeño “malo” y “muy malo” y los predichos por el modelo PLS un desempeño “regular” y “bueno”, por ello, es necesario efectuar un análisis de las señales obtenidas en los experimentos para elaborar una etapa de preprocesamiento de las mismas, la exploración y generación de nuevos métodos que permitan obtener mejores resultados de los que se han obtenido hasta el momento en el proyecto del que forma parte la investigación y realizar tanto la discriminación como la cuantificación de sustancias explosivas. Además, dado que en las pruebas de certificación para canes entrenados se establece que debe realizarse con no menos de 5g de sustancia explosiva, es importante realizar nuevos experimentos de hasta esa cantidad de concentración para analizar el poder predictivo del prototipo e-nose.

De acuerdo con el estado del arte presentado en la sección anterior y el análisis previo de esta sección, diversos modelos han sido utilizados para la discriminación y cuantificación de sustancias explosivas, sin embargo, debido a la existencia de un gran número de técnicas tanto

cuantitativas como cualitativas, es necesario explorar otros métodos que permitan obtener mejores resultados, tomando en cuenta que al combinar modelos cuantitativos y cualitativos será posible decidir si una sustancia explosiva existe o no en una muestra y predecir el rango de concentración de la sustancia, los cuales serán usados como modelos de referencia para el proyecto de investigación 2016-pic-009.

Finalmente, una interfaz en la que se pueda ingresar nuevos datos y visualizar los resultados de la clasificación y cuantificación de sustancias explosivas, o a su vez volver a realizar el entrenamiento de los modelos, es algo que no se presenta en ninguno de los trabajos citados, por lo cual es un aporte significativo al proyecto de investigación en el que se enmarca este trabajo.

## **1.3. Objetivos**

### **1.3.1. GENERAL**

Generar modelos cualitativos y cuantitativos para discriminar y cuantificar sustancias explosivas con el prototipo nariz electrónica del proyecto de investigación 2016-pic-009.

### **1.3.2. ESPECÍFICOS**

- Analizar y evaluar la base de datos existente para realizar el preprocesamiento de información.
- Generar dos modelos lineales y dos no lineales para la discriminación y cuantificación sustancias explosivas.
- Realizar un análisis comparativo de los modelos generados, los cuales servirán como modelos de referencia en el proyecto.

- Desarrollar una interfaz gráfica de usuario que permita ingresar nuevos datos del prototipo nariz electrónica y obtener los resultados de clasificación y predicción de sustancias explosivas.
- Realizar nuevos experimentos con el prototipo nariz electrónica y analizar su poder predictivo.

## 1.4. Descripción del Proyecto

El presente proyecto de investigación busca desarrollar modelos multivariantes de machine learning y deep learning para la discriminación de sustancias explosivas como lo son la pólvora en base doble y el Trinitrotolueno (TNT) y sustancias no explosivas, además de la predicción de la concentración de las mismas. Para ello se seleccionaron dos técnicas lineales: mínimos cuadrados parciales (PLS) y regresión logística, y dos no lineales: red neuronal perceptrón multicapa y red neuronal profunda Long Short Term Memory (LSTM). Los conceptos empleados para el desarrollo del proyecto, el proceso para la creación de los modelos, los resultados obtenidos y las conclusiones y recomendaciones basadas en estos resultados son presentados en nueve capítulos, los cuales se describen a continuación.

### Capítulo 2: Marco Conceptual

En este capítulo se presentan los conceptos para el desarrollo del proyecto de investigación, incluyendo información sobre prototipo nariz electrónica del que cual se obtuvieron los datos para el desarrollo de los modelos multivariantes, sobre los diferentes tipos de algoritmos que se usarán para la discriminación y cuantificación de sustancias explosivas y las herramientas para el desarrollo de la interfaz gráfica de usuario que integre estos modelos.

### **Capítulo 3: Metodología Experimental**

El Capítulo 3 presenta una descripción de los datos que se usaran para el desarrollo y validación de los modelos, además de cada una de las etapas de preprocesamiento y sus respectivos resultados.

### **Capítulo 4: Análisis Discriminante y Regresión de Mínimos Cuadrados Parciales (PLS-DA y PLS-R)**

En el Capítulo 4 se describe el algoritmo de mínimos cuadrados parciales, el proceso para el desarrollo de los modelos de clasificación PLS-DA y de regresión PLS-R, y los resultados obtenidos.

### **Capítulo 5: Regresión Logística**

En el Capítulo 5 se presenta una descripción del algoritmo de regresión logística, el proceso para el desarrollo de los modelos de clasificación y los resultados obtenidos con estos.

### **Capítulo 6: Redes Neuronales Artificiales**

El Capítulo 6 describe una red neuronal perceptrón multicapa, el proceso para el desarrollo de los modelos de clasificación y regresión, y los resultados alcanzados.

### **Capítulo 7: Deep Learning**

En este capítulo se presenta una descripción de la red neuronal profunda LSTM, el proceso para el desarrollo de los modelos de clasificación y regresión, y los resultados obtenidos.

**Capítulo 8: Análisis de Resultados y Pruebas**

El Capítulo 8 muestra un análisis comparativo de los resultados obtenidos con las técnicas utilizadas para el desarrollo de los modelos de clasificación y regresión, y los resultados del desempeño de la interfaz gráfica de usuario.

**Capítulo 9: Conclusiones y Recomendaciones**

En el Capítulo 9 se presentan las conclusiones en base a los objetivos planteados y las recomendaciones para investigaciones futuras relacionadas al proyecto de investigación 2016-PIC-009.

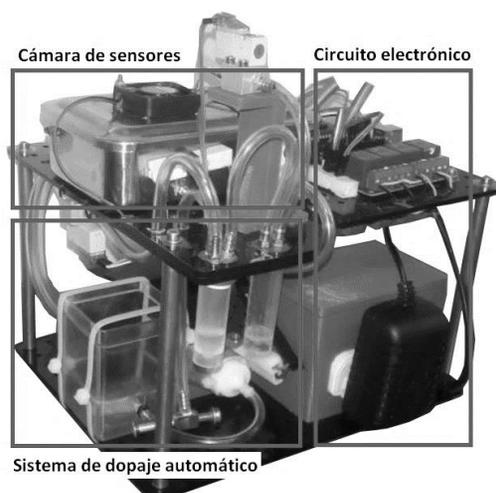
# CAPÍTULO 2

## Marco Conceptual

En este capítulo se presenta una descripción del prototipo del cual se obtuvieron los datos para la creación y evaluación de los modelos multivariantes, incluyendo las partes que lo conforman y el tipo de sensores utilizados. Además, se presentan los conceptos empleados para el desarrollo de los modelos, como lo son las etapas de preprocesamiento, entrenamiento y prueba. Finalmente, se describe la herramienta utilizada para la elaboración de la interfaz gráfica de usuario, la cual integra los modelos desarrollados para su posterior reentrenamiento o prueba.

### 2.1. Prototipo Nariz Electrónica

El prototipo e-nose del cual se obtuvieron los datos para el desarrollo del proyecto de investigación forma parte del proyecto 2016-PIC-009 “Localización de TNT y pólvora de base doble a través de sensado químico en un entorno controlado mediante robótica cooperativa”. Este prototipo elaborado por Vallejo y Zurita (2017) y optimizado por Jacome (2019), es un dispositivo portable para el reconocimiento y clasificación de sustancias como pólvora en base doble y TNT (Ver Figura 1), y está conformado por los siguientes elementos:



**Figura 1.** Prototipo nariz electrónica parte del proyecto de investigación 2016-PIC-009

Fuente: (Vallejo & Zurita, 2017)

- Cámara de sensores:** una cámara de acero inoxidable cuya función es evitar que cambios en el entorno interfieran en la respuesta de un conjunto de sensores encargados de convertir las señales capturadas por los mismos en señales electrónicas que puedan ser utilizadas para el reconocimiento de las características olfativas. Conformada por un ventilador para su limpieza, un sensor de temperatura DTH100 para que las señales sean adquiridas a una temperatura estable (a un valor de referencia de 29°C) y seis sensores químicos T-GS (dos sensores T-GS 822, dos sensores T-GS 826, un sensor T-GS 2610 y un sensor T-GS 825). Los sensores y gases a los que cada uno es sensible se presentan en la Tabla 1.

**Tabla 1**

*Sensores empleados en el prototipo e-nose*

<b>SENSOR</b>	<b>GASES A LOS QUE ES SENSIBLE</b>
TGS 822	Vapores solventes orgánicos y otros vapores volátiles
TGS 825	Sulfuro de hidrogeno
TGS 826	Gas amoniaco
TGS 2610	Gas licuado de petróleo

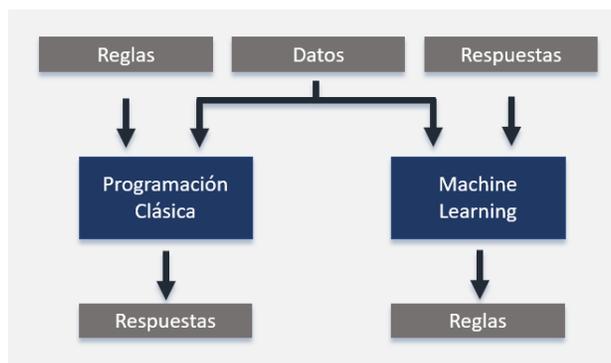
- **Circuito neumático:** cumple la función de respiración (transporte de la sustancia de interés a la cámara de sensores) y limpieza del sistema (activación del ventilador en la cámara de sensores). Este circuito está conformado por un compresor para la entrega de aire al sistema de respiración, mangueras por las cuales atraviesa el aire, acoples y tres electroválvulas: una de ellas para el paso de aire del compresor a los contenedores, otra para el paso de aire contaminado con las sustancias de los contenedores a la cámara de sensores y una última para el paso de aire limpio a la cámara de sensores (Salazar, 2018).
- **Circuitos electrónicos:** utilizados para controlar el funcionamiento de las electroválvulas y para el control de temperatura dentro de la cámara de sensores, se encuentran conformados por un circuito de potencia y uno de control (Vallejo & Zurita, 2017).
- **Sistema de dopaje automático:** sistema para el dopaje automático de los experimentos, conformado por un contenedor de alcohol, una bomba y un depósito de muestras con dos contenedores: el primero en el cual se ingresa de forma automática la sustancia dopante y el segundo para el ingreso manual de la sustancia a analizar.
- **Sistema de adquisición de datos:** sistema constituido por un Arduino UNO para la adquisición de las señales provenientes de los sensores y una interfaz para su almacenamiento y exportación a un archivo. .xlsx.

Su funcionamiento consiste en conectar el prototipo e-nose a la fuente de poder y la tarjeta de control a la computadora. Posterior a ello, se debe verificar que en el contenedor de alcohol exista suficiente contenido para el dopaje automático del contenedor de alcohol del depósito de muestras e ingresar en el contenedor de sustancia explosiva la sustancia de interés. A continuación, en la interfaz del prototipo se selecciona el botón de inicio para empezar el proceso de adquisición de datos. Inicialmente este proceso duraba 650s, dentro de los cuales los 300s iniciales correspondían

a la etapa de limpieza del dispositivo, además, se tomaba una muestra por segundo, es decir que el prototipo adquiriría los datos con una frecuencia de muestreo de 1Hz. Actualmente debido a modificaciones realizadas en (Jacome, 2019), el proceso de adquisición de datos dura 180s, dentro de los cuales los 15 primeros segundos la respuesta de los sensores se estabiliza antes de transportar los gases de los contenedores a la cámara de sensores, además, se toma una muestra cada 0.1s, es decir que el prototipo actualmente adquiere los datos con una frecuencia de muestreo de 10Hz. Finalmente, los datos obtenidos con cada experimento se almacenan en archivos .xlsx., los cuales son utilizados mediante técnicas de aprendizaje automático para la clasificación y cuantificación de sustancias explosivas, cuyos conceptos se presentan a continuación.

## **2.2. Machine Learning**

Machine Learning o aprendizaje de máquina, es un subcampo de la inteligencia artificial (IA) que involucra algoritmos de autoaprendizaje para la construcción de modelos estadísticos basados en datos representativos de los fenómenos analizados (Chollet, 2018). Para elaborar un modelo de machine learning se manejan dos conjuntos de datos: el conjunto de entrenamiento utilizado para descubrir de forma automática las representaciones más adecuadas de los datos de entrada con las que se pueda hacer predicciones sobre nuevos conjuntos de datos (Chollet, 2018), y el conjunto de prueba utilizado para evaluar el desempeño final del modelo en datos que no hayan sido usados en su entrenamiento. Por lo tanto, contrario a la programación clásica, un modelo de machine learning aprenderá de los datos conocidos, descubriendo patrones sin depender de una programación basada en reglas, como se muestra en la Figura 2. El aprendizaje de máquina se puede dividir en tres tipos: aprendizaje supervisado, no supervisado y por refuerzo, cuyas definiciones se detallan continuación.



**Figura 2.** Diferencia de machine learning y programación clásica

Fuente: Basada en (Chollet, 2018)

- **Aprendizaje supervisado:** El objetivo del aprendizaje supervisado es usar un conjunto de datos de entrada (ejemplos de entrenamiento) y sus respectivas salidas para producir un modelo que permita realizar predicciones sobre nuevos datos. Dependiendo del tipo de salida, los algoritmos de aprendizaje supervisado pueden ser de clasificación en el caso de que las salidas sean discretas o de regresión cuando los valores de salida son continuos (Raschka & Mirjalili, 2017).
- **Aprendizaje no supervisado:** El aprendizaje no supervisado permite explorar la estructura de los datos y extraer información significativa de estos sin la guía de una variable de salida (Raschka & Mirjalili, 2017). Dependiendo de la tarea que se desea realizar se divide en dos tipos: clustering y reducción de dimensionalidad. El clustering permite organizar datos en subgrupos sin tener un conocimiento previo del grupo al que pertenecen y la reducción de dimensionalidad permite comprimir los datos en un subespacio dimensional más pequeño mientras se retiene la información más importante (Raschka & Mirjalili, 2017).
- **Aprendizaje reforzado:** En el aprendizaje reforzado los modelos aprenden en base a la experiencia. A través de su interacción con el entorno, el modelo es premiado cuando realiza

bien su tarea o castigado cuando toma una mala decisión, con lo cual el modelo aprende por sí mismo cual es la mejor estrategia para obtener la mayor recompensa con el tiempo (Geron, 2017).

En el presente proyecto de investigación se emplean los algoritmos presentados en la **Tabla 2** para el desarrollo de los modelos multivariantes. El algoritmo de aprendizaje no supervisado denominado análisis de componentes principales (PCA) es utilizado para la detección de señales atípicas dentro de los datos de entrenamiento y prueba, esta técnica es descrita en el Capítulo 3. Además, para la clasificación y cuantificación de sustancias explosivas se emplean algoritmos de aprendizaje supervisado: el método de mínimos cuadrados parciales (PLS), regresión logística y red neuronal perceptrón multicapa (MLP). Los tres algoritmos son empleados para la predicción de salidas categóricas que en este caso son las clases de sustancias explosivas, y únicamente el método de mínimos cuadrados parciales (PLS) y una red neuronal perceptrón multicapa (MLP) para la predicción de salidas continuas que representan los gramos de concentración de las sustancias analizadas.

**Tabla 2**  
*Clasificación de algoritmos de machine learning empleados*

Algoritmos No Supervisados	Algoritmos Supervisados	
	Clasificación	Regresión
Análisis de componentes principales (PCA)	Análisis discriminante de mínimos cuadrados parciales (PLS-DA)	Regresión de mínimos cuadrados parciales (PLS-R)
	Red neuronal perceptrón multicapa (MLP)	Red neuronal perceptrón multicapa (MLP)
	Regresión Logística	

### 2.2.1. Deep Learning

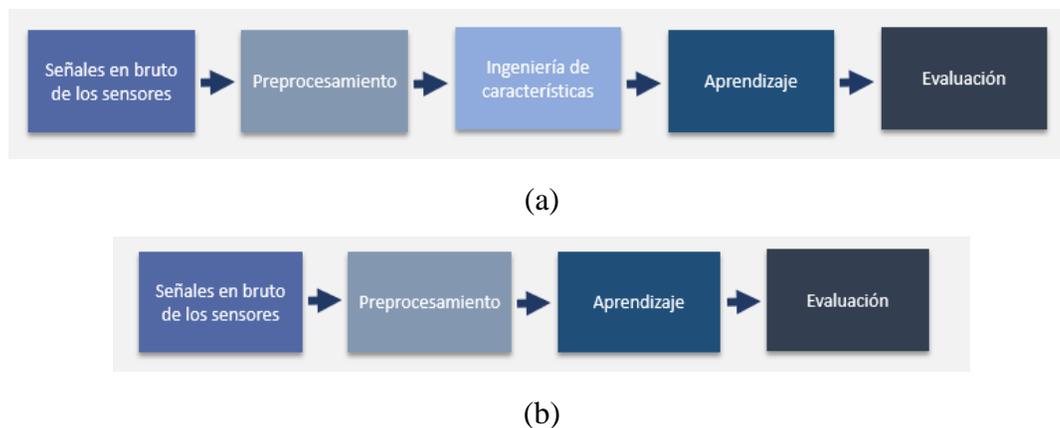
Deep Learning, es un subcampo del Machine Learning, que pone énfasis en el aprendizaje de capas sucesivas cada vez más significativas de representaciones de datos. Implica capas de representaciones sucesivas que son aprendidas de la exposición a los datos de entrenamiento, casi siempre a través de modelos de redes neuronales, mientras que otros enfoques de machine learning, conocidos como aprendizaje superficial, se enfocan únicamente en una capa de representación de datos (Chollet, 2018).

Una de las principales ventajas del deep learning es que automatiza uno de los pasos más importantes en Machine Learning, conocido como feature engineering o ingeniería de características, el cual consiste en la selección o extracción de un conjunto de características de los datos relevantes para el proceso de aprendizaje, este paso implica creatividad y conocimiento de los datos por parte del analista (Burkov, 2019). Por lo tanto, al usar una técnica de deep learning el analista no necesita realizar una selección o extracción de las características de los datos de entrada, sino que la red capa por capa aprenderá representaciones de los datos cada vez más complejas (Chollet, 2018).

En este proyecto de investigación se empleó un algoritmo de aprendizaje profundo tanto para la clasificación como para la cuantificación de sustancias explosivas, denominado, Long Short-Term Memory (LSTM), importante en el aprendizaje profundo para series de tiempo, cuya descripción y resultados se presentan en el Capítulo 7.

## 2.3. Proceso para la Elaboración de Modelos de Machine Learning y Deep Learning

En esta sección se presentan las partes más importantes de un sistema de aprendizaje automático de Machine Learning y Deep Learning. En la Figura 3a, se muestra el flujo de trabajo para el desarrollo de modelos con técnicas de Machine Learning y en la Figura 3b con técnicas de Deep Learning, como se observa, el flujo de trabajo de los modelos de Machine Learning es similar al de Deep Learning, con la diferencia que en Deep Learning no es necesaria la ingeniería de características, tal como se menciona en la Subsección 2.2.1. Las partes del proceso para la elaboración de los modelos se discuten a continuación.



**Figura 3.** Flujo de trabajo de un modelo

(a) de machine learning y (b) de deep learning

### 2.3.1. Preprocesamiento

Para desarrollar modelos predictivos es necesario que los datos con los que se entrenan los modelos permitan un desempeño óptimo. Para cumplir este objetivo los datos en bruto son

preprocesados. Cada una de las etapas de preprocesamiento realizadas en el presente proyecto de investigación son descritas en el Capítulo 3.

### **2.3.2. Ingeniería de Características**

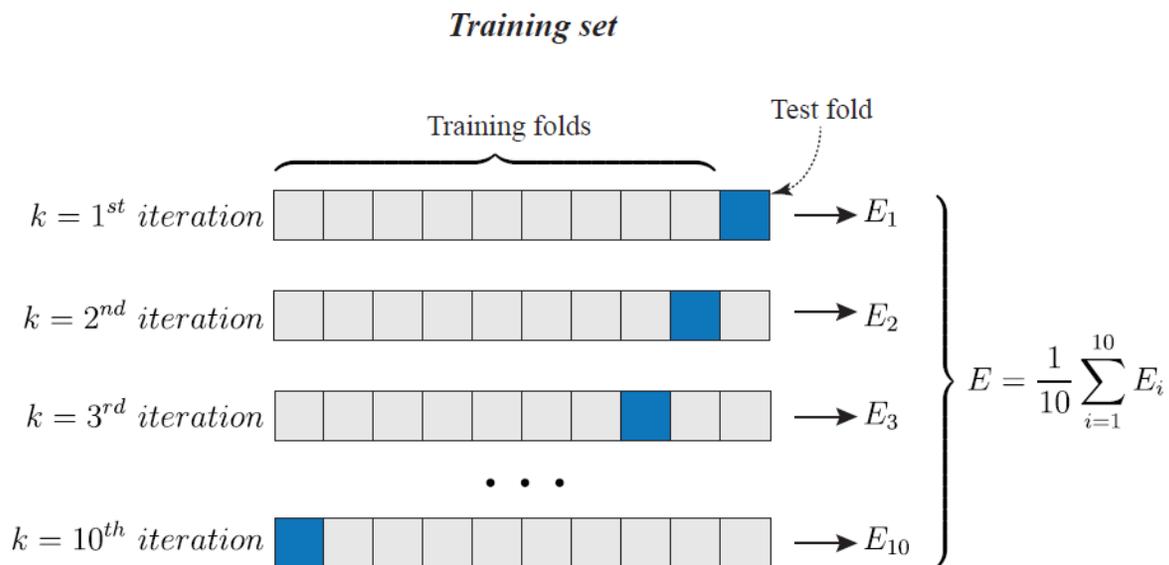
La ingeniería de características consiste en usar el conocimiento de los datos y del algoritmo de aprendizaje para transformar los datos y que sean más fáciles de interpretar por los algoritmos de aprendizaje (Chollet, 2018). Este proceso implica la extracción o la selección de características que son formas de reducción de dimensionalidad (Geron, 2017). Por medio de la selección de características se selecciona un subconjunto de las características originales y por medio de la extracción de características se deriva información del conjunto de características para la construcción de un nuevo subespacio de las mismas (Raschka & Mirjalili, 2017).

### **2.3.3. Aprendizaje**

Antes del proceso de aprendizaje se deben seleccionar los parámetros del algoritmo que no pueden ser aprendidos durante el entrenamiento del modelo, a los cuales se los conoce como hiperparámetros, estos pueden ser el número de variables latentes para la regresión de mínimos cuadrados parciales o el número de épocas en una red neuronal. Los hiperparámetros, al ser configurados manualmente en base al rendimiento en los datos de prueba, corren el riesgo de sobreajustarse a estos datos, es por ello, que es necesario separar el conjunto de datos de entrenamiento en conjuntos de entrenamiento y validación, de forma que el entrenamiento se realice con el primer grupo y la evaluación del desempeño del mismo con el de validación, y únicamente utilizar el conjunto de prueba para evaluar el desempeño final del modelo. La descripción de estos conceptos se presenta a continuación:

### 2.3.3.1. Validación cruzada

Las técnicas de validación cruzada consisten en la división del conjunto de datos de entrenamiento en dos subconjuntos: de entrenamiento y validación. En el presente proyecto de investigación se utilizó la validación cruzada denominada k-fold o de k-iteraciones, en la cual un conjunto de datos, denominado conjunto de prueba es separado del conjunto de datos de entrenamiento para la evaluación final de los modelos y el conjunto de entrenamiento es dividido en  $k$  conjuntos más pequeños.  $K-1$  subconjuntos son usados como datos de entrenamiento y un subconjunto como datos de prueba. La medida de desempeño del algoritmo es calculada mediante el promedio del desempeño de cada uno de los  $k$  modelos entrenados. Como se observa en la Figura 4, se seleccionó un  $k=10$  debido a que en la práctica se ha comprobado que permite estimar modelos que no sufren de subajustes ni sobreajustes a los datos de entrenamiento (J. Brownlee, 2018).



**Figura 4.** Validación cruzada de diez iteraciones k=10 fold

Fuente: (Raschka & Mirjalili, 2017)

### 2.3.3.2. Métricas de desempeño

Las métricas para evaluar el desempeño de los algoritmos deben ser seleccionadas en función de la tarea del modelo, ya sea de clasificación o regresión. Para tareas de clasificación, el desempeño de los algoritmos de aprendizaje puede resumirse en una matriz de confusión, como se observa en la Figura 5. En esta matriz, cada fila representa las clases reales del conjunto de datos evaluado y cada columna las clases predichas por el algoritmo de aprendizaje (Clase 1: Positivo y Clase 2: Negativo). Además, la diagonal de color blanco indica las predicciones correctas, mientras que la otra indica los errores de predicción (Flach, 2012). Si una observación positiva es predicha como positiva se identifica como un verdadero positivo (TP), caso contrario, se define como un falso negativo (FN), en el caso de una observación negativa predicha como negativa esta será un verdadero negativo (TN) y al ser predicha como positiva un falso positivo (FP).

		CLASE PREDICHA	
		Positivo	Negativo
CLASE REAL	Positivo	Verdadero Positivo (TP)	Falso Negativo (FN)
	Negativo	Falso Positivo (FP)	Verdadero Negativo (TN)

**Figura 5.** Matriz de confusión para dos clases

Las filas representan las clases reales, las columnas las clases predichas, la diagonal de color blanco indica las predicciones correctas, mientras que la otra diagonal indica los errores de predicción.

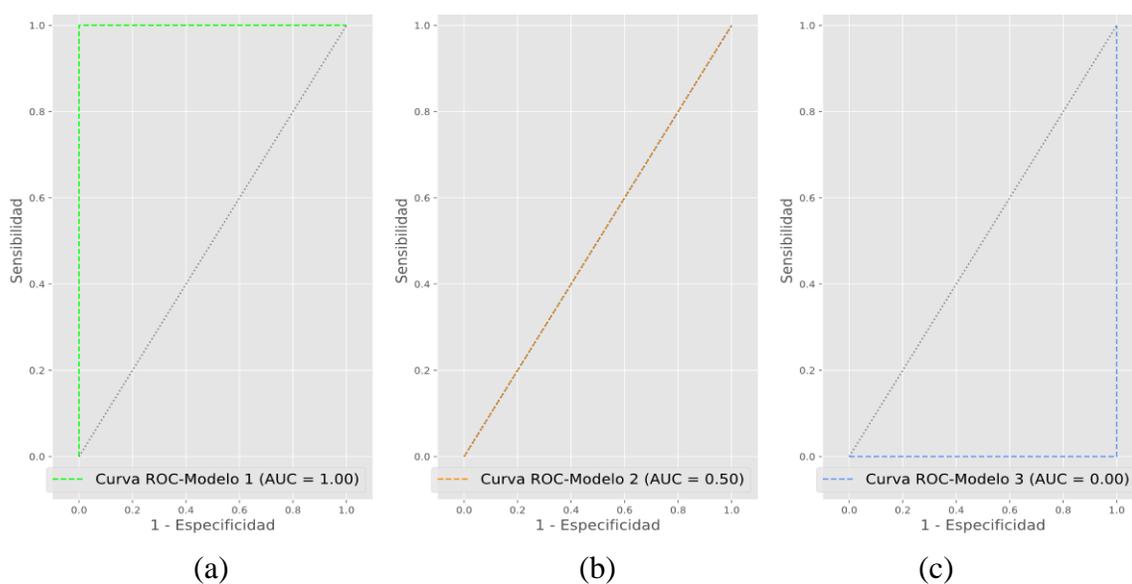
A partir de la matriz de confusión se puede calcular una variedad de métricas de desempeño como aquellas presentadas en la **Tabla 3**, de las cuales la sensibilidad y especificidad se relacionan con la métrica del área bajo la curva (AUC) seleccionada para la evaluación de los modelos de clasificación. La descripción esta métrica se presenta a continuación.

**Tabla 3**  
*Métricas de desempeño para tareas de clasificación*

MÉTRICAS	FÓRMULA	
<b>Error</b>	$(FP + FN)/(TP + TN + FP + FN)$	( 1 )
<b>Exactitud</b>	$(TP + TN)/(TP + TN + FP + FN) = 1 - error$	( 2 )
<b>Tasa de verdaderos positivos (TP rate)</b>	$TP/(TP+FN)$	( 3 )
<b>Tasa de falsos positivos (FP rate)</b>	$FP/(TN + FP)$	( 4 )
<b>Precisión</b>	$TP/(TP + FP)$	( 5 )
<b>Exhaustividad (recall)</b>	$TP/(TP + FN) = TP rate$	( 3 )
<b>F1</b>	$2 \times (precisión \times recall)/(precisión + recall)$	( 6 )
<b>Sensibilidad</b>	$TP/(TP + FN) = TP rate$	( 3 )
<b>Especificidad</b>	$TN/(TN + FP) = 1 - FP rate$	( 7 )

- **Sensibilidad/Exhaustividad/Tasa de verdaderos positivos:** Definida mediante la Ecuación 3, es la proporción de observaciones positivas que el modelo clasifica correctamente (Geron, 2017). Por ejemplo, en la clasificación de sustancias explosivas informara que tan bien las sustancias explosivas son identificadas como tal, sin importar el desempeño del modelo en la clasificación de sustancias no explosivas.
- **Especificidad:** La especificidad presentada en la Ecuación 7 explica que tan bien se detectan las observaciones negativas (Alpaydın, 2012). Por ejemplo, en la clasificación de sustancias explosivas informara que tan bien las sustancias no explosivas son detectadas como tal, sin importar el desempeño del modelo en la clasificación de sustancias que si son explosivas.
- **Área bajo la curva (AUC):** El área bajo la curva de la curva ROC (Característica Operativa del Receptor), es usada para evaluar modelos de clasificación en base a su desempeño respecto a la tasa de falsos positivos (1-especificidad) y a la tasa de verdaderos positivos (sensibilidad/exhaustividad) para diferentes umbrales de decisión del clasificador (Raschka & Mirjalili, 2017). Un clasificador perfecto, como aquel presentado en la Figura 6a, tendrá

un AUC de 1, es decir que su tasa de TP será uno y de FP cero. En la Figura 6b, se observa el resultado de un modelo con un  $AUC=0.50$ , en el cual la tasa de TP es proporcional a la de FP, por lo cual se establece que el modelo no es capaz de distinguir entre una clase y otra. Finalmente, en la Figura 6c, se observa un modelo con el peor desempeño posible ( $AUC=0$ ) en el cual la tasa de FP es uno y de TP cero, por lo cual se puede concluir que el modelo clasifica a las observaciones negativas como positivas y a las positivas como negativas. Es por ello que para seleccionar entre un modelo y otro se preferirá aquel cuya curva ROC se encuentra más cerca de la esquina superior izquierda (con una tasa de TP mayor a la de FP), con lo cual se obtendrá un AUC cercano a 1 (Alpaydm, 2012).



**Figura 6.** Curvas ROC para diferentes clasificadores

Una de las ventajas de usar esta métrica es que se podrá evaluar el desempeño del modelo para diferentes valores de umbral, por lo tanto, la evaluación del modelo no se verá limitada por la elección de este.

Para evaluar los modelos de regresión, al igual que para los modelos de clasificación, existen diversas métricas, de las cuales se seleccionaron dos: el coeficiente de determinación (R2) y el error cuadrático medio (MSE). Cuya descripción se presenta a continuación.

- **Error cuadrático medio (MSE):** Calcula el error cuadrático medio del modelo. Es útil para comparar diferentes modelos de regresión o ajustar sus parámetros a través de una validación cruzada ya que se normaliza la suma de errores cuadrados por el tamaño de la muestra (Raschka & Mirjalili, 2017). La fórmula para su cálculo se presenta a en la Ecuación 8.

$$MSE = \frac{1}{n} - \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad ( 8 )$$

Donde  $y_i$  se define como el valor real de salida de la observación  $i$ ,  $\hat{y}_i$  como el valor predicho de la observación  $i$  y  $n$  el número de muestras.

- **Coficiente de determinación (R2):** Es la versión estandarizada del MSE, presenta la proporción de varianza explicada por las variables independientes del modelo y proporciona una medida de que tan bien es probable que el modelo prediga nuevos datos a través de la proporción de varianza explicada. La mejor puntuación posible es 1, y puede ser negativa cuando el modelo es arbitrariamente peor (Pedregosa et al., 2011). La fórmula para su cálculo es presentada en la Ecuación 9.

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad ( 9 )$$

Donde  $y_i$  se define como el valor real de salida de la muestra  $i$ ,  $\hat{y}_i$  como el valor predicho de la muestra  $i$ ,  $\bar{y}$  el promedio de los valores reales de salida, y  $n$  el número de muestras.

### 2.3.4. Evaluación

Esta etapa es utilizada para evaluar el desempeño final del modelo después de la etapa de aprendizaje en datos que no han sido vistos por el mismo para posteriormente poder seleccionar el modelo con mejor desempeño. La evaluación del desempeño final de los modelos del proyecto de investigación se realizó mediante la comparación del desempeño de los modelos en base a las métricas presentadas en la **Tabla 4**.

**Tabla 4**

*Métricas de desempeño del proyecto*

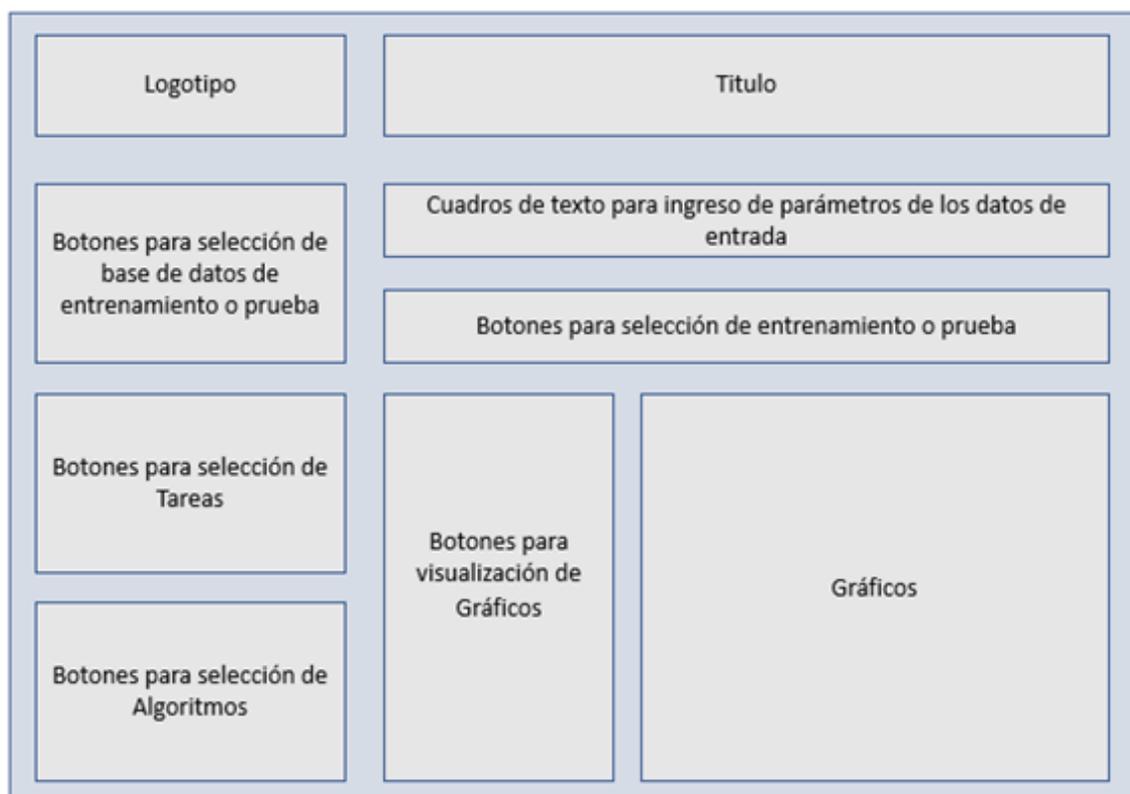
Algoritmos de clasificación	Algoritmos de Regresión
Área bajo la curva (AUC)	Coeficiente de determinación (R <sup>2</sup> )
	Error cuadrático medio (MSE)

## 2.4. Desarrollo de Interfaz Gráfica de Usuario

Para el desarrollo de la interfaz gráfica de usuario se importó el módulo *Tkinter*, una adaptación de la biblioteca grafica Tcl/Tk para *Python*. El cual proporciona las clases de aplicaciones y las constantes asociadas (Lundh, 1999).

La interfaz del proyecto se estructuró en base a las funciones que debe cumplir, como se observa en la Figura 7, en la parte superior de la pantalla se encuentra el logotipo y el título del proyecto utilizados como presentación de la interfaz. Los botones del lado izquierdo se encuentran divididos en tres grupos: el primero permite seleccionar las bases de datos que serán utilizadas para entrenar o evaluar los modelos, el segundo grupo permite seleccionar la tarea que se desea el modelo realice (tareas definidas en la **Tabla 9**) y el tercer grupo permite seleccionar el algoritmo que será

utilizado (PLS-DA, Regresión Logística, entre otros). Finalmente, el lado derecho de la pantalla se encuentra dividido en tres partes: en la primera se encuentran cuadros de texto para ingresar la información correspondiente a los datos de entrada (número de muestras por experimento, tiempo de respiración y frecuencia de muestreo), en la parte inferior se encuentran los botones para la selección del entrenamiento o prueba del modelo, y la tercera parte en el lado derecho permite visualizar las gráficas seleccionadas con los botones del lado izquierdo durante la etapa de preprocesamiento, después de esta etapa permite visualizar las gráficas resultantes del entrenamiento o prueba de los modelos



**Figura 7.** Estructura de la interfaz gráfica de usuario

En este capítulo se describió el prototipo nariz electrónica del cual se adquirió las bases de datos para el desarrollo de los modelos multivariantes. Además, se presentaron los conceptos

utilizados para la creación de los modelos, como la definición de aprendizaje de máquina y deep learning. Posteriormente, se describieron los conceptos básicos para el desarrollo de los modelos como lo es la etapa de preprocesamiento, aprendizaje, ingeniería de características y evaluación. Finalmente se presentó una descripción de la herramienta utilizada para la elaboración de la interfaz gráfica de usuario y su estructura.

# CAPÍTULO 3

## Metodología Experimental

En este capítulo se presenta una descripción de los datos utilizados para el desarrollo de los modelos multivariantes, tanto de la base de datos inicial realizada en (Salazar, 2018) y de la actual realizada en (Jacome, 2019). Posterior a ello, se describe la forma en que fueron distribuidos para el entrenamiento y evaluación de los modelos. Finalmente se describe cada una de las etapas que forman parte del preprocesamiento de estos datos.

### 3.1 Descripción de los Experimentos

Para el desarrollo de los modelos de clasificación y cuantificación de sustancias explosivas con el prototipo e-nose descrito en el Capítulo 2, se emplearon dos bases de datos ya que a pesar de que el número y tipo de sensores utilizados en la elaboración de la base de datos 1 y 2 es el mismo, las resistencias utilizadas para el acondicionamiento de la señal de salida de los sensores, la frecuencia de muestreo de los datos y el proceso de adquisición de los mismos cambio debido a optimizaciones del prototipo e-nose realizadas en (Jacome, 2019), es por ello que se decidió desarrollar los modelos de clasificación y cuantificación de sustancias explosivas de forma independiente para la base de datos 1 realizada en (Salazar, 2018) y para la base de datos 2 realizada con el prototipo optimizado en (Jacome, 2019). Para poder de esta forma poder comparar el poder predictivo de la nariz electrónica antes de su optimización y después de esta con cantidades superiores de sustancia explosiva.

La base de datos inicial realizada en (Salazar, 2018) cuenta con observaciones de TNT y pólvora en estado puro y combinados con 1 gr de pasta dental y 1gr jabón, con concentraciones entre 0.1 g y 3 g, además de experimentos con únicamente 2ml de sustancia dopante (alcohol), los mismos se encuentran clasificados por concentración en la Tabla 5 y por sustancia en la **Tabla 6**. El objetivo de utilizar esta base de datos es evaluar el desempeño de los modelos para bajas concentraciones de sustancias explosiva y cuando existen mezclas de estas sustancias.

**Tabla 5**

*Base de datos 1- Experimentos clasificados por concentración*

<b>SUSTANCIA</b>	<b>PÓLVORA</b>										
Concentración (g)	0.1	0.2	0.3	0.4	0.5	0.7	1	1.5	2	2.5	3
Cantidad de experimentos	5	3	5	3	5	5	5	3	3	3	3
<b>TOTAL</b>	<b>43</b>										

<b>SUSTANCIA</b>	<b>TNT</b>										
Concentración (g)	0.1	0.2	0.3	0.4	0.5	0.7	1	1.5	2	2.5	3
Cantidad de experimentos	5	5	5	5	5	5	5	5	5	3	3
<b>TOTAL</b>	<b>51</b>										

<b>SUSTANCIA</b>	<b>ALCOHOL</b>
Concentración (g)	0
Cantidad de experimentos	48
<b>TOTAL</b>	<b>48</b>

**Tabla 6**

*Base de datos 1- Experimentos clasificados por sustancia*

<b>SUSTANCIA</b>	<b>CANTIDAD DE EXPERIMENTOS</b>
Alcohol	48
Pólvora en estado puro	33
TNT en estado puro	33
Pólvora y Jabón	5
Pólvora y Pasta Dental	5
TNT y Jabón	9
TNT y Pasta Dental	9
<b>TOTAL</b>	<b>142</b>

La base de datos actual realizada en (Jacome, 2019), cuenta con observaciones de TNT y pólvora en estado puro con concentraciones entre 3 g y 5 g además de experimentos con únicamente 1 ml de sustancia dopante (alcohol), los mismos se encuentran clasificados por sustancia en la Tabla 7 y por concentración en la Tabla 8. El objetivo de realizar estos nuevos experimentos es analizar el poder predictivo del prototipo e-nose con mayores concentraciones de sustancia explosiva.

**Tabla 7**

*Base de datos 2 - Experimentos clasificados por sustancia*

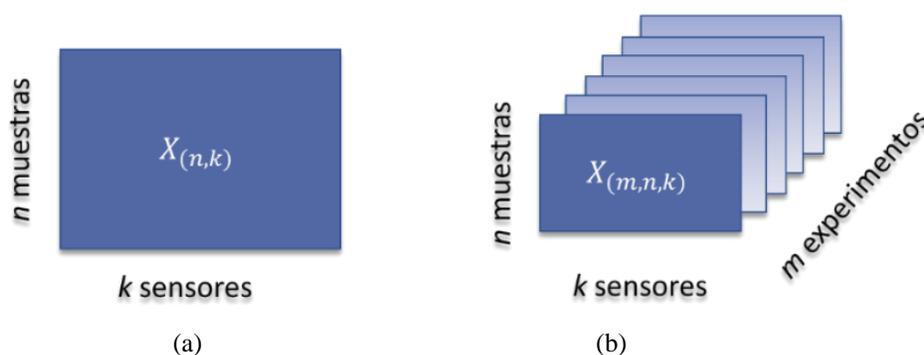
SUSTANCIA	CANTIDAD DE EXPERIMENTOS
Alcohol	21
Pólvora en estado puro	21
TNT en estado puro	21
<b>TOTAL</b>	<b>63</b>

**Tabla 8**

*Base de datos 2 - Experimentos clasificados por concentración*

SUSTANCIA	PÓLVORA			TNT			ALCOHOL
Concentración de sustancia explosiva (g)	3	4	5	3	4	5	0
Cantidad de experimentos	7	7	7	7	7	7	21
<b>TOTAL</b>	<b>63</b>						

Cada uno de los experimentos está conformado por la respuesta de los seis sensores mencionados en el Capítulo 2 durante un tiempo de medición  $tm$  (650s para la base de datos 1 y 180s para la base de datos 2), es por ello que estos experimentos son identificados como series temporales. Por lo tanto, cada observación antes de su preprocesamiento está conformada por una matriz  $X_{(n,k)}$  de dos dimensiones con  $n$  muestras y  $k$  sensores, como se muestra en la Figura 8a. En conjunto todos los experimentos utilizados para el desarrollo de los modelos conforman una matriz  $X_{(n,m,k)}$  de tres dimensiones con  $m$  experimentos,  $n$  muestras y  $k$  sensores, como se muestra en la Figura 8b.

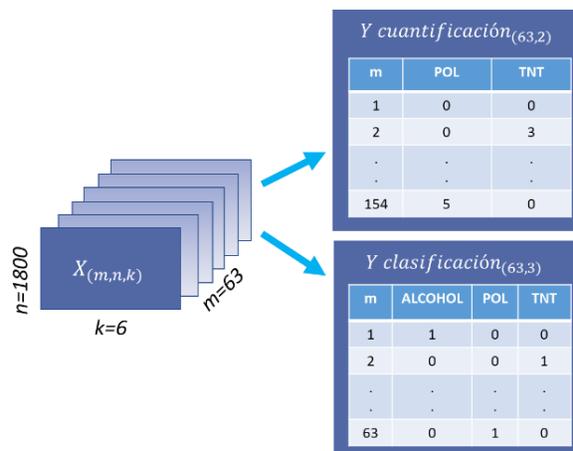


**Figura 8.** Representación de los datos recolectados con el prototipo e-nose

(a) Representación de los datos por experimento, (b) Representación de todo el conjunto de datos

## 3.2. Distribución de los Datos en Conjuntos de Entrenamiento y Prueba

Los experimentos presentados en la sección anterior de la base de datos 1 y 2 se dividieron de en dos conjuntos: uno denominado conjunto de entrenamiento utilizado para entrenar y optimizar los modelos, y otro denominado conjunto de prueba utilizado para evaluar el desempeño de los modelos finales en datos que no hayan sido utilizados durante su desarrollo conocida como validación ciega. Para realizar esta división se asignaron las etiquetas categóricas (Pólvora, TNT, mezcla de TNT, mezcla de pólvora o alcohol) y continuas (concentración de sustancia explosiva en gramos: 0, 1, 2, etc.) a cada experimento, como se observa en la Figura 9.



**Figura 9.** Asignación de etiquetas categóricas y continuas a los experimentos de la base de datos 2

A continuación, se seleccionaron las tareas de los modelos multivariantes, cuyos resultados permitirán analizar el poder predictivo del prototipo e-nose en los escenarios que se muestran en la Tabla 9. Como se observa en la tabla, se seleccionaron ocho tareas de cuantificación y clasificación.

**Tabla 9**

*Selección de las tareas de clasificación y regresión de los modelos multivariantes*

N° DE TAREA	SUSTANCIAS QUE FORMAN PARTE DE LA TAREA DE CLASIFICACIÓN Y REGRESIÓN
1	Alcohol/Sustancias explosivas en estado puro y mezclas
2	Alcohol/Sustancias explosivas en estado puro
3	Alcohol/Pólvora en estado puro
4	Alcohol/TNT en estado puro
5	Pólvora en estado puro/TNT en estado puro
6	Alcohol/Pólvora en estado puro y mezcla de pólvora
7	Alcohol/TNT en estado puro y mezcla de TNT
8	Alcohol/Pólvora en estado puro/TNT en estado puro

Finalmente, los experimentos se permutaron aleatoriamente de forma estratificada y se tomó el 70% como conjunto de entrenamiento y el 30% restante como conjunto de prueba, para cada una de las tareas. El conjunto de entrenamiento de la tarea 1 para la base de datos 1 se presenta clasificado por concentración en la **Tabla 10a** y el conjunto de prueba en la **Tabla 10b**. En estas tablas se comprueba la distribución estratificada de los datos ya que se mantienen aproximadamente

las mismas proporciones de etiquetas de clase y concentración de sustancia tanto en los conjuntos de prueba y entrenamiento como en el conjunto de datos original, que evitará un error en la predicción de los modelos al capturar más la relación entre características de una clase que de otra.

**Tabla 10**

*Base de datos 1- Experimentos clasificados por concentración de sustancia*

*(a) Conjunto de entrenamiento (b) Conjunto de prueba*

<b>SUSTANCIA</b>	<b>PÓLVORA</b>										
Concentración (g)	0.1	0.2	0.3	0.4	0.5	0.7	1	1.5	2	2.5	3
Cantidad de experimentos	3	2	4	2	4	3	4	2	2	2	2
<b>TOTAL</b>	<b>30</b>										

<b>SUSTANCIA</b>	<b>TNT</b>										
Concentración (g)	0.1	0.2	0.3	0.4	0.5	0.7	1	1.5	2	2.5	3
Cantidad de experimentos	4	3	3	3	4	4	4	4	3	2	2
<b>TOTAL</b>	<b>36</b>										

<b>SUSTANCIA</b>	<b>ALCOHOL</b>										
Concentración (g)	0										
Cantidad de experimentos	33										
<b>TOTAL</b>	<b>33</b>										

(a)

<b>SUSTANCIA</b>	<b>PÓLVORA</b>										
Concentración (g)	0.1	0.2	0.3	0.4	0.5	0.7	1	1.5	2	2.5	3
Cantidad de experimentos	2	1	1	1	1	2	1	1	1	1	1
<b>TOTAL</b>	<b>13</b>										

<b>SUSTANCIA</b>	<b>TNT</b>										
Concentración (g)	0.1	0.2	0.3	0.4	0.5	0.7	1	1.5	2	2.5	3
Cantidad de experimentos	1	2	2	2	1	1	1	1	2	1	1
<b>TOTAL</b>	<b>15</b>										

<b>SUSTANCIA</b>	<b>ALCOHOL</b>										
Concentración (g)	0										
Cantidad de experimentos	15										
<b>TOTAL</b>	<b>15</b>										

(b)

### 3.3. Preprocesamiento de las Señales de Entrenamiento

La calidad de los datos de entrenamiento determinará la efectividad con que los algoritmos de machine learning pueden aprender de estos, por lo tanto, es fundamental que sean examinados y preprocesados antes de la etapa de aprendizaje (Raschka & Mirjalili, 2017). Es por ello que primero se examinaron los datos tanto de la base de datos 1 como de la base de datos 2 y posteriormente se elaboró un algoritmo en Python para el preprocesamiento de los datos de entrenamiento que consta de las siguientes etapas: implementación de filtro, corrección de línea base, alineamiento de picos, concatenación de sensores, detección y eliminación de outliers, las mismas se describen a continuación.

#### 3.3.1. Implementación de Filtro

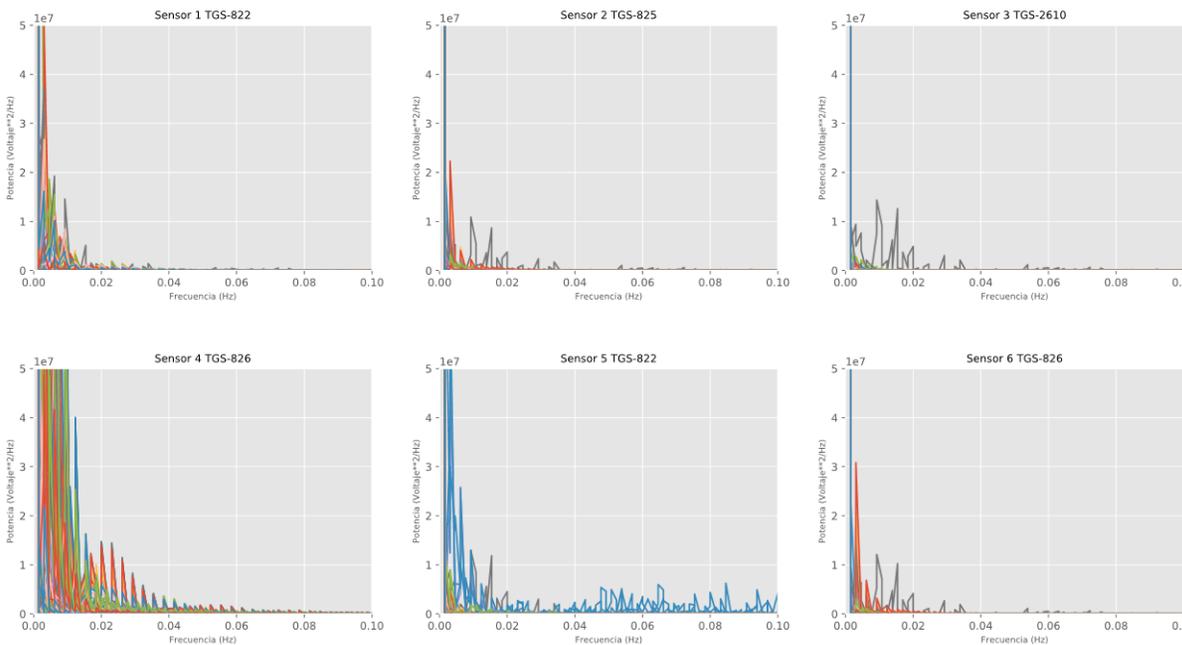
El prototipo e-nose puede estar expuesto a perturbaciones, ya sea por señales externas, vibraciones por parte del dispositivo durante el proceso de adquisición de las señales o por ruido térmico generado por los componentes del mismo (Salazar, 2018). Es por esto que se realizó el análisis espectral de las señales con el fin de seleccionar el tipo de filtro a implementar que permita disminuir el ruido y mejorar la característica de las señales.

El análisis espectral se realizó mediante el cálculo de la transformada rápida de Fourier (FFT) de las señales, con la función `numpy.fft.rfft` de Python con la que se obtuvo los valores de las señales en función de la frecuencia. A continuación, se obtuvo el espectro de potencia para visualizar cuánta potencia está contenida en los componentes de frecuencia de las señales mediante la Ecuación 10 y se calculó las frecuencias asociadas a los componentes mediante la función `np.fft.fftfreq` de la Ecuación 11, en la que  $n$  representa el número de muestras de las señales y  $d$  el espaciado de las mismas (inversa de la frecuencia de muestreo).

$$P = |fft(x)|^2 \quad ( 10 )$$

$$f = \frac{\left[0, 1, \dots, \frac{n}{2} - 1, -\frac{n}{2}, \dots, -1\right]}{(d \times n)} \quad ( 11 )$$

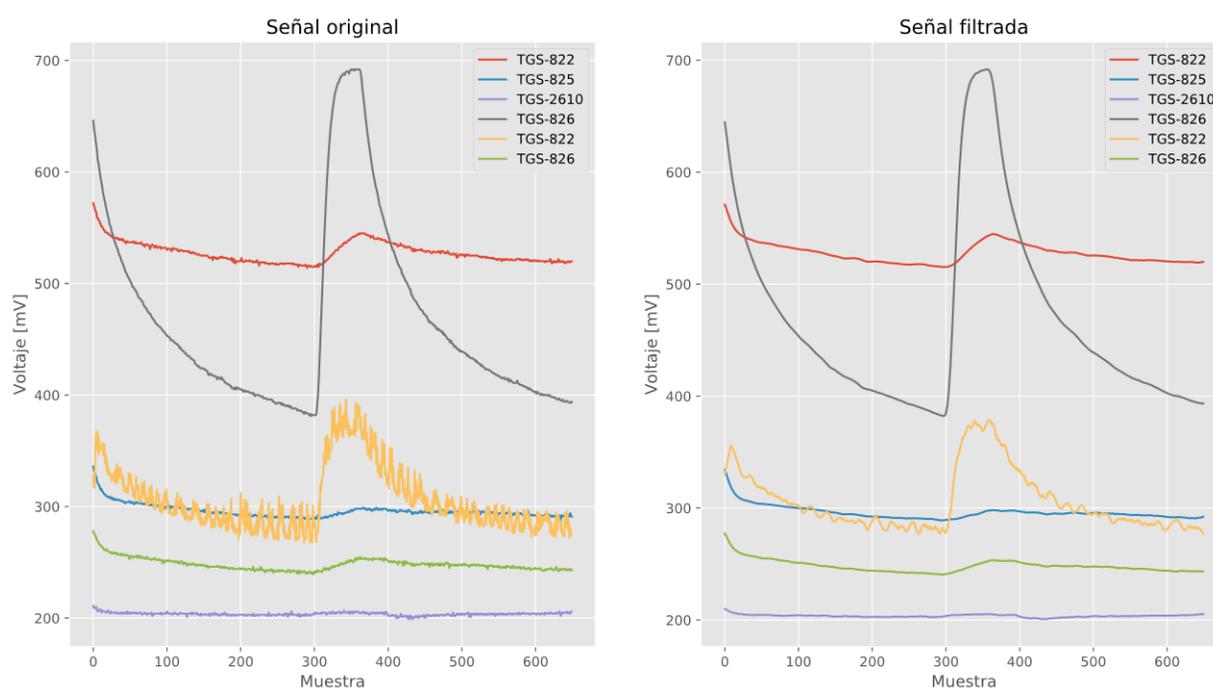
Al analizar el espectro de potencia de cada uno de los sensores para la base de datos 1, se determinó que la potencia de las señales está contenida en bajas frecuencias, como se muestra en la Figura 10. Además, a pesar de que el sensor 1 y 5 son TGS-822, el sensor 5 que previamente se determinó proporcionó señales con ruido en la mayoría de los experimentos, cuenta con componentes en alta frecuencia. Es por ello que se decidió utilizar el filtro de Savitzky–Golay, un tipo de filtro pasa bajo que se caracteriza por reducir el ruido al suavizar las señales y hacerlo sin perder la forma original de la misma, manteniendo la forma y altura de los picos (Zuñiga, 2018).



**Figura 10.** Espectro de potencia de los datos de entrenamiento-Base de datos 1

El filtro de Savitzky–Golay realiza el suavizado de datos mediante una aproximación polinómica de mínimos cuadrados locales, a través del ajuste de un polinomio de grado  $N$  a un

conjunto de  $M$  muestras de entrada y posteriormente la evaluación del polinomio resultante en un solo punto (Schafer, 2011). El filtro fue implementado mediante la función *savgol\_filter* de la librería *spicy*, se seleccionó un polinomio de grado  $N=1$  para que la línea resultante disminuya el ruido de la señal y una ventana de tamaño  $M=11$  para que exista una cantidad considerable de puntos a los que deba ajustarse el polinomio que evite que la señal se distorsione. La Figura 11, presenta una muestra de alcohol de la base de datos 1 y la señal resultante después de la aplicar el filtro Savitzky–Golay.



**Figura 11.** Experimento de 2ml de alcohol filtrado - Base de datos 1.

### 3.3.2. Corrección de Línea Base

La línea base de los sensores cambia constantemente debido a fluctuaciones en el entorno, estas fluctuaciones se conocen como deriva de la línea base, y en caso de no corregirla disminuirá el desempeño de los modelos de machine learning ya que su valor variará de un experimento en

otro (Zhang & Peng, 2016). Es por ello que es importante realizar una compensación de la deriva de la línea base de las señales antes del desarrollo de los modelos.

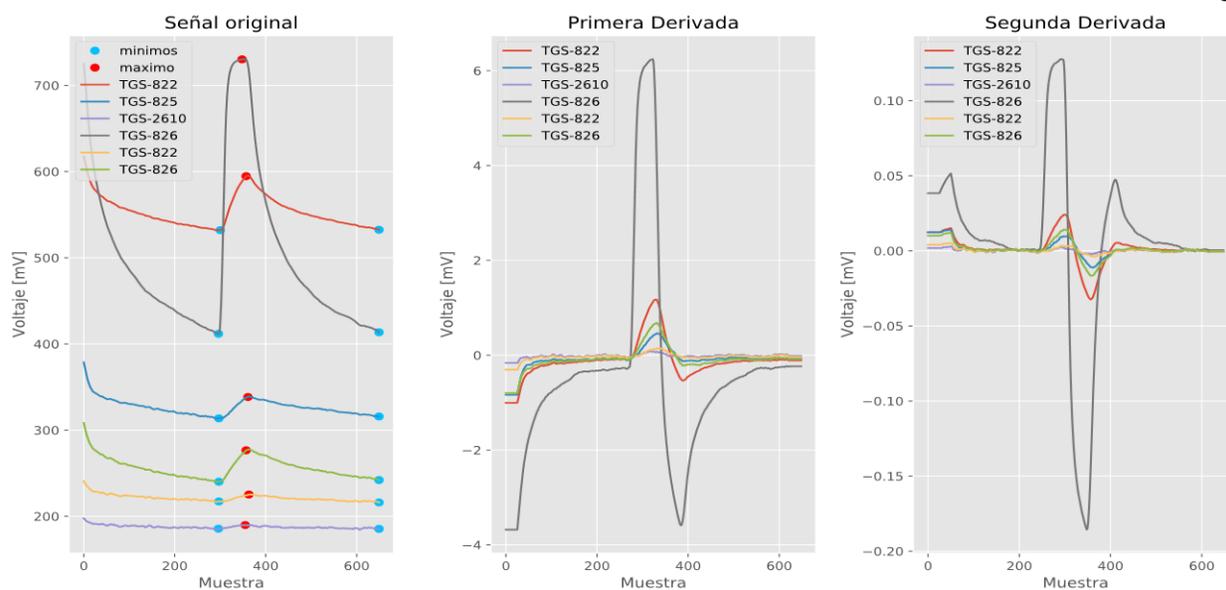
La compensación de la deriva de la línea base se realizó mediante la resta de la señal original con una estimación de línea base de la misma, en este caso un polinomio de grado 1, de forma que todas las señales tengan su línea base en el origen. Los coeficientes del polinomio de la Ecuación 12 se obtuvieron con la función *polyfit* de la librería *numpy* la cual realiza un ajuste del polinomio a los puntos  $(x,y)$  de forma que los coeficientes  $(a,b)$  minimicen el error cuadrado de la Ecuación 13 en la Ecuación 14.

$$p(x) = ax + b \quad ( 12 )$$

$$E = \sum_{i=0}^k (p(x[i]) - y[i])^2 \quad ( 13 )$$

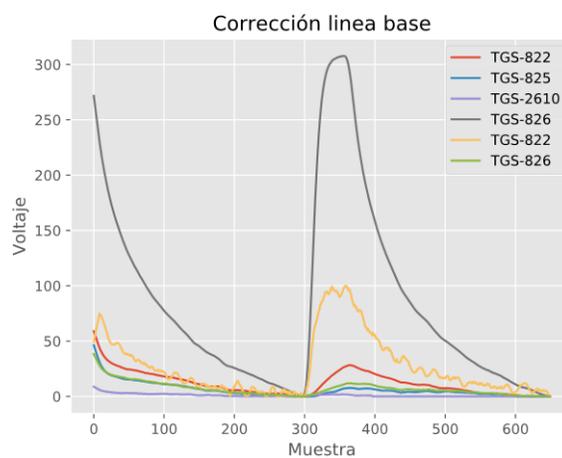
$$\begin{cases} ax[0] + b = y[0] \\ ax[1] + b = y[1] \end{cases} \quad ( 14 )$$

Para obtener los puntos  $(x,y)$  a los que el polinomio debe ajustarse se calculó la derivada de todas las señales de entrenamiento ya que esta muestra la evolución de la pendiente a lo largo de la curva, por lo tanto, permite calcular los puntos dónde la pendiente es cero ( $X_{FILTRADA}' = 0$ ), es decir, los máximos y mínimos de las señales. Además, se calculó la segunda derivada de la señal para identificar si el punto crítico en el cual  $X_{FILTRADA}' = 0$  es un punto mínimo o máximo. El resultado de la identificación de los mínimos y el máximo de cada una de las señales para un experimento se muestra en la Figura 12.



**Figura 12.** Identificación de puntos máximos y mínimos de las señales de un experimento - Base de datos 1

Con los puntos mínimos identificados, se procedió a utilizar estos valores para calcular la línea base estimada de cada una de las señales y restarla de la señal original, obteniendo las señales con una misma línea base en el origen, como se observa en la Figura 13.



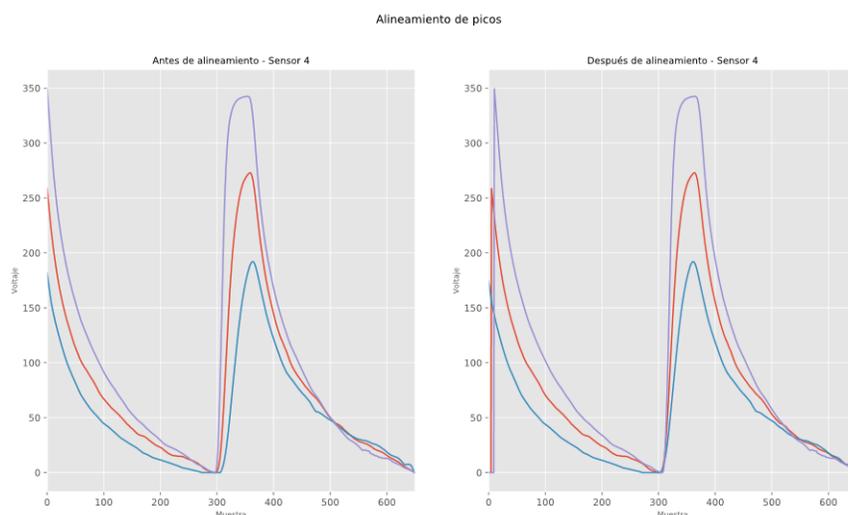
**Figura 13.** Señales de un experimento con líneas base corregidas - Base de datos 1

### 3.3.3. Alineamiento de Picos

Durante el tiempo de respiración el prototipo e-nose envía las muestras de olor de las sustancias de interés a la cámara de sensores, lo cual produce un incremento de la conductividad de los mismos hasta que el aire dentro de la cámara es saturado por la sustancia de interés, esto es representado por el pico de las señales de cada uno de los sensores. Estos picos pueden encontrarse desalineados entre un experimento y otro debido a desplazamientos temporales durante la adquisición de las señales de los sensores. Por lo tanto, se decidió alinear los picos en función del sensor al que pertenecen de forma que la respuesta de los sensores cuando se encuentran saturados por la sustancia de interés este ubicada en el mismo punto y puedan ser comparables entre un experimento y otro.

Para realizar el alineamiento de los picos, en primer lugar, se calculó la muestra promedio en la que se encuentran los picos por sensor para establecer como punto de alineamiento de las señales. A continuación, se comparó en cada uno de los experimentos el pico del sensor con el punto de alineamiento correspondiente, si el pico se encontraba en el punto de alineamiento la señal se mantenía en la misma posición, caso contrario, si el pico se encontraba a la izquierda del punto de alineamiento se realizaba un desplazamiento hacia la derecha o un desplazamiento hacia la izquierda si se encontraba a la derecha del punto.

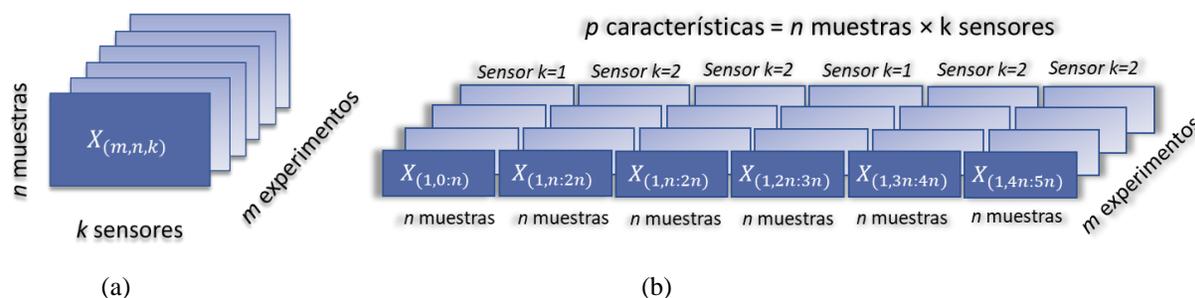
El desplazamiento de las señales se realizó mediante el comando *np.roll* de la librería *numpy*, el cual permite desplazar las señales dado un valor de desplazamiento, positivo para un desplazamiento hacia la derecha y negativo para un desplazamiento hacia la izquierda. El resultado del alineamiento del sensor cuatro para tres experimentos se observa en la Figura 14.



**Figura 14.** Señales correspondientes al sensor 4 alineadas- Base de datos 1

### 3.3.4. Concatenación de Sensores

Para poder ingresar los datos de entrenamiento en los algoritmos de machine learning es necesario agrupar los datos de cada experimento que inicialmente estaba formado por una matriz  $X_{(n,k)}$  de dos dimensiones con  $n$  muestras y  $k$  sensores, como se muestra en la Figura 15a, en un vector  $X_{(1,n \times k)}$  de  $n \times k$  características, como se muestra en la Figura 15b.

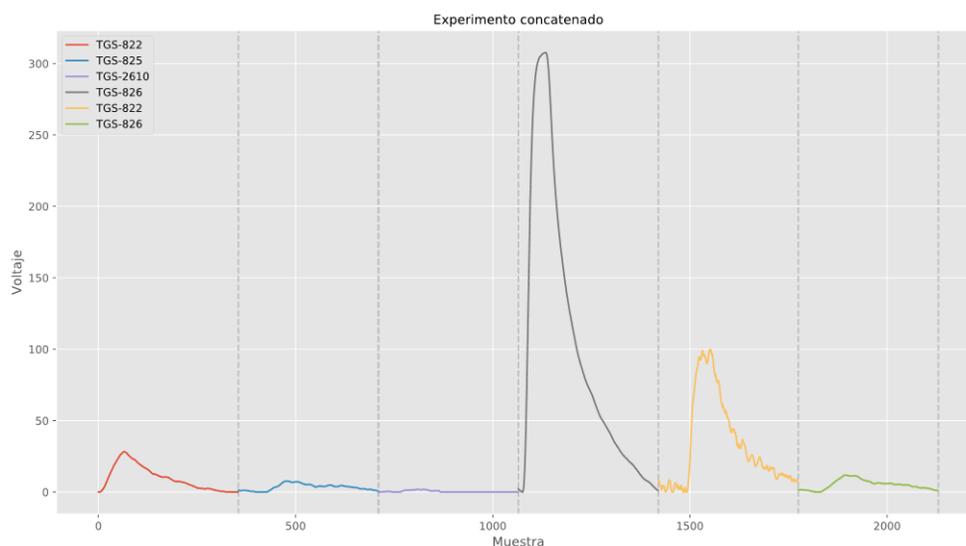


**Figura 15.** Representación del conjunto de datos de entrenamiento

(a) agrupados como una matriz de tres dimensiones  $X_{(m,n,k)}$  de tres dimensiones,

(b) agrupados como una matriz de dos dimensiones  $X_{(m,n \times k)}$

Esta agrupación se realizó mediante la concatenación de las señales de los sensores de cada uno de los experimentos, además, por cada sensor se eliminaron las muestras anteriores al tiempo de respiración del prototipo. El resultado de la concatenación para un experimento se observa en la Figura 16.



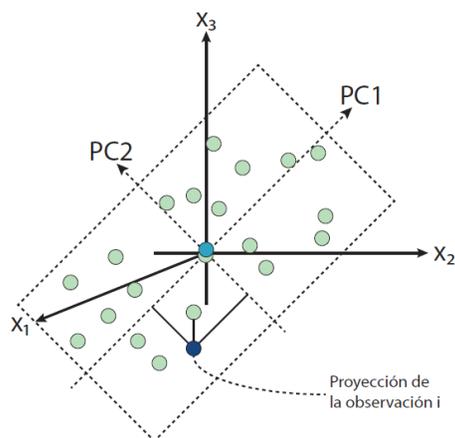
**Figura 16.** Señales de los seis sensores de un experimento concatenadas - Base de datos 1

### 3.3.5. Detección y Eliminación de Outliers

Finalmente, es importante conocer si dentro del conjunto de datos de entrenamiento existen señales que se desvían del resto (a estas señales se las conoce con el nombre de outliers), ya que con esto se podrá identificar si existe alguna característica anormal en el prototipo e-nose que afecta el proceso de generación de datos (Aggarwal, 2013). Por consiguiente, se decidió utilizar el análisis de componentes principales (PCA) para detectar outliers, analizar si un sensor en específico produjo estas señales atípicas y caso de ser necesario eliminarlas antes del desarrollo de los modelos. El proceso de creación de los modelos para la detección de outliers y los conceptos relacionados a su elaboración se definen en la subsección siguiente.

### 3.3.5.1. Análisis de Componentes Principales (PCA)

El análisis de componentes principales (PCA) es un método de transformación lineal no supervisada que permite identificar patrones en los datos de entrenamiento, en base a la correlación entre las características de estos datos. Su objetivo es proyectar los datos en un nuevo subespacio con iguales o menores dimensiones que el original, este nuevo subespacio está formado por ejes ortogonales entre sí, denominados componentes principales (PC), que representan las direcciones de máxima varianza de los datos (Raschka & Mirjalili, 2017). Una explicación geométrica de este análisis se presenta en la Figura 17, donde  $x_1$ ,  $x_2$  y  $x_3$  representan los ejes de las características originales,  $PC1$  el componente principal que va en dirección de la máxima varianza de las proyecciones de cada observación en  $PC1$  y  $PC2$  el segundo componente principal que de igual forma ira en dirección de la máxima varianza dada la restricción de que debe ser ortogonal al primero (Dunn, 2019). La distancia del origen del sistema al punto proyectado de cada observación en los nuevos ejes  $PC1$  y  $PC2$ , se los conoce como scores. Los loadings o vectores dirección, son vectores unitarios que definen la dirección de los componentes principales en el sistema coordenado original.



**Figura 17.** Representación geométrica de modelo PCA

Fuente: (Eriksson et al., 2006)

La explicación matemática del análisis de componentes principales (PCA). se representa mediante la combinación lineal de la Ecuación 15.

$$T_{(N,A)} = X_{(N,K)}P_{(K,A)} \quad ( 15 )$$

Donde  $T_{(N,A)}$  representa la matriz de scores con  $N$  observaciones y  $A$  componentes principales,  $X$  la matriz original de observaciones con  $N$  observaciones y  $K$  características, y finalmente la matriz de vectores dirección, conocido también como loadings, con  $K$  características y  $A$  componentes principales.

Es importante tomar en cuenta que antes del análisis de componentes principales es necesario acondicionarlo previamente, para que los datos tengan una forma más adecuada para el análisis. En la literatura se ha demostrado que el preprocesamiento previo de los datos puede hacer la diferencia en la obtención de un modelo útil de uno que no lo sea (Eriksson et al., 2006).

### 3.3.5.2. $T^2$ de Hotelling

El valor  $T^2$ -de Hotelling es la medida de la variación de cada observación dentro del modelo PCA, cuyo cálculo se realiza mediante la siguiente ecuación:

$$T_i^2 = \sum_{a=1}^A \frac{t_{i,a}^2}{s_a} \quad ( 16 )$$

Donde  $t_{i,a}$  se define como el score del componente principal  $a$  para la observación  $i$  y  $s_a$  como la desviación estándar del componente  $a$  (Dunn, 2019; Mujica, Rodellar, Fernández, & Güemes, 2011).

Para calcular, el valor  $T^2$ -de Hotelling que se encuentra en un nivel de confianza seleccionado se utilizó la Ecuación 17 en la que se establece que el valor de  $T^2$ -de Hotelling es proporcional a la distribución-F.

$$t^2 \sim T_{p,n-1}^2 = \frac{p(n-1)}{n-p} F_{p,n-p} \quad ( 17 )$$

Donde  $n$  es el número de observaciones,  $p$  el número de componentes principales y  $F_{p,n-p}$  el valor crítico de la distribución-F que se calculó con la función *scipy.stats.f.ppf*.

### 3.3.5.3. Desarrollo del Análisis de Componentes Principales (PCA)

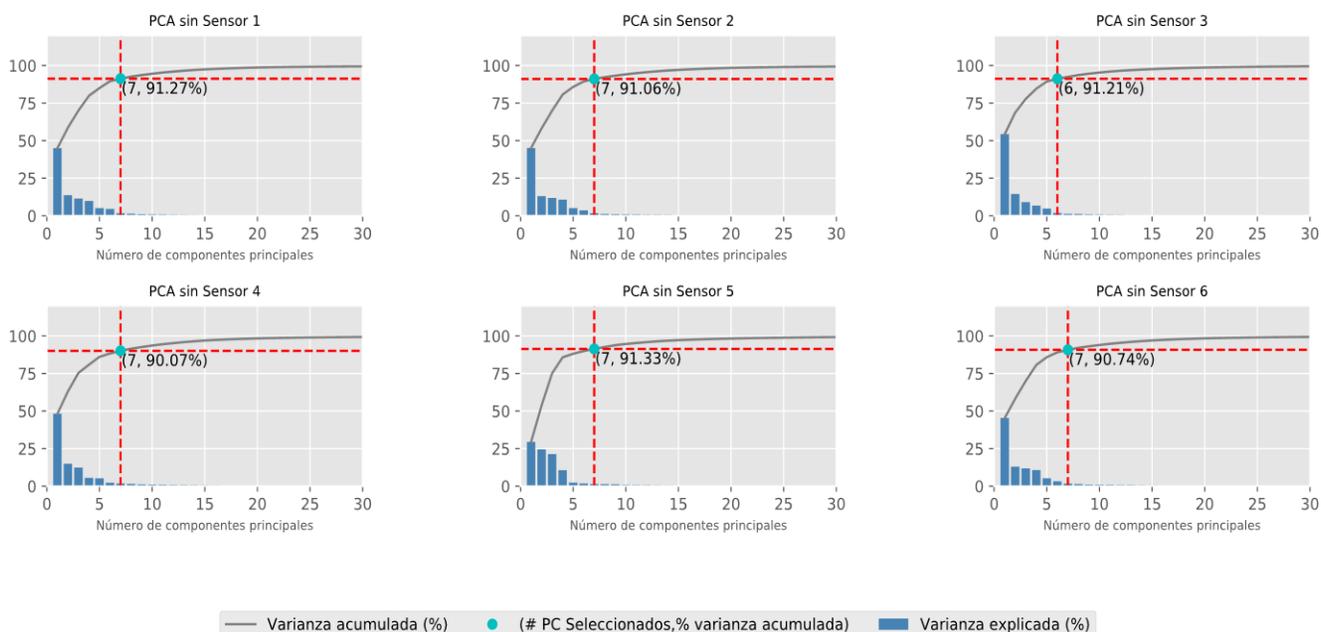
En primer lugar, se escalaron los datos de forma robusta para resaltar los outliers, este escalamiento centra cada una de las características de los datos en la mediana y los escala de acuerdo con el rango intercuantil IQR (diferencia entre el cuantil  $Q_1$  y  $Q_3$ ), su cálculo se muestra en la Ecuación 18.

$$X_{centrada} \text{ } i \text{ característica} = \frac{X_{i \text{ característica}} - \text{mediana}(X_{i \text{ característica}})}{Q_{1i \text{ característica}} - Q_{3i \text{ característica}}} \quad ( 18 )$$

A continuación, se distribuyeron los datos en seis grupos cada uno de ellos sin la señal de uno de los sensores, esto con el fin de analizar si un sensor en específico es el que influye negativamente en los experimentos. Posteriormente, con el algoritmo no supervisado de análisis de componentes principales (PCA) se elaboraron modelos con los seis grupos de datos.

Se determinó el número de componentes principales que se utilizarían en cada uno de los modelos, para ello se elaboraron gráficas de sedimentación, como se muestra en la Figura 18, que muestra el porcentaje de varianza explicado por cada uno de los componentes principales y el porcentaje acumulado. La gráfica permitió identificar el punto en el que el descenso del porcentaje

de varianza y el ascenso del porcentaje acumulado empezó a estabilizarse, llamado también punto codo, este se encontró en el sexto componente principal para el tercer modelo y en el séptimo componente principal en los demás, en conjunto explican más del 90% de varianza de los datos.



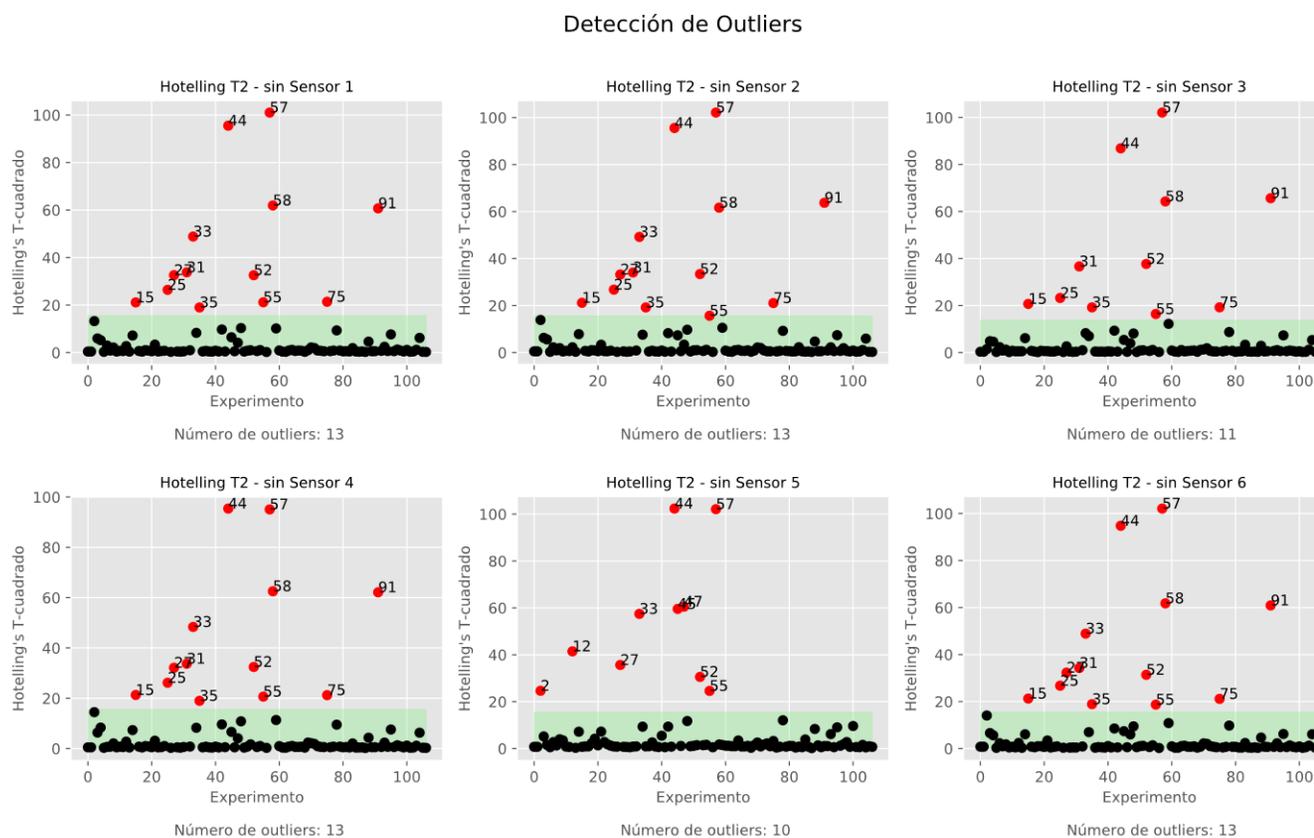
**Figura 18.** Graficas de sedimentación de los componentes de los modelos PCA - Base de datos 1

Después de seleccionar el número de componentes principales para cada uno de los modelos, se definió la métrica utilizada para la detección de outliers, en este caso el valor  $T^2$ -de Hotelling que es la medida de la variación de cada observación dentro del modelo PCA. A continuación, se fijó un nivel de confianza de 95%, con el cual se establece que los experimentos fuera de este serán considerados como outliers.

Con el valor  $T^2$ -de Hotelling de cada uno de los experimentos y el nivel de confianza calculado se identificaron como outliers aquellos experimentos con una variación mayor al límite de confianza, tal como se muestra en la Ecuación 19.

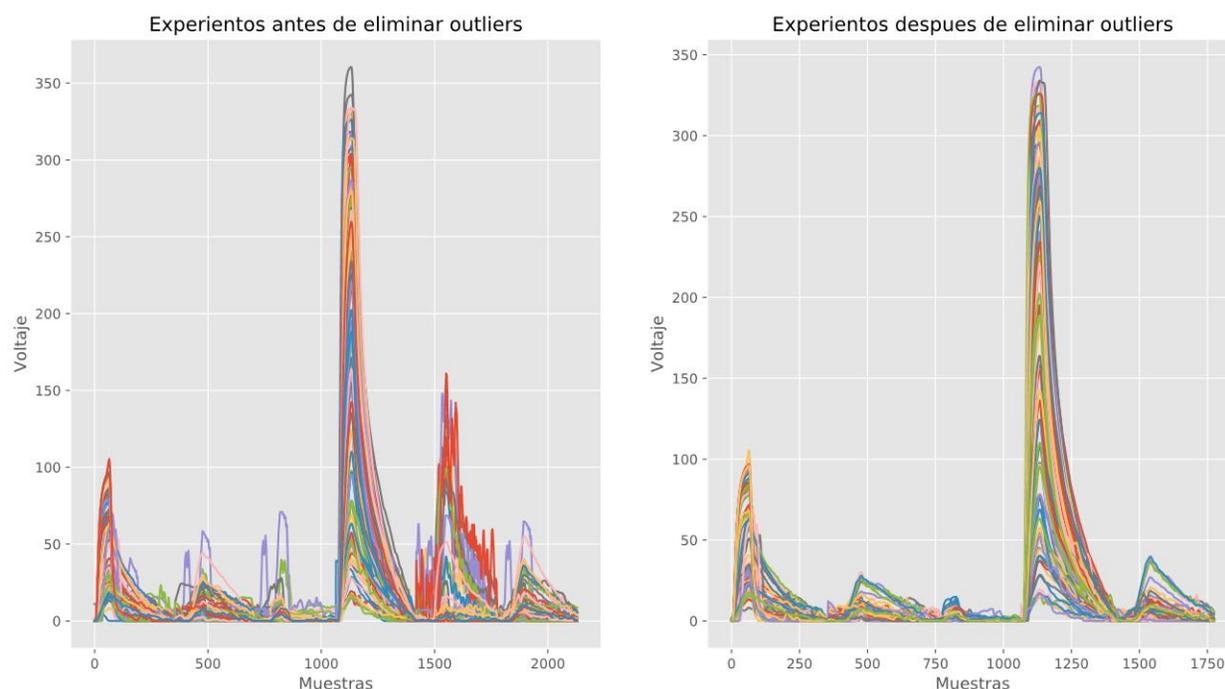
$$T_i^2 > t^2 \rightarrow \text{Outlier} \quad ( 19 )$$

En la Figura 19 se observa el valor  $T^2$ -de Hotelling de cada uno de los experimentos para cada uno de los modelos de la base de datos 1, en el cual se identificó con color rojo las observaciones que se encontraban fuera del límite de confianza y con negro aquellas que no. Como se observa en la figura, los experimentos considerados como outliers son los mismos en todos los modelos menos en el modelo sin el sensor cinco, con lo cual se puede concluir que este sensor es el causante de la mayoría de los experimentos atípicos. Es por ello que en esta etapa se agregaron dos opciones, una permite al usuario eliminar todos los experimentos que en alguno de los seis modelos fueron identificados como outliers y otra eliminar las señales del sensor causante de las señales atípicas y eliminar las señales identificadas como outliers sin este sensor.



**Figura 19.** Detección de outliers en cada uno de los modelos PCA para las observaciones de la base de datos 1

En este caso se seleccionó la opción número dos, es decir, se eliminó la señal correspondiente al sensor cinco y se eliminaron únicamente los experimentos identificados como outliers sin este sensor, como se observa en la Figura 20. En el caso de la base de datos 2, no se encontró a un sensor causante de la mayoría de los experimentos atípicos, por lo cual se decidió seleccionar la opción número 1, es decir, eliminar todos los experimentos que en alguno de los seis modelos fueron identificados como outliers, mas no las señales de un sensor específico.



**Figura 20.** Observaciones previo a la eliminación de outliers y después de esta - Base de datos 1

En la **Tabla 11**, se muestra el número de experimentos tanto del conjunto de entrenamiento como de prueba antes de la detección de outliers, el número de experimentos identificados como outliers y el número de experimentos después de eliminarlos que serían usados tanto para el desarrollo como para la evaluación de los modelos.

**Tabla 11***Numero de experimentos previo y posterior a la eliminación de outliers -Base de datos 1 y 2*

Numero de experimentos	Base de datos 1			Base de datos 2		
	Datos de entrenamiento	Datos de prueba	Total	Datos de entrenamiento	Datos de prueba	Total
<b>Antes de detección de outliers</b>	99	43	142	44	19	63
<b>Outliers</b>	10	6	16	8	2	10
<b>Después de eliminar outliers</b>	89	37	126	36	17	53

En este capítulo se presentó una descripción de la cantidad y concentración de los experimentos utilizados para el desarrollo de los modelos multivariantes, además, se describió la forma en que se distribuyeron los datos en conjuntos de entrenamiento y prueba. Finalmente, se indicó cada una de las etapas que forman parte del preprocesamiento de los datos.

# CAPÍTULO 4

## Análisis Discriminante y Regresión de Mínimos Cuadrados Parciales (PLS-DA y PLS-R)

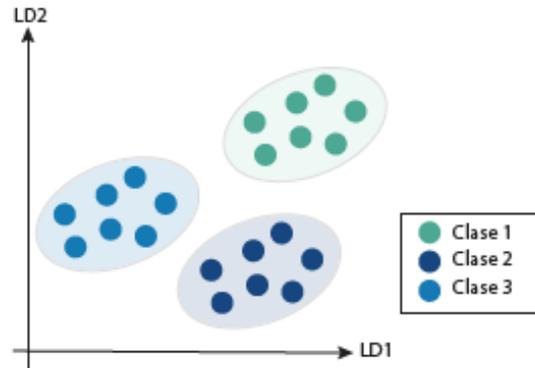
En este capítulo en primer lugar se presenta una descripción de los conceptos bases para el desarrollo del análisis discriminante y la regresión de mínimos cuadrados parciales (PLS-DA) y (PLS-R), respectivamente. En segundo lugar, se presenta el procedimiento para el desarrollo de los modelos, en el que se describe el tipo de tratamiento previo al ingreso de los datos al algoritmo PLS y los hiperparámetros seleccionados en la validación cruzada. Para terminar, se describen los resultados obtenidos con los modelos para las diferentes tareas en base al desempeño de los modelos evaluado con las métricas mencionadas en el Capítulo 2.

### 4.1. Conceptos Básicos

#### 4.1.1. Análisis Discriminante Lineal (LDA)

El análisis discriminante lineal es una técnica supervisada de clasificación, usada para maximizar la separabilidad entre clases (Raschka & Mirjalili, 2017). Durante el proceso de entrenamiento el análisis discriminante aprende los ejes más discriminatorios entre clases,

denominados *LD*, los cuales definen un hiperplano en donde son proyectados los datos de entrada (Geron, 2017), como se observa en la Figura 21.



**Figura 21.** Ejemplo de análisis discriminante lineal

En esta figura se presentan los discriminantes lineales LD1 y LD2 que representan los ejes más discriminativos entre clases. Se muestran tres clases linealmente separables, donde las observaciones representadas con color verde pertenecen a la clase 1, con color naranja a la clase 2 y con color azul a la clase 3

#### 4.1.2. Regresión Lineal

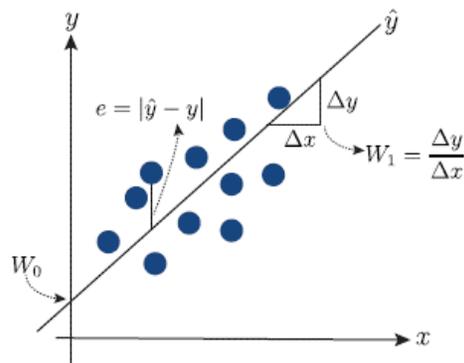
Como se mencionó en el Capítulo 2, la regresión es una tarea de aprendizaje supervisado, en la cual se busca modelar la relación entre una o múltiples variables independientes y una variable dependiente continua. En la regresión lineal, el modelo es una función lineal de los datos de entrada, cuyos pesos y termino de bias (conocido también como intercepción en el eje y del modelo) son los parámetros aprendidos durante el proceso de entrenamiento (Alpaydin, 2012). La ecuación de un modelo de regresión lineal se expresa mediante la suma ponderada de las  $M$  características de los datos de entrada y el termino de bias  $w_0$ , como se observa en la Ecuación 20.

$$\hat{y} = w_0 + w_1x_1 + \dots + w_Mx_M = \sum_{i=0}^M w_i x_i \quad ( 20 )$$

La ecuación presentada es el producto punto entre los valores de los pesos  $W$  y los datos de entrada  $X$ , por lo cual, puede expresarse en forma vectorizada, como se muestra en la Ecuación 21, para un dato de entrada.

$$\hat{y} = [w_0 \quad w_1 \quad \dots \quad w_M]_{(1,M)} \begin{bmatrix} x_0 \\ x_1 \\ \cdot \\ \cdot \\ x_M \end{bmatrix}_{(M,1)} = W_{(M,1)}^T X_{(M,1)} \quad ( 21 )$$

El entrenamiento de un modelo de regresión lineal consiste en obtener valores óptimos de pesos  $W$  que permitan que el modelo  $\hat{y}$  se ajuste de la mejor manera posible al conjunto de datos de entrenamiento, con el fin de que posteriormente el modelo pueda ser usado para predecir las respuestas continuas de nuevos datos de entrada. Existen diversos métodos para optimizar los pesos del modelo de regresión lineal durante el proceso de aprendizaje, uno de los cuales se define en la subsección 4.1.3. Un ejemplo de regresión lineal se presenta en la Figura 22, en la cual el modelo de una sola variable de entrada (univariante) se define como una línea de regresión (definido como un plano en el caso dos variables o como un hiperplano para tres o más variables de entrada), sus elementos se detallan a continuación.



**Figura 22.** Ejemplo de regresión lineal univariante

En esta figura se presentan un ejemplo de regresión lineal univariante, donde  $\hat{y}$  representa el modelo obtenido en el proceso de entrenamiento,  $e$  el error de predicción de cada observación,  $w_1$  el peso de la característica de la entrada  $x$  y  $w_0$  la interceptación en el eje  $y$ .

- $\hat{y}$  : es el modelo que mejor se ajusta a través del conjunto de observaciones.
- $e$ : residuo u offset, representa el error de predicción de cada observación (línea vertical de cada observación a la línea de regresión) y se calcula como la diferencia de la respuesta predicha con el modelo y respuesta la real  $|\hat{y} - y|$ .
- $w_1$ : peso de la característica de la entrada  $x$ , representa la pendiente de la línea de regresión.
- $w_0$ : termino de bias o intercepción en el eje  $y$  del modelo.

#### 4.1.3. Método de Mínimos Cuadrados

El método de mínimos cuadrados consiste en encontrar un modelo que minimice la suma de los residuos cuadrados de la Ecuación 22, con el cual se obtiene el resultado de la suma de las penalizaciones al aplicar el modelo de regresión en los  $N$  datos de entrenamiento.

$$\sum_{i=1}^N (\hat{y}_i - y_i)^2 \quad ( 22 )$$

Donde  $\hat{y}_i$  representa de la salida predicha por el modelo y  $y_i$  la salida real. La ecuación anterior puede expresarse de forma matricial como se observa en la Ecuación 23 y 24:

$$(X_{(N,M)} W_{(M,1)} - Y_{(N,1)})^2 \quad ( 23 )$$

$$(X_{(NXM)} W_{(MX1)} - Y_{(NX1)})^T (X_{(NXM)} W_{(MX1)} - Y_{(NX1)}) \quad ( 24 )$$

Donde,  $X$  representa el conjunto de datos de entrenamiento con  $N$  observaciones y  $M$  características,  $W$  los  $M$  pesos del modelo, y  $Y$  las salidas de las  $N$  observaciones. Para minimizar la suma de los residuos cuadrados de la Ecuación 22, se debe establecer el gradiente de la función en cero, como se observa en la Ecuación 25. Esto significa obtener el valor de los pesos en los extremos

de la función, es decir, el valor de los pesos con los cuales el resultado de la suma de los residuos cuadrados es mínimo.

$$\frac{\delta_{((X_{(N,M)} W_{(M,1)} - Y_{(N,1)})^2)}}{\delta_W} = 0 \quad ( 25 )$$

$$X^T X W = X^T Y \quad ( 26 )$$

Asumiendo que la matriz  $X^T X$  es invertible, se multiplican ambos términos de la Ecuación 26 por  $(X^T X)^{-1}$ , obteniendo la ecuación normal con los valores de pesos óptimos que se muestra en la Ecuación 27.

$$W = (X^T X)^{-1} X^T Y \quad ( 27 )$$

Este método es útil cuando el número de entradas es mayor al de características ( $N > M$ ), ya que al contar con más observaciones que incógnitas no existe una solución única de pesos óptimos, pero esta puede ser estimada al ajustar el modelo al conjunto de observaciones de entrenamiento de forma que la suma del cuadrado de los errores sea mínima. Sin embargo, en casos en los que el número de características es mayor al de entradas ( $M > N$ ) o cuando las características de entrada son linealmente dependientes (cuando alguna de las características es una combinación lineal de las demás) este método no puede ser usado ya que se tendrían más incógnitas a ser estimadas que observaciones, por lo tanto, infinitas soluciones (Geladi & Kowalski, 1986).

#### 4.1.4. Regresión de Mínimos Cuadrados Parciales (PLS-R)

Como se menciona en la subsección anterior los parámetros óptimos para un modelo de regresión lineal se obtienen usualmente con el método de mínimos cuadrados mediante la Ecuación 27, sin embargo  $(X^T X)^{-1}$  puede no existir cuando el número de características es mayor al de

entradas ( $M > N$ ) o cuando las características son linealmente dependientes. PLS-R soluciona este problema al proyectar las características de  $X$  en un subespacio de menor dimensión, formando un nuevo conjunto de variables ortogonales entre sí, denominadas variables latentes, que maximicen la covariancia entre  $X$  y  $Y$ . De forma que estas variables latentes sean usadas para obtener un modelo matemático que permita predecir salidas continuas, que en este caso serán las concentraciones de sustancia explosiva. El proceso para la construcción de un modelo PLS-R consta de dos etapas principales: La construcción de las variables latentes (Componentes PLS) y la construcción del modelo de predicción (Lee, Liong, & Jemain, 2018), como se muestra en la Figura 23, su descripción se presenta en la subsección 4.1.6 y 4.1.7.



**Figura 23.** Proceso para la construcción de un modelo PLS-R

El proceso de construcción de un modelo PLS-R consta de cuatro etapas: el ingreso de los datos de entrada y salida, la construcción de variables latentes, la construcción del modelo de predicción y la obtención de las salidas predichas.

#### 4.1.5. Análisis Discriminante de Mínimos Cuadrados Parciales PLS-DA

El análisis discriminante de mínimos cuadrados parciales (PLS-DA), es un método de clasificación lineal que combina la reducción de dimensionalidad de los datos (construcción de componentes PLS) y el análisis discriminante (construcción del modelo de predicción cualitativo) en un solo algoritmo (Lee et al., 2018). Este método puede definirse como una variante de la regresión lineal PLS utilizada cuando la salida del modelo es categórica, su objetivo es encontrar un modelo matemático que permita predecir la pertenencia de cada observación a una clase u otra.

El proceso para la construcción de un modelo PLS-DA es el mismo del modelo PLS-R y consta de igual manera de dos etapas principales: La construcción de las variables latentes (Componentes PLS) y la construcción del modelo de predicción, como se muestra en la Figura 24, su descripción se presenta en la subsección 4.1.6 y 4.1.7.



**Figura 24.** Proceso para la construcción de un modelo PLS-DA

El proceso de construcción de un modelo PLS-DA consta de cuatro etapas: el ingreso de los datos de entrada y salida, la construcción de variables latentes, la construcción del modelo de predicción y la obtención de las salidas predichas.

#### 4.1.6. Construcción de Componentes PLS

Los componentes PLS, se obtienen mediante la descomposición de  $X$  en  $T$  scores ortogonales,  $P$  loadings y  $E$  error como se observa en la Ecuación 28, y de  $Y$  en  $U$  scores,  $Q$  loadings y  $F$  error, de acuerdo con la Ecuación 29.

$$X_{(N,M)} = T_{(N,A)}P_{(M,A)}^T + E_{(N,M)} \quad ( 28 )$$

$$Y_{(N,P)} = U_{(N,A)}Q_{(P,A)}^T + F_{(N,P)} \quad ( 29 )$$

Donde,  $A$  representa el número de variables latentes, los loadings los vectores dirección unitarios de las variables latentes a lo largo de cada eje y los scores los vectores con las proyecciones de cada observación en los vectores de loadings. Dado que los scores  $T$  y  $U$  son combinaciones lineales de  $X$  y  $Y$  respectivamente, pueden expresarse como el producto de la matriz de entrada por

una matriz de pesos  $W$  en el caso de  $X$ , como se observa en la Ecuación 30 y el producto de la salida real por una matriz de pesos  $C$  en el caso de  $Y$ , como se observa en la Ecuación 31.

$$T_{(N,A)} = X_{(N,M)}W_{(M,A)} \quad ( 30 )$$

$$U_{(N,A)} = Y_{(N,Q)}C_{(Q,A)} \quad ( 31 )$$

La covarianza de los componentes  $T$  y  $U$  debe ser máxima, como se muestra en la Ecuación 32. Además, deben ser ortogonales entre sí mediante la restricción de ortogonalidad del vector  $W$  que se indica en la Ecuación 33.

$$[cov(XW, U)]^2 = [W^T cov(X, U)]^2 = W^T cov(X, U) cov(X, U)^T W \quad ( 32 )$$

$$W^T W = 1 \quad ( 33 )$$

El vector de pesos  $W$  que cumple con las restricciones de la Ecuación 32 y 33, es el vector de covarianzas normalizado de la Ecuación 34.

$$W = \frac{cov(X, U)}{\|cov(X, U)\|} = \frac{X_{(NXM)}^T U_{1(NX1)}}{\|U_{1(NX1)}^T U_{1(NX1)}\|} \quad ( 34 )$$

#### 4.1.7. Construcción del Modelo de Predicción

Con los componentes PLS calculados, se puede obtener el modelo de predicción de la Ecuación 35, el cual modela la relación de las salidas  $Y$  con las variables latentes de  $X$ . El coeficiente de regresión del modelo es representado por la Ecuación 36, con lo cual al remplazarlo en la Ecuación 35, el modelo de predicción se define finalmente como la función de la Ecuación 37.

$$\hat{Y} = UQ^T = TQ^T = XW(P^TW)^{-1}Q^T \quad ( 35 )$$

$$B_{(MXP)} = W_{(MXA)}(P_{MXA}^TW_{MXA})^{-1}Q_{PXA}^T \quad ( 36 )$$

$$\hat{Y} = XB \quad ( 37 )$$

Donde valor de salida predicho  $\hat{Y}$  será un valor entre 0 y 1 para indicar la pertenencia a una clase u otra de sustancia explosiva en el caso de PLS-DA y valores continuos relacionados a la concentración de sustancia explosiva para PLS-R.

## 4.2. Generación de Modelos PLS-DA y PLS-R

El proceso de desarrollo de un modelo PLS-DA o de un modelo PLS-R, como se mencionó en la sección anterior, consta de dos etapas principales: la construcción de los componentes PLS en la etapa de reducción de dimensionalidad y la construcción del modelo de predicción. Todo esto se realizó con la función *PLSRegression* de la librería *sklearn*, la cual utiliza el algoritmo de proceso iterativo no lineal de mínimos cuadrados parciales (NIPALS) para el cálculo de los componentes PLS. Antes de ingresar los datos de entrenamiento al algoritmo *PLSRegression* se realizó un escalamiento de las características de los datos para que sean más fáciles de interpretar por el algoritmo de aprendizaje. Posterior a ello se seleccionaron el número de variables latentes de los modelos mediante la validación cruzada y finalmente se entrenaron los modelos con el número de variables latentes seleccionado.

### 4.2.1. Escalamiento

Un paso previo al entrenamiento de los modelos es la selección del método de escalamiento adecuado para el conjunto de datos de entrenamiento y prueba. El efecto del escalamiento de las

características de los datos depende del tipo de algoritmo aprendizaje utilizado. En el caso del método de mínimos cuadrados parciales, este si depende de la escala de los datos ya que características con mayor magnitud que otras tendrán mayor importancia en el proceso de aprendizaje. Por lo tanto, es importante seleccionar un escalamiento adecuado que permita enfocar el modelo en las características más importantes o en ausencia de conocimiento sobre la importancia relativa de las características dar la misma importancia a todas (Wold, Sjöström, & Eriksson, 2001).

#### ▪ **Modelo PLS-DA**

Para seleccionar el método de escalamiento adecuado se realizaron pruebas con los diferentes tipos de escalamientos proporcionados por Python, con cada uno de ellos se evaluó el desempeño promedio de los modelos generados durante la validación cruzada en función de la métrica AUC definida en el Capítulo 2, como se muestra en la **Tabla 12**. Con estos resultados se concluyó que el desempeño del algoritmo de mínimos cuadrados parciales como establece la teoría depende del tipo de escalamiento utilizado, en este caso el mejor desempeño durante la validación cruzada se obtuvo al centrar los datos en la media y escalarlos con una varianza igual a uno, es decir, mediante el método de autoescalamiento (*StandardScaler* en Python).

**Tabla 12**

*Desempeño del modelo de clasificación de alcohol, pólvora y TNT para la base de datos 1 y 2 en función del tipo de escalamiento utilizado*

<b>Escalamiento</b>	<b>BASE DE DATOS 1</b>		<b>BASE DE DATOS 2</b>	
	<b>#LV</b>	<b>AUC CV</b>	<b>#LV</b>	<b>AUC CV</b>
Sin escalamiento	10	0.735	2	0.607
StandardScaler	4	0.763	4	0.526
MinMaxScaler	5	0.751	7	0.488
MaxAbsScaler	5	0.767	10	0.446
RobustScaler	5	0.747	9	0.415

Para realizar este escalamiento, se debe eliminar la media y escalar las características a la varianza unitaria, como se observa en la Ecuación 38.

$$x_{escalado(i,j)} = \frac{x_{(i,j)} - \bar{x}_j}{\delta(x_j)} \quad ( 38 )$$

Donde,  $x_i$  representa la característica  $i$  de la observación  $j$ . La media  $\bar{x}_j$  y la desviación estándar  $\delta(x_j)$  de cada una de las características  $j$ , obtenidas del conjunto de entrenamiento fueron aplicadas posteriormente en el conjunto de datos de prueba para que los valores de los datos de entrenamiento y de prueba sean comparables entre sí.

- **Modelo PLS-R**

De igual manera para seleccionar el método de escalamiento adecuado para el método PLS-R se realizaron pruebas con los diferentes tipos de escalamientos proporcionados por Python, con cada uno de ellos se evaluó el desempeño promedio de los modelos generados durante la validación cruzada en función de la métrica  $R^2$  definida en el Capítulo 2, como se muestra en la **Tabla 13**. Con estos resultados se concluyó que se obtuvo mejores resultados durante la validación cruzada tanto para la base de datos 1 como la 2 al únicamente centrar los datos en la media, es decir, al no utilizar ningún método de escalamiento. El centrado de los datos se realizó mediante la Ecuación 39.

$$x_{centrado(i,j)} = x_{(i,j)} - \bar{x}_j \quad ( 39 )$$

Donde,  $x_i$  representa la característica  $i$  de la observación  $j$ . La media de cada una de las características  $\bar{x}_j$  obtenida del conjunto de entrenamiento fue aplicada posteriormente en el conjunto de datos de prueba para que los valores de los datos de entrenamiento y de prueba sean comparables entre sí.

**Tabla 13**

*Desempeño del modelo de cuantificación de alcohol, pólvora y TNT para la base de datos 1 y 2 en función del tipo de escalamiento utilizado*

Escalamiento	BASE DE DATOS 1		BASE DE DATOS 2	
	#LV	R2 Test CV	#LV	R2 Test CV
Sin escalamiento	8	1.520	2	4.657
StandardScaler	8	1.528	2	6.155
MinMaxScaler	6	1.652	2	6.264
MaxAbsScaler	6	1.653	2	6.151
RobustScaler	3	1.733	2	5.660

#### 4.2.2. Balanceo de Clases

Al tener un mayor número de observaciones de una clase que de otra, el algoritmo de aprendizaje aprenderá implícitamente un modelo que optimiza las predicciones basadas en la clase más abundante en el conjunto de datos (Raschka & Mirjalili, 2017). Una forma de solucionar este problema es balancear la cantidad de observaciones en cada clase ya sea mediante una técnica de sobre muestro o un submuestreo de los datos. En este caso se utilizó el algoritmo *SMOTE*, el cual sobre muestrea la clase minoritaria creando observaciones sintéticas (Chawla, Bowyer, Hall, & Kegelmeyer, 2002).

#### 4.2.3. Validación Cruzada

Se utilizó la validación cruzada de  $k=10$  iteraciones mencionada en el Capítulo 2 para seleccionar el número de variables latentes tanto para el modelo PLS-R como para el modelo PLS-DA, con la diferencia de que la métrica utilizada para la selección del número de variables latentes del modelo PLS-R se realizó con el valor del error cuadrático medio (MSE) y con el valor del área bajo la curva (AUC) para el modelo PLS-DA. Finalmente, los modelos se entrenaron con el número de variables latentes seleccionado mediante la validación cruzada. Los resultados del desempeño de los modelos se presentan en la sección 4.3.

## 4.3. Resultados del Desempeño de los Modelos PLS-DA y PLS-R

### 4.3.1. Modelo PLS-DA

En la **Tabla 14** y

**Tabla 15**, se presentan los resultados obtenidos con los modelos de clasificación para la base de datos 1 y 2, en los cuales se describe el número de variables latentes (LV) utilizado en cada uno de los modelos y el AUC de los datos de entrenamiento y prueba para la evaluación de su desempeño. Estos resultados fueron graficados para el análisis que se realizó a continuación.

**Tabla 14**

*Base de datos 1- Resultados de los modelos PLS-DA*

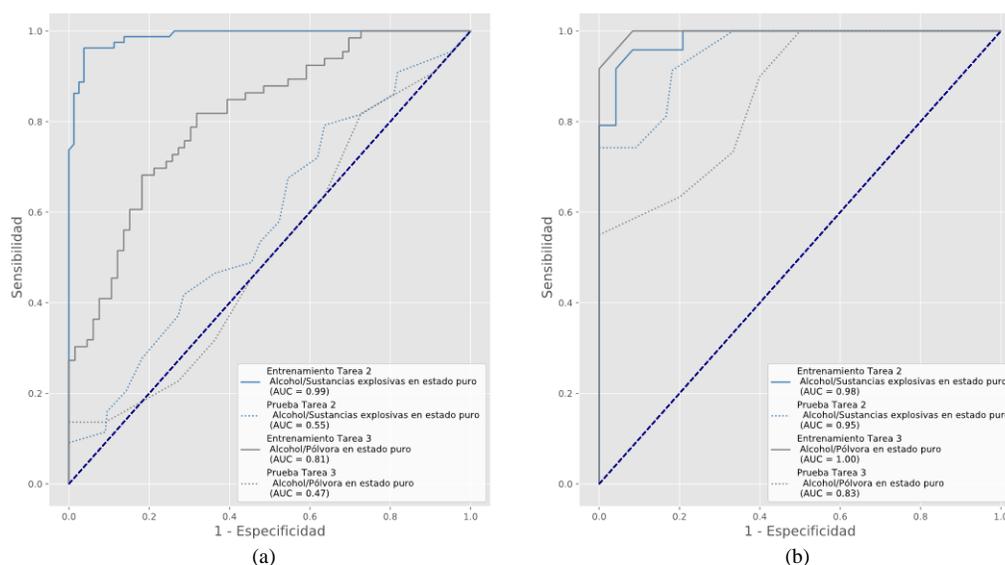
N°	TAREA	#LV	AUC	
			TRAIN	TEST
1	Alcohol/Sustancias explosivas en estado puro y mezclas	5	0.88	0.51
2	Alcohol/Sustancias explosivas en estado puro	9	0.99	0.55
3	Alcohol/Pólvora en estado puro	2	0.81	0.47
4	Alcohol/TNT en estado puro	2	0.84	0.56
5	Pólvora en estado puro/TNT en estado puro	2	0.88	0.35
6	Alcohol/Pólvora en estado puro y mezcla de pólvora	2	0.72	0.29
7	Alcohol/TNT en estado puro y mezcla de TNT	8	0.99	0.52
8	Alcohol/Pólvora en estado puro/TNT en estado puro	5	0.87	0.51

**Tabla 15**

*Base de datos 2- Resultados de los modelos PLS-DA*

N°	TAREA	#LV	AUC	
			TRAIN	TEST
2	Alcohol/Sustancias explosivas en estado puro	6	0.98	0.95
3	Alcohol/Pólvora en estado puro	10	1.00	0.83
4	Alcohol/TNT en estado puro	5	1.00	0.79
5	Pólvora en estado puro/TNT en estado puro	2	0.88	0.40
8	Alcohol/Pólvora en estado puro/TNT en estado puro	8	0.99	0.65

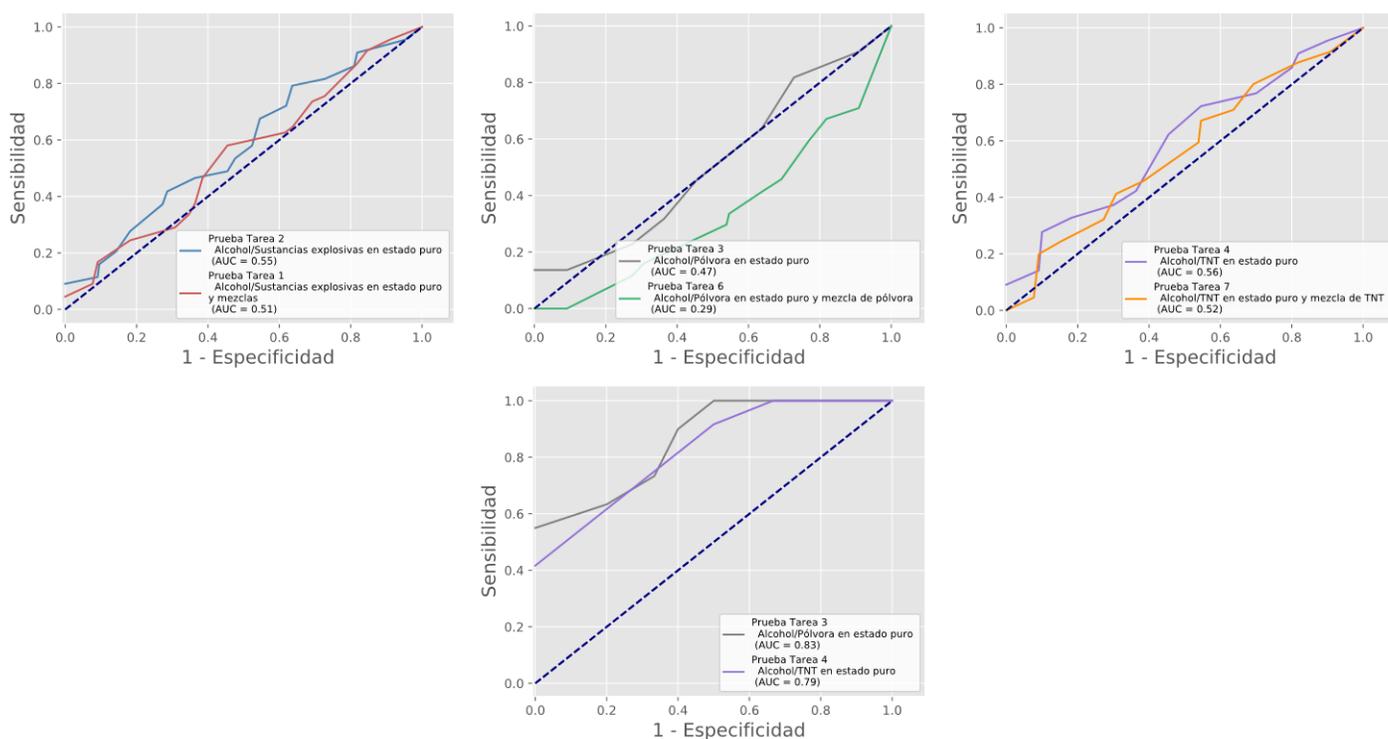
En la Figura 25a, se presentan las curvas ROC y el AUC de dos de los modelos de la base de datos 1, en los que se puede apreciar que el desempeño de los modelos en los datos de entrenamiento es mucho mejor que el de los datos de prueba, ya que la curva ROC se encuentra mucho más cerca de la esquina superior izquierda ( $AUC \gg 0.5$ ) para los datos de entrenamiento representada por líneas continuas, y muy cerca de la región de indecisión ( $AUC \approx 0.5$ ) para los datos de prueba representada por líneas discontinuas. En el caso de la base de datos 2, cómo se observa en la Figura 25b, la diferencia del desempeño entre los datos de entrenamiento y prueba no es significativa, a excepción de los resultados obtenidos en la tarea 5 y 8 cuyo análisis se presenta posteriormente.



**Figura 25.** Curvas ROC para la descripción del desempeño de los modelos en datos de entrenamiento y prueba

En cuanto al desempeño de los modelos de la base de datos 1 para clasificación de sustancias explosivas puras y con mezclas de jabón o pasta dental, este es inferior al de clasificación de únicamente sustancias explosivas puras, como se puede observar en los resultados del AUC y gráficamente en las curvas ROC de la Figura 26. Además, como se muestra en la Figura 26b y Figura 26c, el desempeño de los modelos de clasificación de TNT es superior a los de clasificación

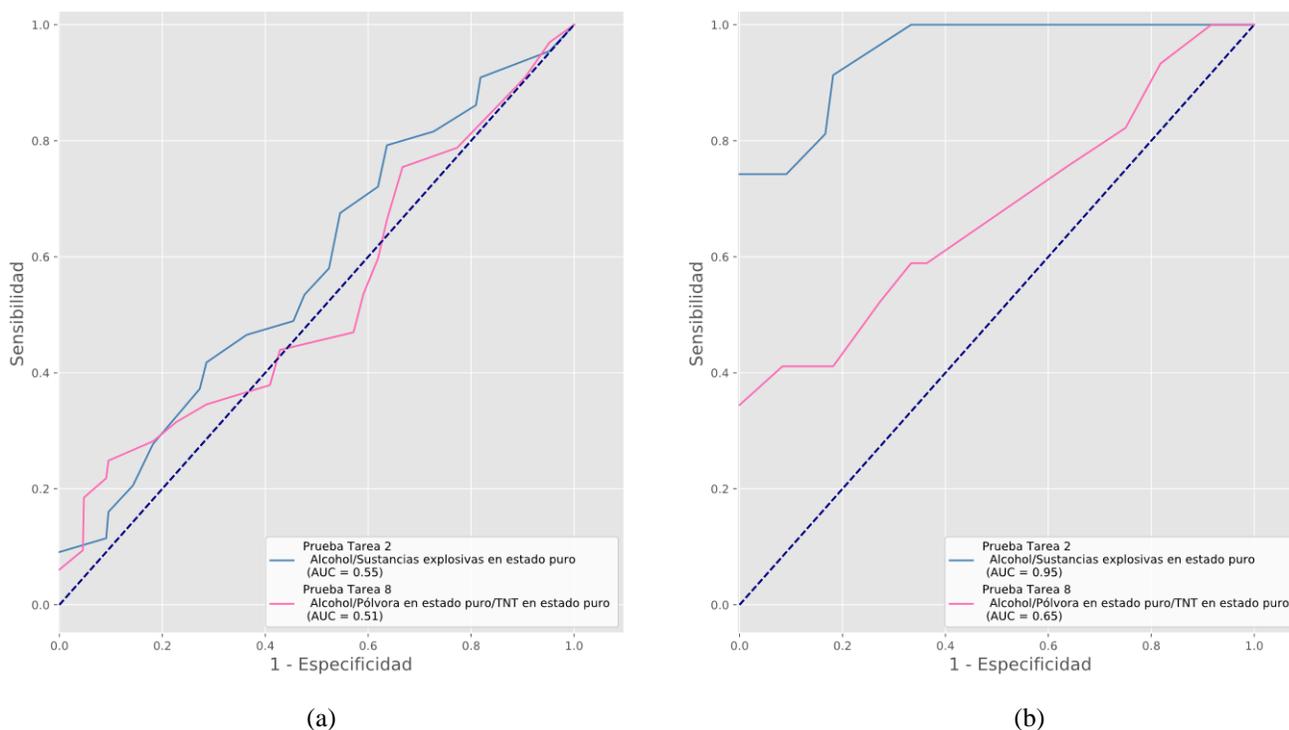
de pólvora para la base de datos 1 (AUC:  $0.56 > 0.47$ ), y como se muestra en la Figura 26d para la base de datos 2, el desempeño del modelo de clasificación de pólvora es superior al de TNT (AUC:  $0.83 > 0.79$ ). Por consiguiente, se espera que bajo las condiciones actuales del prototipo y con concentraciones entre 3 y 5gr de sustancias explosivas este sea capaz de clasificar correctamente en mayor porcentaje observaciones de pólvora que de TNT con los modelos PLS-DA.



**Figura 26.** Curvas ROC para descripción del desempeño de los modelos para clasificación de sustancias explosivas puras y mezclas

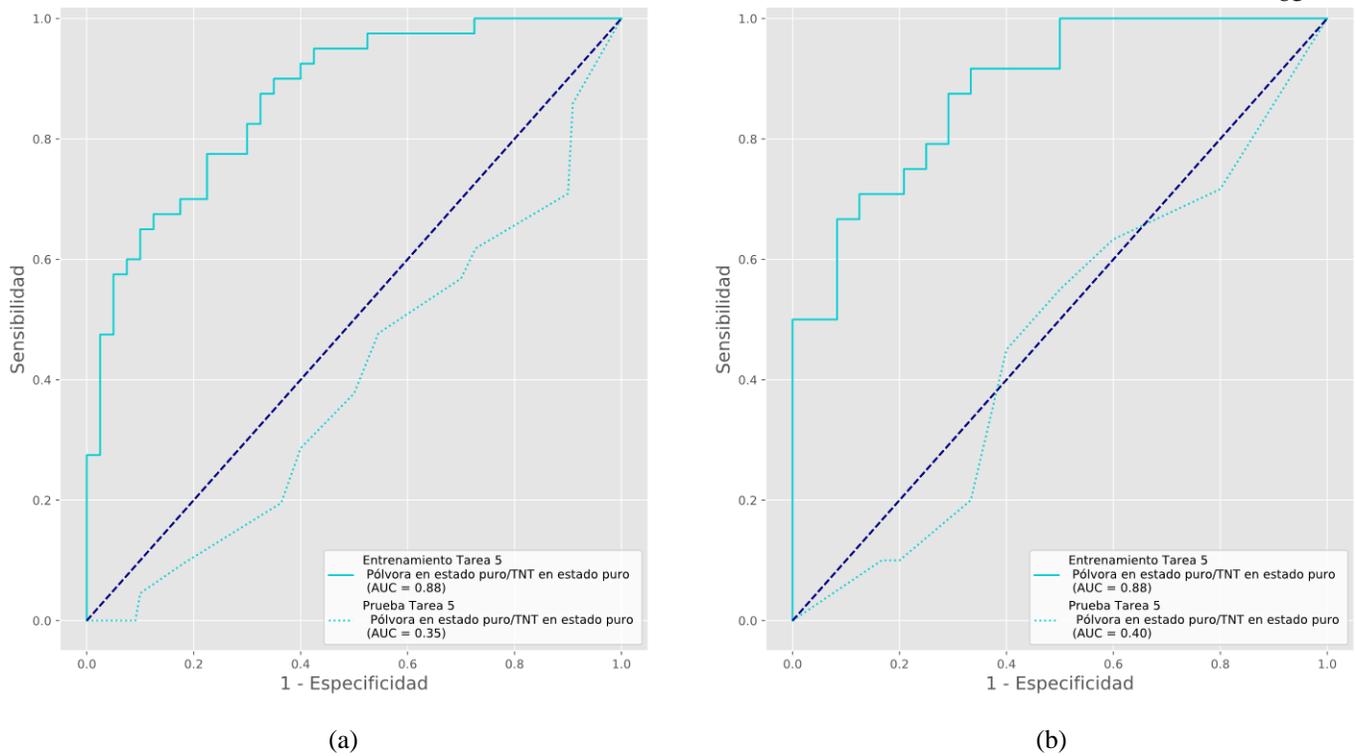
Respecto a la diferencia entre un modelo que determina si una observación es sustancia explosiva o no y uno que determina si una observación no es sustancia explosiva y si lo es el tipo entre TNT y pólvora. Se pudo observar, como se muestra en la Figura 27a para la base de datos 1 y en la Figura 27b para la base de datos 2, que el desempeño de un modelo que únicamente determina si una sustancia es o no explosiva es superior al modelo que determina además el tipo de sustancia explosiva que es. Además, al comparar el desempeño entre los modelos realizados con bajas

concentraciones de sustancia explosiva (experimentos de la base de datos 1 con concentraciones entre 0.1 y 3gr de pólvora y TNT) y los realizados con mayores concentraciones (experimentos de la base de datos 2 con concentraciones entre 3 y 5gr de pólvora y TNT), se pudo observar que se obtuvo mejores resultados al aumentar el nivel de concentración de sustancia explosiva.



**Figura 27.** Curvas ROC para descripción del desempeño de los modelos de clasificación entre dos y tres clases de sustancias

Finalmente, al analizar los resultados de los modelos de clasificación entre TNT y pólvora de la Figura 28a para la base de datos 1 y Figura 28b para la base de datos 2, se determinó que el prototipo e-nose no es capaz de clasificar correctamente entre estas sustancias, a pesar de haber aumentado la concentración de sustancia explosiva en el modelo de la base de datos 2 ya que  $AUC < 0.5$ .



**Figura 28.** Curvas ROC para descripción del desempeño de los modelos de clasificación entre pólvora y TNT

#### 4.3.2. Modelo PLS-R

En la

**Tabla 16** y **Tabla 17**, se presentan los resultados de los modelos de regresión para la base de datos 1 y 2, respectivamente, en los cuales se describe el número de variables latentes (LV) utilizado en cada uno de los modelos, y el MSE y R2 de los datos de entrenamiento y prueba para la evaluación de su desempeño. Estas métricas en conjunto con las gráficas de los valores reales vs los valores predichos por los modelos fueron utilizadas para el análisis que se realizó a continuación.

**Tabla 16***Base de datos 1- Resultados de modelos PLS-R*

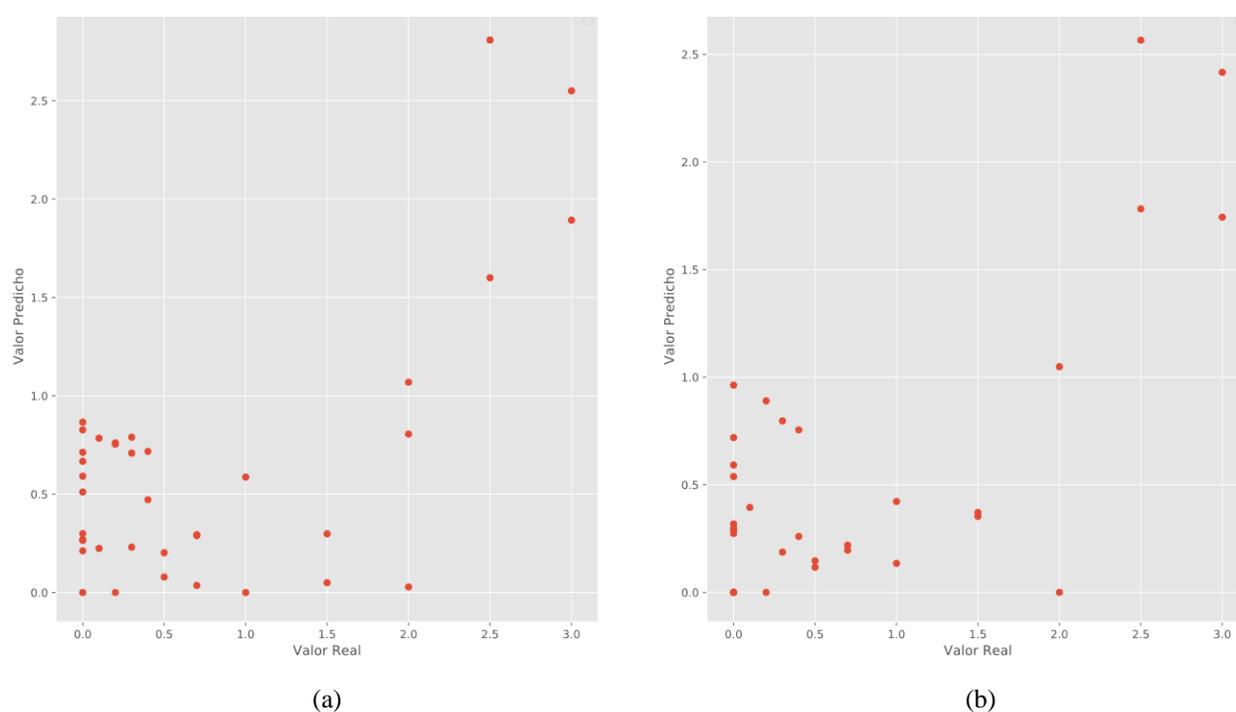
N°	TAREA	#LV	MSE		R2	
			TRAIN	TEST	TRAIN	TEST
1	Alcohol/Sustancias explosivas en estado puro y mezclas	6	0.34	0.51	0.51	0.38
2	Alcohol/Sustancias explosivas en estado puro	5	0.38	0.47	0.52	0.46
3	Alcohol/Pólvora en estado puro	4	0.33	0.52	0.42	0.31
4	Alcohol/TNT en estado puro	7	0.20	0.36	0.69	0.53
5	Pólvora en estado puro/TNT en estado puro	2	0.44	0.42	0.51	0.52
6	Alcohol/Pólvora en estado puro y mezcla de pólvora	2	0.45	0.61	0.14	0.12
7	Alcohol/TNT en estado puro y mezcla de TNT	6	0.30	0.46	0.52	0.40
8	Alcohol/Pólvora en estado puro/TNT en estado puro	5	0.36	0.46	0.24	0.20

**Tabla 17***Base de datos 2- Resultados de modelos PLS-R*

N°	TAREA	#LV	MSE		R2	
			TRAIN	TEST	TRAIN	TEST
2	Alcohol/Sustancias explosivas en estado puro	2	2.85	3.40	0.26	0.17
3	Alcohol/Pólvora en estado puro	2	3.02	3.89	0.26	0.09
4	Alcohol/TNT en estado puro	2	3.04	3.35	0.29	0.22
5	Pólvora en estado puro/TNT en estado puro	2	0.47	0.98	0.24	-0.34
8	Alcohol/Pólvora en estado puro/TNT en estado puro	2	3.34	3.40	0.09	0.08

En base a los resultados obtenidos en los datos de entrenamiento y prueba, se determinó que los modelos de la base de datos 1 y 2, no se sobreajustan en gran medida a los datos de entrenamiento ya que el resultado de R2 y MSE para el conjunto de entrenamiento es similar al obtenido en el conjunto de prueba. En la Figura 29a, se presenta la gráfica de los valores reales vs los valores predichos por el modelo de regresión de sustancias explosivas puras y mezclas, y en la Figura 29b del modelo de regresión de únicamente sustancias explosivas puras, los dos para la base de datos 1.

Estos modelos se encargan de predecir la concentración de sustancia explosiva sin importar el tipo de sustancia explosiva que es, por lo tanto, proporcionan una única salida (concentración predicha por el modelo) por cada observación. El modelo con el que se obtuvo mejores resultados fue aquel encargado de predecir la concentración de sustancias explosivas puras, esto puede verificarse en las gráficas ya que los valores predichos se acercan más a los reales, sobre todo con las concentraciones más altas de sustancia.

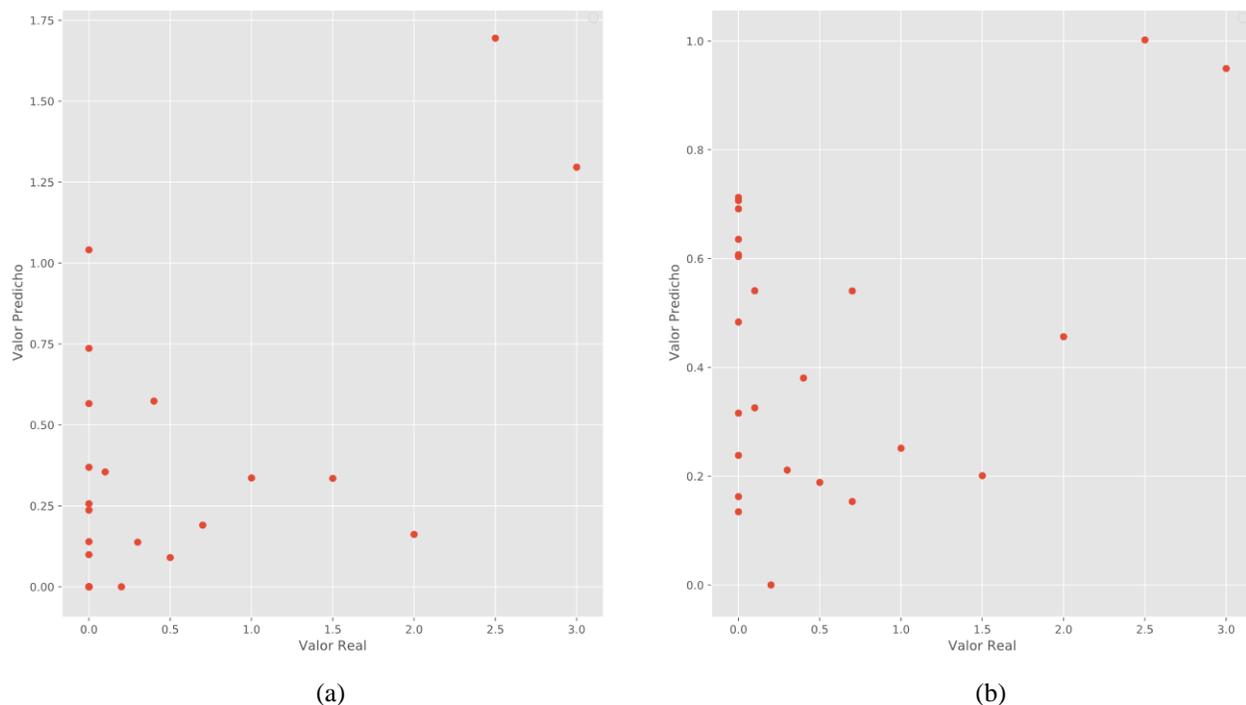


**Figura 29.** Valores predichos por los modelos de regresión

(a) Regresión de sustancias explosivas puras y mezclas (b) Regresión de sustancias explosivas puras

En la Figura 30a, se presenta la gráfica de los valores reales vs los valores predichos por el modelo de regresión de pólvora en estado puro, y en la Figura 30b del modelo de regresión de pólvora en estado puro y mezclas, los dos para la base de datos 1. El modelo con el que se obtuvo mejores resultados al igual que con los modelos analizados en la parte superior, fue aquel encargado

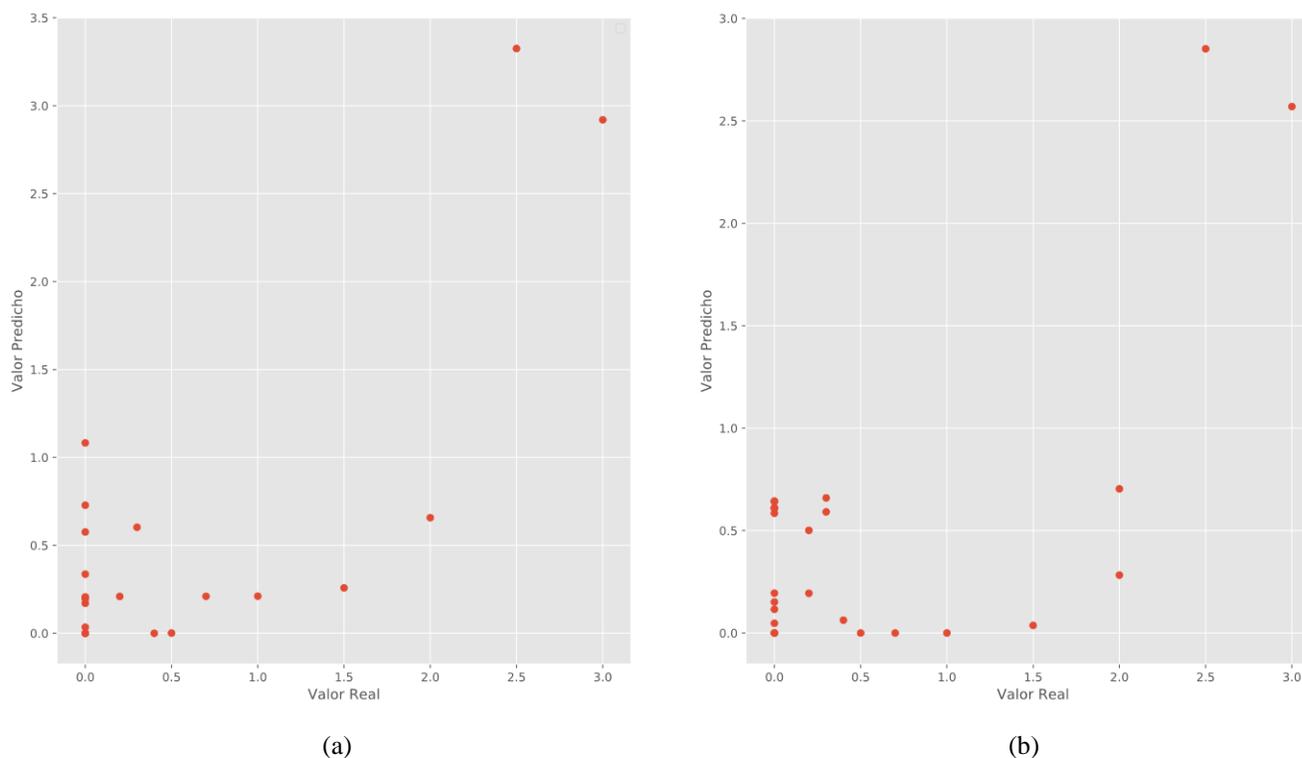
de predecir la concentración de pólvora en estado puro, sin embargo, el desempeño no es bueno, ya que como se observa en las gráficas los valores predichos se alejan mucho de los reales.



**Figura 30.** Valores predichos por los modelos de regresión

(a) Regresión de pólvora en estado puro (b) Regresión de pólvora en estado puro y mezclas

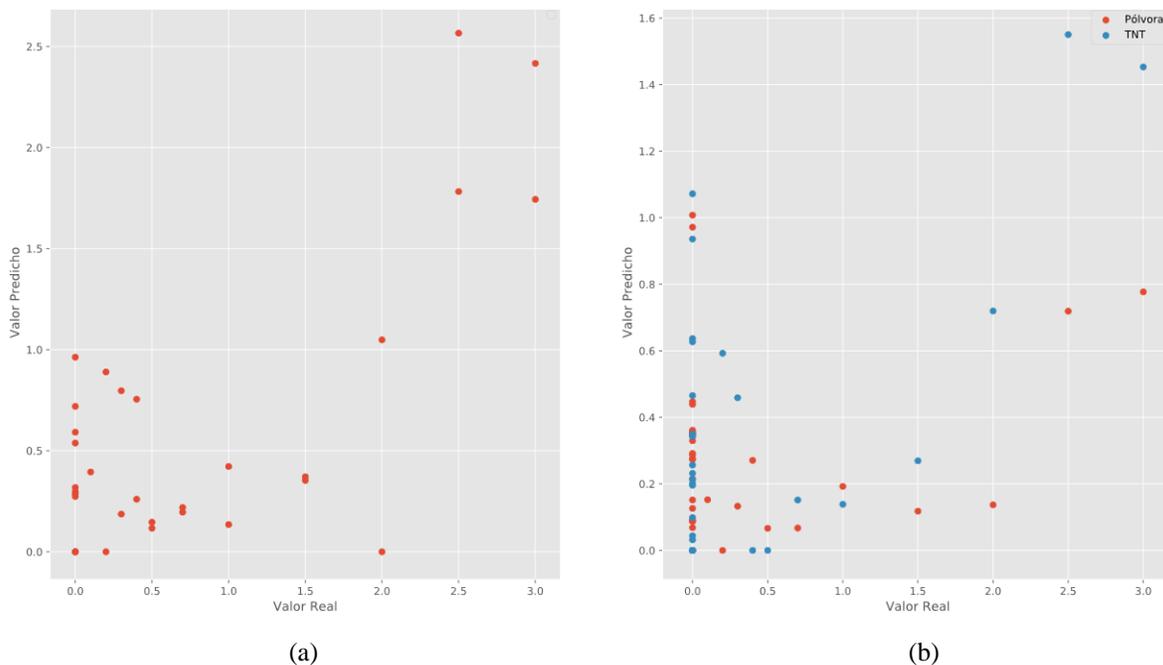
En la Figura 31a, se presenta la gráfica de los valores reales vs los valores predichos por el modelo de regresión de TNT en estado puro, y en la Figura 31b del modelo de regresión de TNT en estado puro y mezclas, los dos para la base de datos 1. El modelo con el que se obtuvo mejores resultados al igual que con los modelos analizados anteriormente, fue aquel encargado de predecir la concentración de TNT en estado puro. A diferencia del encargado de predecir la concentración de pólvora, el desempeño de este modelo es bueno, ya que como se aprecia en las gráficas para observaciones sin sustancia explosiva el valor predicho es muy cercano a cero y para concentraciones altas cercano al real.



**Figura 31.** Valores predichos por los modelos de regresión

(a) Regresión de TNT en estado puro (b) Regresión de TNT en estado puro y mezclas

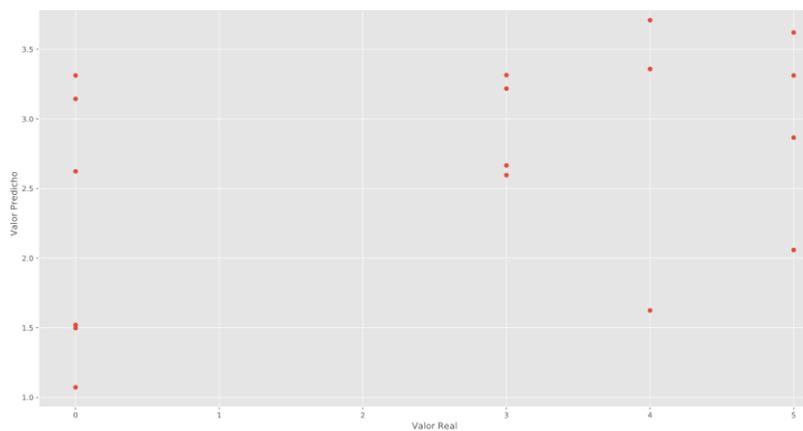
En la Figura 32a, se presenta la gráfica de los valores reales vs los valores predichos por el modelo de regresión de sustancias explosivas puras, y en la Figura 32b la gráfica del modelo de regresión de TNT y pólvora en estado puro, que a diferencia del primero este si cuenta con dos salidas: uno para la concentración predicha de la clase 1 correspondiente a pólvora y otro para la clase 2 correspondiente a TNT, ambos para la base de datos1. De estos resultados se determinó que el prototipo bajo las condiciones iniciales en que se encontraba y con bajas concentraciones de sustancia explosiva, era capaz de predecir mejor la concentración de las sustancias con un modelo en el que no se toma en cuenta si la observación que ingresa al modelo pertenece a una clase de explosivo u otra que con un modelo de dos clases que predice la concentración por cada clase de sustancia.



**Figura 32.** Valores predichos por los modelos de regresión

(a) Regresión de sustancias explosivas puras (b) Regresión de TNT y pólvora en estado puro

Con relación a los modelos de la base de datos 2, el prototipo tuvo un mal desempeño en todos los casos, como se observa en la Figura 33 de los valores reales vs los valores predichos por el modelo de regresión de sustancias explosivas puras. Las observaciones sin sustancia explosiva son incorrectamente identificadas con valores entre 1 y 3.5 gr y en el caso de las sustancias con el más alto nivel de concentración: 5gr, los resultados predichos varían mucho y se alejan del real.



**Figura 33.** Valores predichos por el modelo de regresión de sustancias explosivas puras para la base de datos 2

En este capítulo se describieron los conceptos utilizados para el desarrollo de los modelos de clasificación PLS-DA y de regresión PLS-R como lo es el método de mínimos cuadrados parciales para la reducción de dimensionalidad de los datos y la construcción de los modelos de predicción. A continuación, se describió el proceso para la creación de los modelos en el que se definió el tipo de escalamiento usado y cuál fue el uso que se dio a la técnica de validación cruzada. Finalmente se presentó una descripción de los resultados obtenidos para cada una de las tareas.

# CAPÍTULO 5

## Regresión Logística

En este capítulo se presenta una descripción de los conceptos utilizados para el desarrollo del modelo de regresión logística. A continuación, se presenta el procedimiento para el desarrollo del modelo, en el que se describe el tipo de tratamiento previo al ingreso de los datos al algoritmo de aprendizaje y los hiperparámetros seleccionados en la validación cruzada. Para terminar, se describen los resultados obtenidos con el modelo para las diferentes tareas en base al desempeño de los modelos evaluados con las métricas mencionadas en el Capítulo 2.

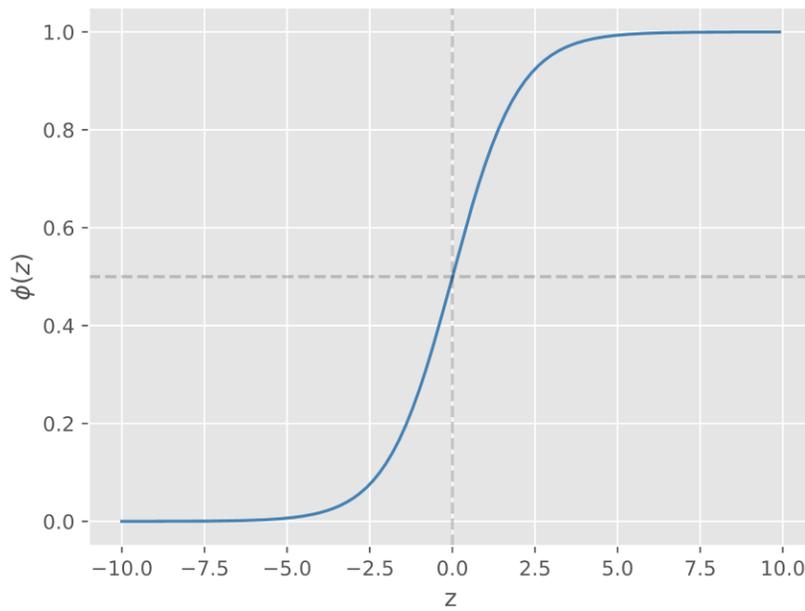
### 5.1. Conceptos Básicos

#### 5.1.1. Función Logística Sigmoide

La función logística sigmoide, es una curva en forma de S que puede tomar cualquier valor real y asignarlo a un valor entre 0 y 1 (Jason Brownlee, 2016). La ecuación que define a esta función se presenta a continuación:

$$\phi(z) = \frac{1}{1 + e^{-z}} \quad ( 40 )$$

Al representar gráficamente la función logística sigmoide, como se observa en la Figura 34, se puede apreciar que la función transforma los valores reales de entrada  $z$  a valores entre 0 y 1 (para valores de entrada grandes  $\phi(z)$  se acercará a 1, caso contrario a 0), su intersección se encuentra en  $\phi = 0.5$ .



**Figura 34.** Función logística sigmoide

La función transforma los valores reales de entrada  $z$  a valores entre 0 y 1, su intersección se encuentra en  $\phi = 0.5$ .

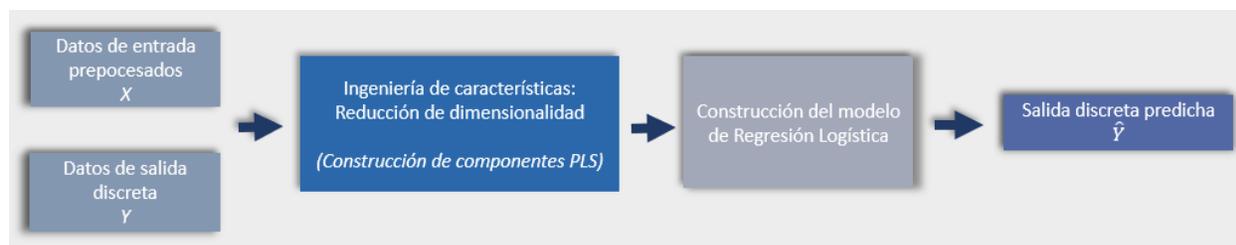
### 5.1.2. Regularización L2

La regularización L2, es un método que permite evitar el sobreajuste de los modelos a los datos de entrenamiento. Consiste en introducir un término a la función de costo para penalizar valores extremos de los parámetros optimizados durante el proceso de entrenamiento, que en este caso son los pesos de los modelos (Raschka & Mirjalili, 2017). La regularización L2, se define mediante la Ecuación 41, en la cual  $\lambda$  representa el parámetro de regularización y  $w_j$  el peso de la característica  $j$ . Mediante la manipulación de  $\lambda$  se controla la fuerza de regularización, mientras mayor sea el valor de  $\lambda$  mayor será la penalización de pesos de gran magnitud.

$$\frac{\lambda}{2} \|w\|^2 = \frac{\lambda}{2} \sum_{j=1}^m w_j^2 \quad ( 41 )$$

### 5.1.3. Regresión Logística

La regresión logística es un método de clasificación lineal que analiza la relación entre múltiples variables independientes y una o múltiples variables dependientes categóricas. Su objetivo es predecir la probabilidad de ocurrencia de un evento categórico mediante el ajuste de los datos a la función logística sigmoide (Park, 2013), definida en la subsección 5.1.1. Con este método el modelo resultante será capaz de predecir la pertenencia de las observaciones a las clases existentes y evitará sobre ajustarse a los datos de entrenamiento mediante el uso del parámetro de regularización L2, definido en la subsección 5.1.2. Para la construcción de los modelos de regresión logística se propuso dos etapas principales: La reducción de la dimensionalidad de los datos de entrada mediante la construcción de componentes PLS y la construcción del modelo de predicción, como se muestra en la Figura 35, estructura propuesta en (Fort & Lambert-Lacroix, 2005). Su descripción se presenta en las subsecciones siguientes.



**Figura 35.** Proceso para la construcción del modelo de regresión logística

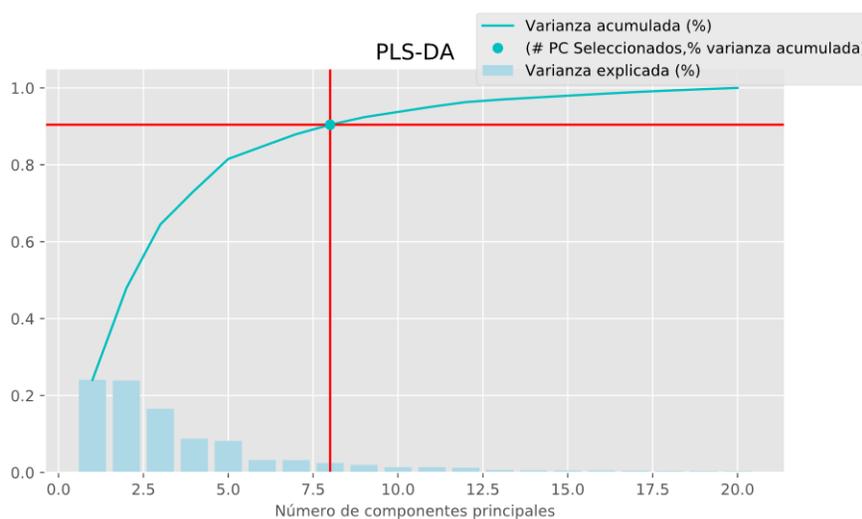
El proceso de construcción del modelo de Regresión Logística consta de cuatro etapas: el ingreso de los datos de entrada y salida, la construcción de variables latentes, la construcción del modelo de predicción y la obtención de las salidas predichas.

### 5.1.4. Construcción de Componentes PLS

Antes de ingresar los datos de entrenamiento al algoritmo de regresión logística, se decidió utilizar el método de mínimos cuadrados parciales (PLS) definido en el Capítulo 4 para reducir la

dimensionalidad de los datos de forma supervisada. Una de las ventajas de utilizar este tipo de reducción de dimensionalidad es que al ser supervisada considera los datos de entrada y sus respectivas salidas para intentar maximizar la separabilidad de las clases en un espacio lineal de características (Raschka & Mirjalili, 2017) y al reducir la dimensionalidad de los datos eliminar características redundantes que no aporten con información importante al algoritmo de aprendizaje.

La selección del número de variables latentes que ingresarán al algoritmo de Regresión Logística se realizó mediante la elaboración de gráficas de sedimentación, como se muestra en la Figura 36, la cual muestra el porcentaje de varianza explicado por cada una de las variables latentes y el porcentaje acumulado. La gráfica permitió identificar el punto en el que el descenso del porcentaje de varianza y el ascenso del porcentaje acumulado empezó a estabilizarse, este se encontró en la octava variable latente para el modelo de discriminación entre alcohol, pólvora y TNT de la base de datos 2, las ocho variables latentes explican alrededor del 90% de varianza de los datos.



**Figura 36.** Grafica de sedimentación para selección de variables latentes

La grafica muestra mediante barras la varianza explicada por cada una de las variables latentes, mediante una línea azul la varianza acumulada en cada una y el número de variables latentes seleccionados con un punto.

### 5.1.5. Construcción de Modelo de Regresión Logística

Como se observa en la Figura 37 la construcción de un modelo de regresión logística implica cinco etapas. En la primera, se realiza la suma ponderada de las  $m$  características de cada una de las observaciones pertenecientes al conjunto de entrenamiento.



**Figura 37.** Pasos para construcción de modelo de Regresión Logística

La suma ponderada ingresa a la función de activación que en este caso es la función logística sigmoide de la Ecuación 42, la cual transforma la suma ponderada de las características de entrada en valores entre 0 y 1.

$$\phi(z) = \frac{1}{1 + e^{-z}} = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + \dots + w_mx_m)}} \quad (42)$$

Para seleccionar los parámetros óptimos del modelo ( $w_0, w_1, \dots, w_m$ ), se debe definir una función objetivo que debe ser optimizada durante el proceso de aprendizaje, esta es la función de costo denominada pérdida logarítmica o log loss y se calcula mediante la Ecuación 43.

$$J(w) = \frac{1}{n} \sum_{k=1}^n -y_k \log(\phi(z_k)) - (1 - y_k) \log(1 - \phi(z_k)) \quad (43)$$

Donde,  $n$  representa el número de datos de entrenamiento,  $\phi(z_k)$  la función de activación para la observación  $k$  y  $y_k$  la salida respectiva de esa observación.

El resultado de la función de costo debe ser minimizado con ayuda de una función de optimización, el algoritmo seleccionado como función de optimización es el algoritmo LIBLINEAR que minimiza la función objetivo y a su vez utiliza la regularización L2 para evitar sobre ajustar los modelos a los datos de entrenamiento. Finalmente, con los pesos óptimos seleccionados durante la etapa de entrenamiento, la probabilidad de ocurrencia predicha (valores entre 0 y 1 para predecir probabilidades) después de que la suma ponderada de las características de los datos de entrada pase por la función de activación es utilizada para evaluar el modelo con diferentes umbrales mediante la métrica AUC definida en el Capítulo 2.

## 5.2. Generación de Modelos de Regresión Logística

El proceso de desarrollo de un modelo de regresión logística, como se mencionó en la sección anterior, consta de dos etapas principales: la reducción de la dimensionalidad de los datos mediante la construcción de variables latentes y la construcción del modelo de predicción. Las variables latentes se obtuvieron con ayuda de la función *PLSRegression* de la librería *sklearn*, la cual utiliza el algoritmo de proceso iterativo no lineal de mínimos cuadrados parciales (NIPALS) para el cálculo de los componentes PLS. Antes de ingresar los datos de entrenamiento al algoritmo *PLSRegression* se realizó un escalamiento de las características de los datos para que sean más fáciles de interpretar por el algoritmo. A continuación, se seleccionaron el número de variables latentes que ingresarían al algoritmo de regresión logística, mediante el proceso indicado en la subsección 5.1.4. Finalmente, se utilizó la función *LogisticRegression* de la librería *sklearn* para entrenar los modelos de clasificación con los parámetros de regularización inverso seleccionados durante la validación cruzada.

### **5.2.1. Escalamiento**

Debido a que el desempeño del método de mínimos cuadrados parciales depende de la escala de los datos se seleccionó el mismo método de escalamiento empleado en el Capítulo 5. El cual centra las características de los datos en la media y los escala con una varianza igual a 1.

### **5.2.2. Balanceo de Clases**

Al tener un mayor número de observaciones de una clase que de otra, el algoritmo de aprendizaje aprenderá implícitamente un modelo que optimiza las predicciones basadas en la clase más abundante en el conjunto de datos (Raschka & Mirjalili, 2017). Una forma de solucionar este problema es balancear la cantidad de observaciones en cada clase ya sea mediante una técnica de sobre muestro o un submuestreo de los datos. En este caso se utilizó el algoritmo *SMOTE*, el cual sobre muestrea la clase minoritaria creando observaciones sintéticas (Chawla et al., 2002).

### **5.2.3. Validación Cruzada**

Se utilizó la validación cruzada de  $k=10$  iteraciones para seleccionar el parámetro de regularización L2 inverso de los modelos de regresión logística, el cual definirá la fuerza de regularización de pesos de gran magnitud. Con el parámetro de regularización inverso seleccionado para cada modelo se realizó el entrenamiento de los modelos. Los resultados de su desempeño se presentan en la siguiente sección.

### 5.3. Resultados del Desempeño de los Modelos de Regresión

#### Logística

En la Tabla 18 y **Tabla 19**, se presentan los resultados obtenidos con los modelos de clasificación para la base de datos 1 y 2, en los cuales se describe el número de variables latentes (LV) utilizado en cada uno de los modelos, el parámetro inverso de regularización (C) seleccionado durante la validación cruzada y el AUC de los datos de entrenamiento y prueba para la evaluación desempeño de los modelos. Estos resultados fueron graficados para el análisis que se realizó a continuación.

**Tabla 18**

*Base de datos 1- Resultados de modelos de regresión logística*

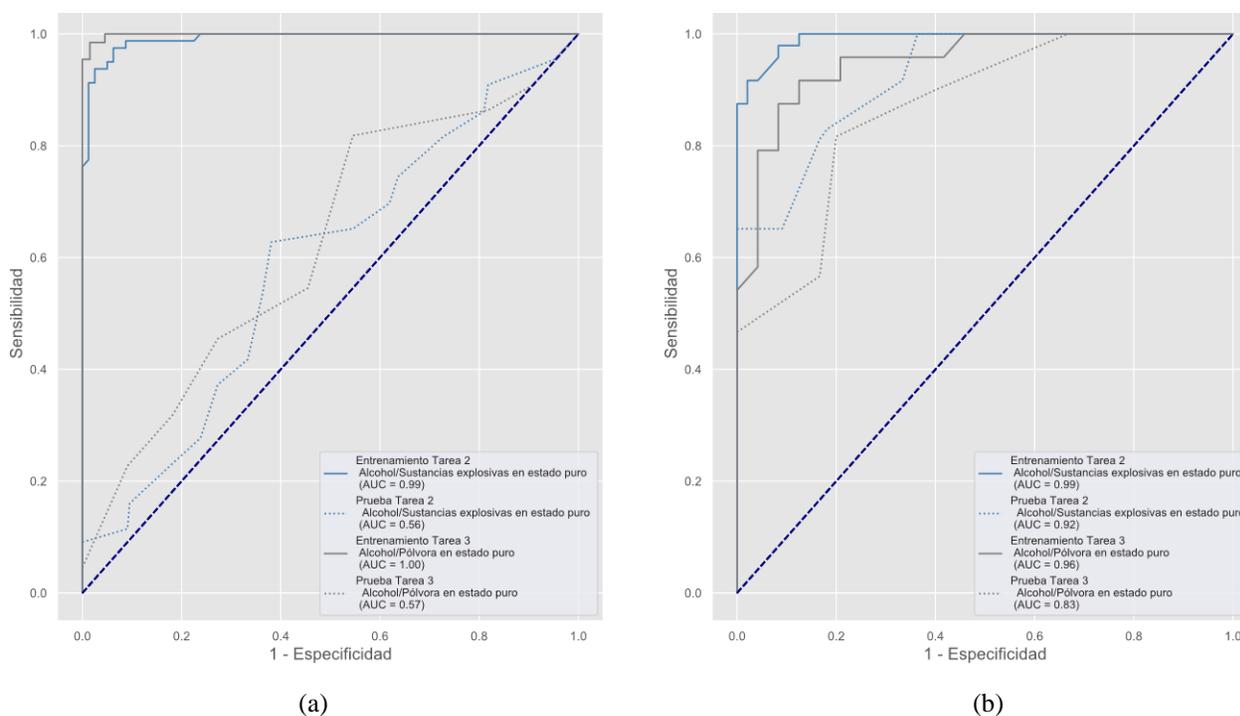
N°	TAREA	#LV	C	AUC	
				TRAIN	TEST
1	Alcohol/Sustancias explosivas en estado puro y mezclas	9	0.01	0.96	0.41
2	Alcohol/Sustancias explosivas en estado puro	9	1	0.99	0.56
3	Alcohol/Pólvora en estado puro	8	0.1	0.99	0.57
4	Alcohol/TNT en estado puro	8	0.01	0.99	0.65
5	Pólvora en estado puro/TNT en estado puro	9	10	1	0.29
6	Alcohol/Pólvora en estado puro y mezcla de pólvora	9	0.1	0.98	0.34
7	Alcohol/TNT en estado puro y mezcla de TNT	8	10	0.99	0.56
8	Alcohol/Pólvora en estado puro/TNT en estado puro	7	1	0.95	0.57

**Tabla 19**

*Base de datos 2- Resultados de modelos de regresión logística*

N°	TAREA	#LV	C	AUC	
				TRAIN	TEST
2	Alcohol/Sustancias explosivas en estado puro	7	0.001	0.99	0.92
3	Alcohol/Pólvora en estado puro	8	0.0001	0.96	0.83
4	Alcohol/TNT en estado puro	8	0.001	1	0.88
5	Pólvora en estado puro/TNT en estado puro	9	1	1	0.43
8	Alcohol/Pólvora en estado puro/TNT en estado puro	8	0.01	0.99	0.68

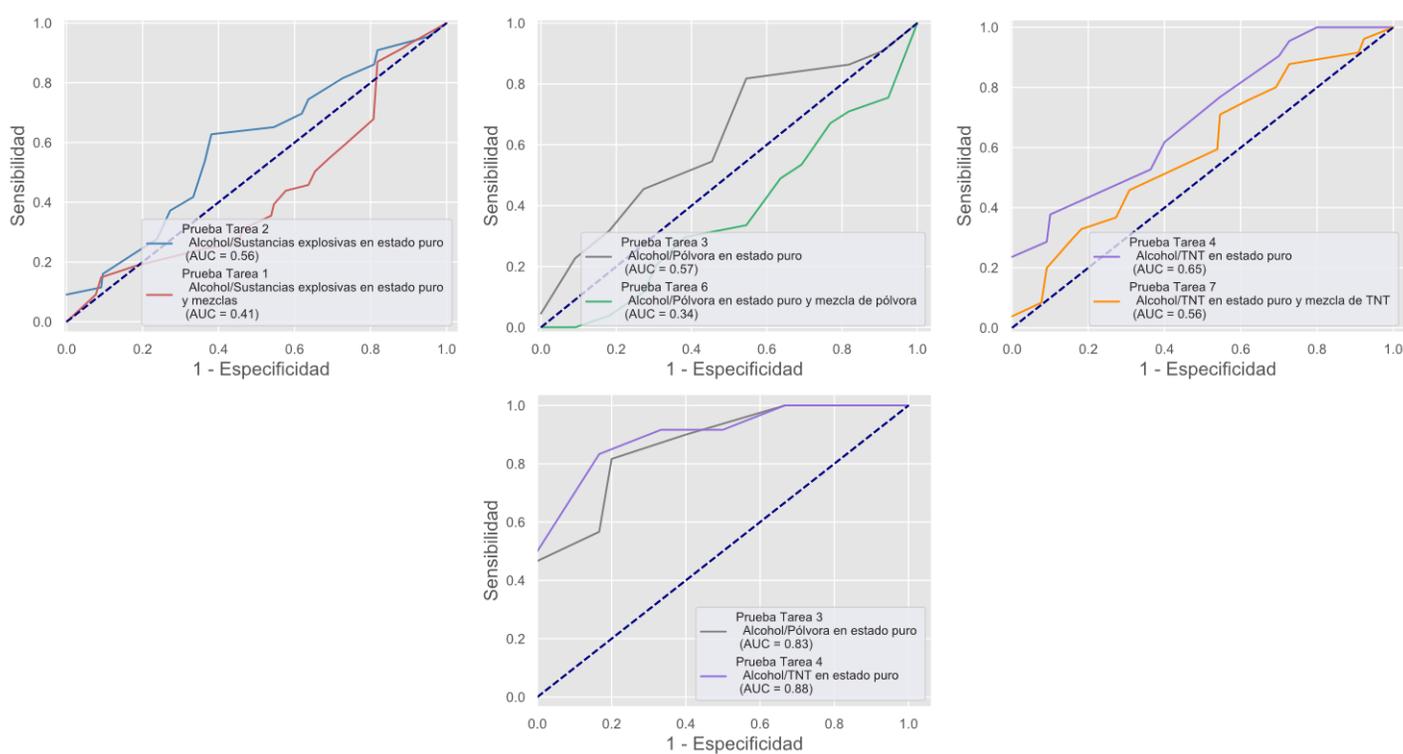
En la Figura 38a, se presentan las curvas ROC y el AUC de dos de los modelos de la base de datos 1, en los cuales se puede apreciar que al igual que en los modelos PLS-DA, el desempeño de los modelos en los datos de entrenamiento es mucho mejor que el de los datos de prueba, ya que la curva ROC se encuentra mucho más cerca de la esquina superior izquierda ( $AUC \gg 0.5$ ) para los datos de entrenamiento representada por líneas continuas, y muy cerca de la región de indecisión ( $AUC \approx 0.5$ ) para los datos de prueba representada por líneas discontinuas. En el caso de la base de datos 2, cómo se observa en la Figura 38b, la diferencia del desempeño entre los datos de entrenamiento y prueba no es significativa (a excepción de los resultados obtenidos en la tarea 5 y 8 cuyos análisis se presentan posteriormente), por lo tanto, se puede asumir que los modelos no se sobreajustan a los datos de entrenamiento.



**Figura 38.** Curvas ROC para la descripción del desempeño de los modelos en datos de entrenamiento y prueba

En cuanto al desempeño de los modelos de la base de datos 1 para clasificación de sustancias explosivas puras y con mezclas este es inferior al de los modelos de clasificación de únicamente

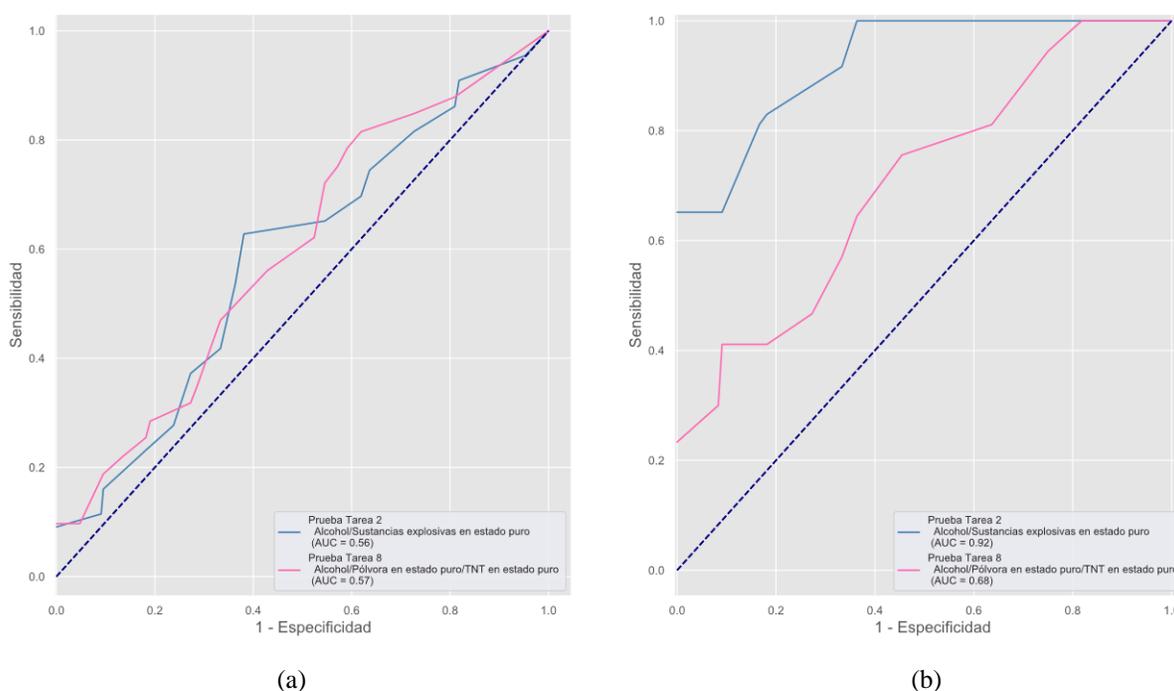
sustancias explosivas puras, como se puede observar en los resultados del AUC y gráficamente en las curvas ROC de la Figura 39. Además, como se muestra en la Figura 39b y Figura 39c, el desempeño de los modelos de clasificación de TNT es superior a los de clasificación de pólvora. En el caso de la base de datos 2, como se muestra en la Figura 39d, el desempeño del modelo de clasificación de TNT es superior al de pólvora, contrario al resultado obtenido con el modelo PLS-DA.



**Figura 39.** Curvas ROC para descripción del desempeño de los modelos para clasificación de sustancias explosivas puras y mezclas

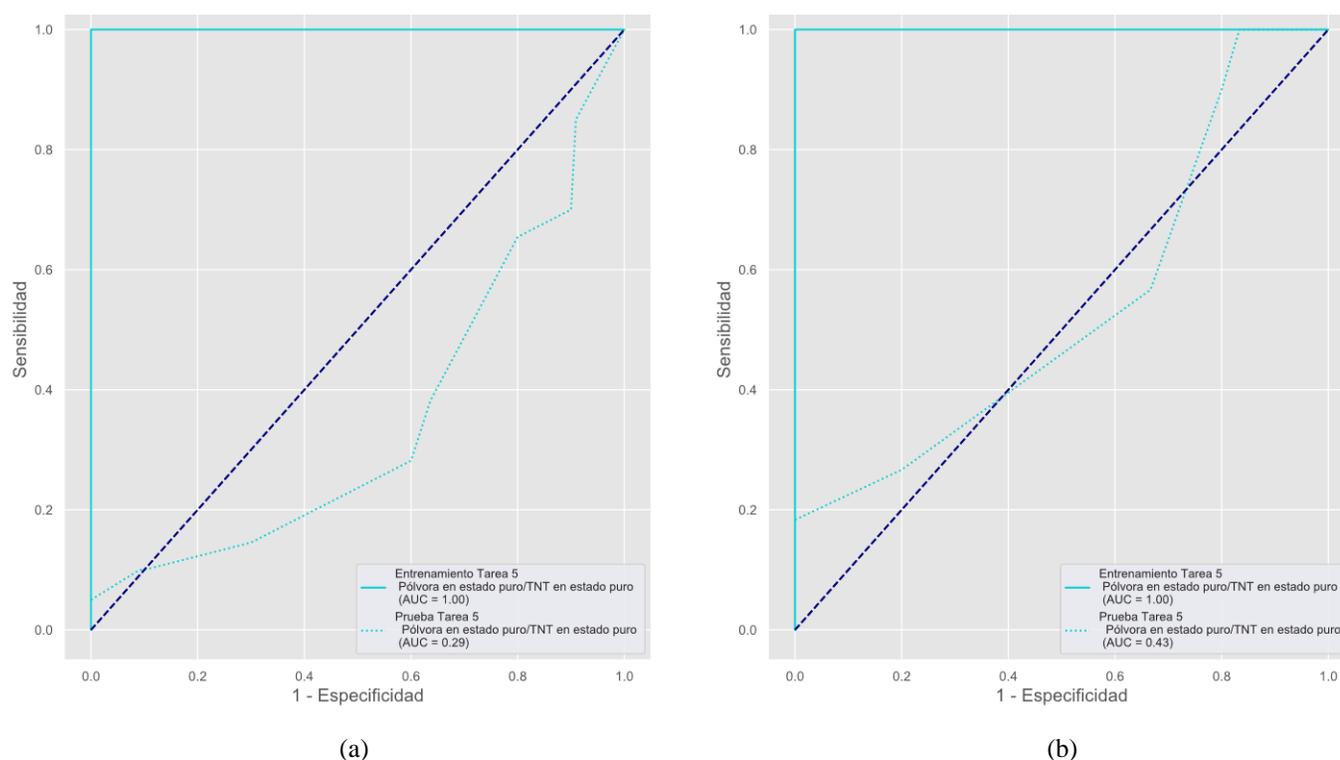
Respecto a la diferencia entre un modelo que determina si una observación es sustancia explosiva o no y uno que determina si una observación no es sustancia explosiva y si lo es el tipo entre TNT y pólvora. Se pudo observar, como se muestra en la Figura 40a para la base de datos 1, que el desempeño de un modelo que únicamente determina si una sustancia es o no explosiva es

bajo y similar al modelo que determina además el tipo de sustancia explosiva que es. Por lo cual, si estos modelos hubieran sido utilizados en las condiciones de los experimentos realizados para elaborar la base de datos 1, el prototipo tendría una baja capacidad de clasificación entre sustancias explosivas y no explosivas con el modelo de la Tarea 2 y por ende una baja capacidad de clasificación entre TNT, pólvora y alcohol con el modelo de la Tarea 8. En el caso de la base de datos 2, como se muestra en la Figura 40b, el desempeño de un modelo que únicamente determina si una sustancia es o no explosiva es superior al del modelo que determina además el tipo de sustancia explosiva que es. Además, al comparar el desempeño entre los modelos realizados con bajas concentraciones de sustancia explosiva (experimentos de la base de datos 1 con concentraciones entre 0.1 y 3gr de pólvora y TNT) y de los realizados con mayores concentraciones (experimentos de la base de datos 2 con concentraciones entre 3 y 5gr de pólvora y TNT), se pudo observar que se obtuvo mejores resultados al aumentar el nivel de concentración de sustancia.



**Figura 40.** Curvas ROC para descripción del desempeño de los modelos de clasificación entre dos y tres clases de sustancias

Por último, al analizar los resultados de los modelos de clasificación entre TNT y pólvora de la Figura 41a para la base de datos 1 y Figura 41b para la base de datos 2, se determinó que los modelos no son capaces de clasificar entre TNT y pólvora ( $AUC < 0.5$ ), a pesar de haber aumentado la concentración de sustancia explosiva para elaborar el modelo de la base de datos 2. Además, en las figuras se identifica que el resultado de la clasificación de los datos de entrenamiento alta ( $AUC = 1$ ) en comparación a los de prueba, es decir, los datos de entrenamiento se sobreajustan a los datos de prueba a pesar de tener un parámetro de regularización L2 que busca penalizar el sobreajuste de los modelos a los datos de entrenamiento.



**Figura 41.** Curvas ROC para descripción del desempeño de los modelos de clasificación entre pólvora y TNT

En este capítulo se presentó una descripción de los conceptos utilizado para el desarrollo de los modelos de clasificación con el método de regresión logística. A continuación, se describió el

proceso para la generación de los modelos en el que se definió el tipo de escalamiento usado y cuál fue la aplicación que se dio a la técnica de validación cruzada. Para terminar, se presentó una descripción de los resultados obtenidos para cada una de las tareas.

# CAPÍTULO 6

## Redes Neuronales Artificiales

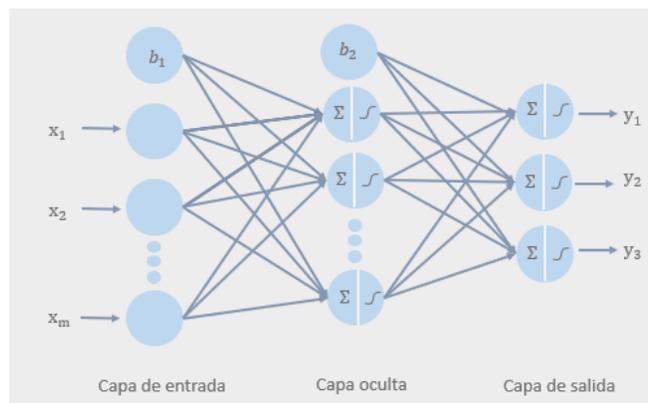
En el presente capítulo se presenta una descripción de los conceptos utilizados en el desarrollo de los modelos de clasificación y regresión mediante redes neuronales artificiales denominadas perceptrón multicapa (MLP). Posteriormente, se presenta el procedimiento para la elaboración de los modelos, en el que se describe el tipo de tratamiento previo al ingreso de los datos a los algoritmos y los hiperparámetros seleccionados durante la validación cruzada. Finalmente, se muestran los resultados obtenidos con los modelos desarrollados en base al desempeño evaluado con las métricas mencionadas en el Capítulo 2.

### 6.1. Conceptos Básicos

#### 6.1.1. Red Perceptrón Multicapa (MLP)

La red perceptrón multicapa (MLP), es un estimador no lineal que puede ser usado tanto para tareas de clasificación como de regresión, cuya estructura se basa en una conexión completa de la red (Alpaydm, 2012). Como se observa en el ejemplo de la Figura 42, una red MLP está compuesta por una capa de entrada, una o múltiples capas ocultas y una capa de salida. En la capa de entrada se encuentra un conjunto de  $m$  neuronas, en las cuales se realiza el producto de las  $m$  características de los datos de entrada  $x$  con sus pesos  $w$  correspondientes. En cada una de las neuronas de la capa oculta se realiza una suma ponderada  $\Sigma$  de sus entradas  $(x_1, x_2, \dots, x_m)$ , seguida

por una función de activación no lineal. La capa de salida recibe los valores de la última capa oculta y los transforma en las salidas ( $y_1, y_2, y_3$ ).



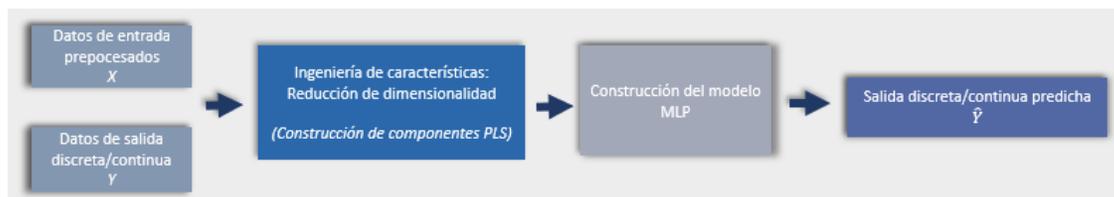
**Figura 42.** Red Perceptrón Multicapa con una capa oculta

La grafica muestra un ejemplo de una red perceptrón multicapa, los círculos celestes representan las neuronas y las flechas grises las conexiones entre ellas.

El entrenamiento de una red MLP consiste en ingresar en la red las características de cada observación y calcular la salida de cada neurona en las capas consecutivas, posteriormente, se debe medir el error de salida de la red que se define como la diferencia de la salida deseada y la real, calcular cuánto contribuyo cada neurona de la última capa oculta al error y así sucesivamente con cada capa oculta hasta llegar a la capa de entrada. Este proceso de entrenamiento en el cual el cálculo del error se lo realiza desde la capa de salida hacia la de entrada se lo conoce con el nombre de retropropagación. Finalmente, se modifica ligeramente los pesos del modelo para reducir el error y se actualiza el modelo. Todo el proceso se realiza por un numero de épocas definidas por el diseñador (Geron, 2017).

Para la construcción de los modelos de redes neuronales se definieron dos etapas principales: La reducción de la dimensionalidad de los datos de entrada mediante la construcción de

componentes PLS y la construcción del modelo de predicción, como se muestra en la Figura 43, su descripción se presenta en las subsecciones siguientes.



**Figura 43.** Proceso para la construcción de modelo de red neuronal MLP

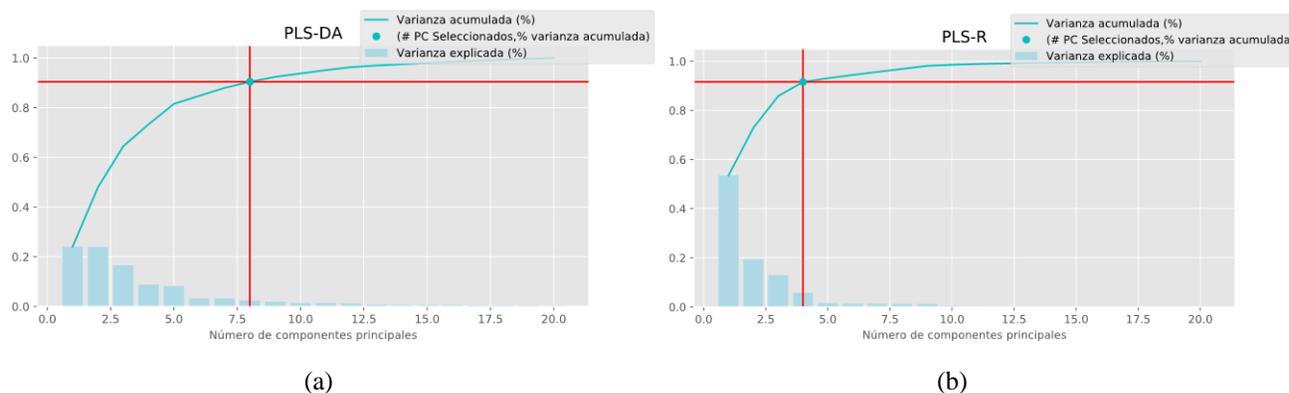
El proceso de construcción del modelo consta de cuatro etapas: El ingreso de los datos de entrada y salida, la construcción de variables latentes, la construcción del modelo de predicción y la obtención de las salidas predichas.

Se realizó la reducción de dimensionalidad de los datos de entrada de la red neuronal debido a que son series temporales, por tanto, sus características son altamente correlacionadas. El uso de este tipo de características correlacionadas en la elaboración de modelos puede provocar un sobreajuste de los modelos a los datos de entrenamiento, además, al no utilizar una técnica de reducción de dimensionalidad se proporcionan características de entrada redundantes a la red neuronal, ocasionando que se realicen cálculos redundantes con neuronas y pesos redundantes afectando así el tiempo de entrenamiento de la red. En consecuencia, al utilizar el método PLS se podrá reducir la dimensionalidad de las características de entrada de forma supervisada, para transformar las variables correlacionadas en menos variables no correlacionadas donde se incorporen las correlaciones originales (Samarasinghe, 2006).

### 6.1.2. Construcción de Componentes PLS

Para reducir la dimensionalidad de los datos de entrenamiento se utilizó el método de mínimos cuadrados parciales (PLS) definido en el Capítulo 4 con el fin de eliminar características redundantes y mantener un modelo simple que se generalice adecuadamente a los datos de

entrenamiento. La selección del número de variables latentes que posteriormente ingresarían al algoritmo MLP se realizó mediante la elaboración de gráficas de sedimentación, como se muestra en la Figura 44 a y b, que muestra el porcentaje de varianza explicado por cada una de las variables latentes y el porcentaje acumulado. La gráfica permitió identificar el punto en el que el descenso del porcentaje de varianza y el ascenso del porcentaje acumulado empezó a estabilizarse, este se encontró en la octava variable latente para el modelo de discriminación entre alcohol, pólvora y TNT de la base de datos 2 y en la cuarta para el modelo de cuantificación, estas variables latentes explican alrededor del 90% de varianza de los datos.



**Figura 44.** Grafica de sedimentación para selección de variables latentes

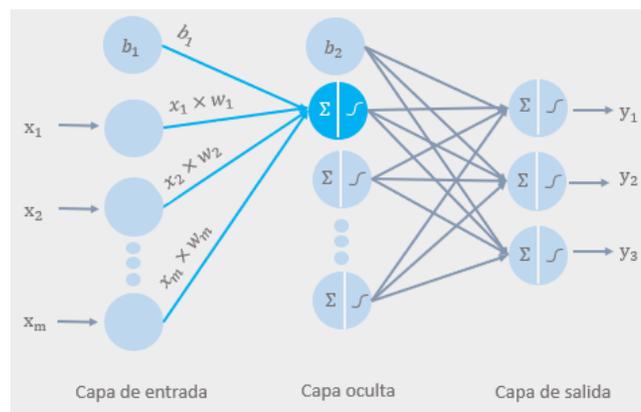
Las gráficas muestran el número de variables latentes seleccionadas para el modelo de clasificación en la izquierda y de regresión en la derecha mediante puntos. Las barras muestran la varianza explicada por cada una de las variables latentes y las líneas azul la varianza acumulada en cada una.

### 6.1.3. Construcción del Modelo MLP

Para la construcción de la red neuronal MLP se debe definir la arquitectura de la red, las funciones de activación de las neuronas de las capas ocultas y de salida, la función objetivo que va a ser minimizada durante el proceso de aprendizaje y la función de optimización que minimizara la función objetivo, cada uno de estos parámetros se describen a continuación.

### 6.1.3.1. Arquitectura de la Red

Para la selección de la arquitectura de la red, se definió el número de capas ocultas que tendría la red, en este caso se seleccionó únicamente una capa oculta para posteriormente poder comparar el desempeño de este tipo de red neuronal clásica con una red neuronal profunda, presentada en el Capítulo 7. Además, debido a que es una red MLP todas las neuronas de la capa de entrada se conectarán a la capa oculta y todas estas a la capa de salida, como se observa en la Figura 45.



**Figura 45.** Número de capas ocultas seleccionadas para el modelo MLP

La grafica muestra la red MLP con una capa oculta seleccionada para el desarrollo de los modelos de clasificación y cuantificación, con celeste se identifica a la neurona cuyas operaciones se describen en el texto.

### 6.1.3.2. Funciones de Activación

Como se observa en la Figura 45 en cada una de las neuronas de la capa oculta y de salida se realiza la suma ponderada  $\Sigma$  de sus entradas. En la Ecuación 44 se presenta la suma ponderada de las entradas de la neurona identificada con color celeste.

$$\Sigma: x_1w_1 + x_2w_2 + \dots + x_mw_m + b_1 \quad ( 44 )$$

Donde  $x_m$ , representa la entrada  $m$  de la neurona y  $b_1$  el termino de bias o intercepción en el eje  $y$  de la capa. A continuación, la suma ponderada  $\Sigma$  es transformada mediante una función de activación  $f(\Sigma)$  en la salida de la neurona, como se observa en la Ecuación 45.

$$f(\Sigma): \quad f(x_1w_1 + x_2w_2 + \dots + x_mw_m + b_1) \quad ( 45 )$$

La función de activación en las capas ocultas debe ser diferenciable, ya que el proceso de aprendizaje involucra el método de retropropagación que buscará la combinación de pesos que minimice la función objetivo mediante un método de optimización. Dado que este método requiere el cálculo de la gradiente de la función objetivo se debe garantizar su diferencialidad por lo que al tener una función de activación diferenciable el error también lo será (Rojas, 1996). Además, es importante hacer uso de funciones de activación en la capa oculta que no sean únicamente diferenciables sino también no lineales, ya que la suma de funciones lineales producirá una función lineal por lo que no se podrá beneficiar del uso de capas ocultas y no será posible obtener representaciones más complejas de los datos de entradas (Chollet, 2018). Existen diversos tipos de funciones de activación para las capas ocultas, entre las más usadas se encuentran las siguientes:

- **Función logística sigmoide:** Transforma la entrada  $z$  en un valor entre 0 y 1 mediante la Ecuación 46.

$$\phi(z) = \frac{1}{1 + e^{-z}} \quad ( 46 )$$

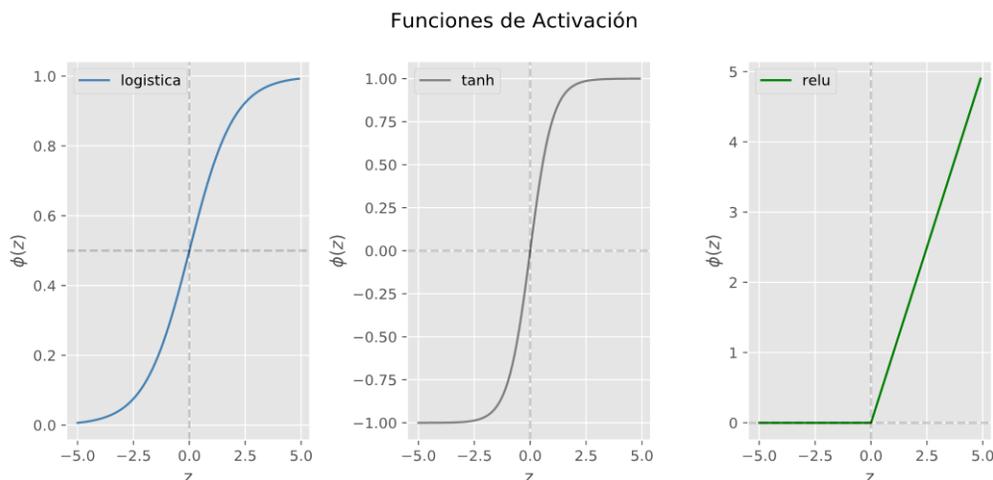
- **Función tangente hiperbólica:** Transforma la entrada  $z$  en un valor entre 1 y -1 mediante la Ecuación 47.

$$\phi(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}} \quad ( 47 )$$

- **Función ReLU:** La unidad lineal rectificada (ReLU) es una función continua definida por la Ecuación 48.

$$\phi(z) = \begin{cases} z, & z \geq 0 \\ 0, & z < 0 \end{cases} \quad ( 48 )$$

En la Figura 46, se presenta la representación gráfica de las funciones de activación para las capas ocultas definidas anteriormente. Como se puede observar, las tres funciones son continuas y diferenciables, sin embargo, para valores extremos la curva de las funciones logística sigmoide y tangente hiperbólica tienden a responder muy poco a los cambios en  $z$ , por lo cual, en estos extremos su derivada será muy cercana a cero provocando que el entrenamiento sea muy lento, la función relu por el contrario es menos costosa computacionalmente ya que utiliza funciones matemáticas simples y su derivada para valores positivos será 1 (Raschka & Mirjalili, 2017).



**Figura 46.** Funciones de activación para las capas ocultas de una red neuronal MLP

Las funciones de activación para la capa de salida dependen de la tarea del modelo, para tareas de clasificación la función más usada es softmax, definida mediante la Ecuación 49, la cual calcula la probabilidad de que la entrada  $z$  pertenezca la clase  $i$ , donde  $k$  representa el número de clases.

- **Función softmax**

$$\phi(z) = \frac{e^{z_i}}{\sum_{j=1}^P e^{z_j}} \quad ( 49 )$$

Para tareas de regresión cuya salida son valores continuos arbitrarios la función de activación para la capa oculta es la función identidad de la Ecuación 50.

- **Función identidad**

$$\phi(z) = z \quad ( 50 )$$

### 6.1.3.3. Función Objetivo

Para poder entrenar la red neuronal MLP es necesario utilizar una función objetivo que mida el error entre las predicciones realizadas con la red y el objetivo verdadero (salidas que se quería que la red produjera), obteniendo una medida de que tan bien se desempeñó la red para una entrada específica. Esta medida será utilizada como señal de retroalimentación para ajustar el valor de los pesos y de esta forma disminuir la medida obtenida con la función objetivo (Chollet, 2018).

La selección de la función objetivo depende de la tarea del modelo, para tareas de regresión una función objetivo común es el error cuadrático medio (MSE) presentado en la Ecuación 51, el cual calcula el promedio del cuadrado de la diferencia entre la salida predicha  $\hat{y}$  y el objetivo  $y$  (Rao & McMahan, 2019).

$$J_{MSE}(y, \hat{y}) = \frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2 \quad ( 51 )$$

Para tareas de clasificación la métrica comúnmente usada es la entropía categórica cruzada presentada en la Ecuación 52, con la cual se calcula el promedio del producto del logaritmo negativo de la salida predicha  $\hat{y}$  con el objetivo  $y$  (Rao & McMahan, 2019).

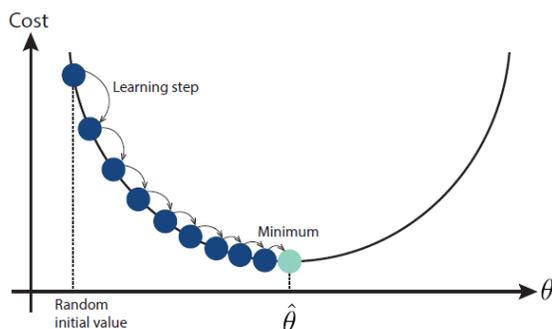
$$J(\theta) = -\frac{1}{n} \sum_{i=0}^n \sum_{j=0}^m y_{(i,j)} \times \log (\hat{y}_{(i,j)}) \quad ( 52 )$$

Donde  $n$  representa el número de datos de entrenamiento,  $m$  el número de clases,  $y_{(i,j)}$  la salida real de la observación  $i$  para la clase  $j$  y  $h_{\theta}(x_{i,j})$  la salida predicha por el modelo.

#### 6.1.3.4. Función de Optimización

La función de optimización es aquella que implementa el algoritmo de retropropagación para ajustar los pesos del modelo según la información proporcionada por la función objetivo (Chollet, 2018). Existen diversos tipos de funciones de optimización para redes neuronales, las principales se detallan a continuación.

- **Descenso de gradiente (GD):** Este algoritmo de optimización se encarga de ajustar los pesos del modelo de forma iterativa para minimizar la función objetivo. Inicia con valores aleatorios para los pesos y gradualmente mediante pasos intenta disminuir la función objetivo hasta que el algoritmo converge en el mínimo, como se observa en la Figura 47. En cada iteración o época se dará un paso en dirección opuesta del gradiente, esta dirección depende del gradiente local de la función objetivo con respecto a los pesos y el tamaño de los pasos por cada iteración del hiperparámetro de la tasa de aprendizaje, si esta es muy pequeña el modelo tendrá que realizar muchas iteraciones para converger o si es muy grande el algoritmo puede divergir sin encontrar una buena solución saltando de un lado a otro de la función objetivo (Geron, 2017).



**Figura 47.** Función de optimización: Descenso de gradiente

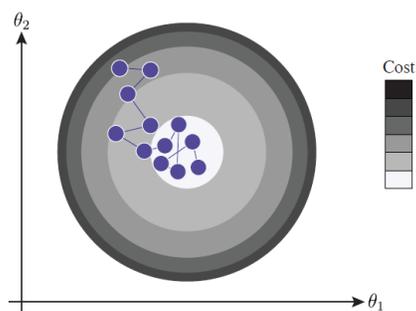
Fuente: (Geron, 2017)

Para implementar el descenso de gradiente se debe calcular por cada iteración la derivada parcial de la función objetivo con respecto a cada peso para todas las observaciones del conjunto de entrenamiento. Con el valor del gradiente calculado  $\nabla J(w)$  se puede actualizar los pesos mediante la suma de estos con el cambio de peso  $\Delta w$ , como se observa en la Ecuación 53. El cambio de peso se define como el producto del gradiente negativo de la función objetivo  $\nabla J(w)$  por la tasa de aprendizaje  $\eta$ , como se muestra en la Ecuación 54 (Geron, 2017).

$$w = w + \Delta w \quad ( 53 )$$

$$\Delta w = -\eta \nabla J(w) \quad ( 54 )$$

- **Descenso de gradiente estocástico (SGD):** Este algoritmo selecciona una instancia aleatoria del conjunto de entrenamiento en cada paso y calcula los gradientes basados únicamente en esta instancia, esto permite que el algoritmo sea más rápido ya que tiene menos datos que manipular en cada iteración. Además, debido a su naturaleza aleatoria, la función objetivo rebotará de arriba y abajo, como se observa en la Figura 48, por lo cual habrá una mayor posibilidad de encontrar el mínimo global cuando la función objetivo es muy irregular (Geron, 2017).



**Figura 48.** Función de optimización: Descenso de gradiente estocástico

Fuente: (Geron, 2017)

Para el desarrollo de los modelos con redes neuronales se decidió utilizar la función de optimización SGD y las funciones objetivo y de activación para la capa de salida propuestas en la teoría para las tareas de clasificación y regresión. El proceso para el desarrollo de los modelos se detalla en la siguiente sección.

## 6.2. Generación de Modelos de Redes Neuronales MLP

El proceso de desarrollo de los modelos MLP, como se mencionó en la sección anterior, consta de dos etapas principales: la reducción de la dimensionalidad de los datos mediante la construcción de variables latentes y la construcción del modelo de predicción. Las variables latentes se obtuvieron con ayuda de la función *PLSRegression* de la librería *sklearn*, la cual utiliza el algoritmo de proceso iterativo no lineal de mínimos cuadrados parciales (NIPALS) para el cálculo de los componentes PLS. Antes de ingresar los datos de entrenamiento al algoritmo *PLSRegression* se realizó un escalamiento de las características de los datos para que sean más fáciles de interpretar por el algoritmo. A continuación, se seleccionaron el número de variables latentes que ingresarían al algoritmo MLP, mediante el proceso indicado en la subsección 6.1.2. Finalmente, se utilizó las funciones *MLPClassifier* y *MLPRegressor*, de la librería *sklearn* para entrenar los modelos de

clasificación y regresión, respectivamente, con los hiperparámetros seleccionados durante la validación cruzada.

### **6.2.1. Escalamiento**

Debido a que el desempeño del método de mínimos cuadrados parciales depende de la escala de los datos se seleccionó los métodos de escalamiento empleados en el Capítulo 5. Los cuales centran las características de los datos en la media y los escala con una varianza igual a 1 para la reducción de dimensionalidad para tareas de clasificación y únicamente los centra en la media para tareas de regresión.

### **6.2.2. Balanceo de Clases**

Al tener un mayor número de observaciones de una clase que de otra, el algoritmo de aprendizaje aprenderá implícitamente un modelo que optimiza las predicciones basadas en la clase más abundante en el conjunto de datos (Raschka & Mirjalili, 2017). Una forma de solucionar este problema es balancear la cantidad de observaciones en cada clase ya sea mediante una técnica de sobre muestro o un submuestreo de los datos. En este caso se utilizó el algoritmo *SMOTE*, el cual sobre muestrea la clase minoritaria creando observaciones sintéticas (Chawla et al., 2002).

### **6.2.3. Validación Cruzada**

Se utilizó la validación cruzada de  $k=10$  iteraciones mencionada en el Capítulo 2 para seleccionar el número de neuronas en la capa oculta y su función de activación respectiva, la tasa de aprendizaje para la función de optimización y el número de épocas para el entrenamiento del modelo, con la diferencia de que la métrica utilizada para la selección de estos hiperparámetros se realizó con el valor del error cuadrático medio (MSE) para los modelos de regresión y con el valor

del área bajo la curva (AUC) para los de clasificación. Finalmente, los modelos se entrenaron con los hiperparámetros seleccionados mediante la validación cruzada. Los resultados del desempeño de los modelos se presentan en la siguiente sección.

## 6.3. Resultados del Desempeño de los Modelos de Redes

### Neuronales MLP

#### 6.3.1. Modelos de Clasificación

En la **Tabla 20** y **Tabla 21**, se presentan los resultados obtenidos con los modelos de clasificación para la base de datos 1 y 2, en los cuales se describe el número de variables latentes (LV) utilizado para la reducción de dimensionalidad, el número de neuronas de la red, la tasa de aprendizaje y el número de épocas, además del AUC de los datos de entrenamiento y prueba para la evaluación desempeño de los modelos. Estos resultados fueron graficados para el análisis que se realizó a continuación.

**Tabla 20**

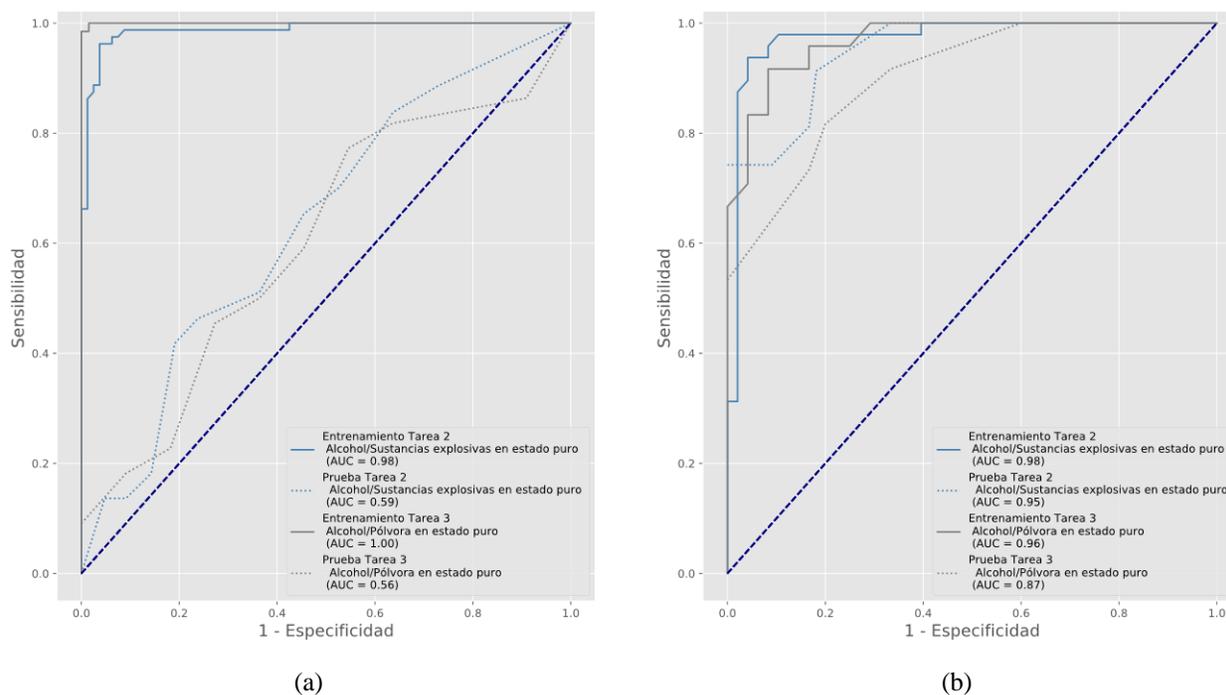
*Base de datos 1- Resultados de modelos de clasificación MLP*

N°	TAREA	#LV	#NEURONAS	TASA DE APRENDIZAJE	ÉPOCAS	AUC	
						TRAIN	TEST
1	Alcohol/Sustancias explosivas en estado puro y mezclas	9	6	0.01	60	0.97	0.43
2	Alcohol/Sustancias explosivas en estado puro	9	4	0.1	40	0.98	0.59
3	Alcohol/Pólvora en estado puro	8	6	0.01	200	1	0.56
4	Alcohol/TNT en estado puro	8	200	0.01	10	0.99	0.62
5	Pólvora en estado puro/TNT en estado puro	9	6	0.1	40	1	0.33
6	Alcohol/Pólvora en estado puro y mezcla de pólvora	9	80	0.01	20	0.99	0.32
7	Alcohol/TNT en estado puro y mezcla de TNT	8	80	0.01	200	0.99	0.50
8	Alcohol/Pólvora en estado puro/TNT en estado puro	7	80	0.1	10	0.96	0.56

**Tabla 21***Base de datos 2- Resultados de modelos de clasificación MLP*

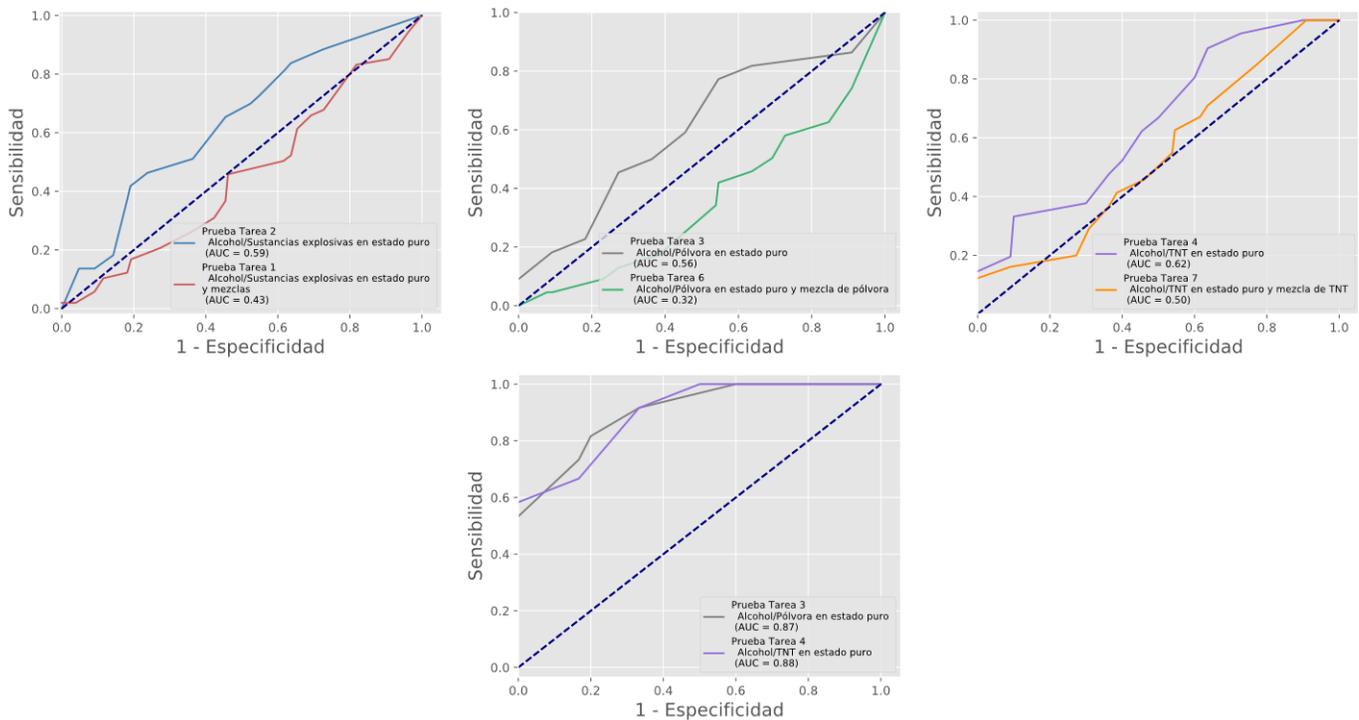
N°	TAREA	#LV	#NEURONAS	TASA DE APRENDIZAJE	ÉPOCAS	AUC	
						TRAIN	TEST
2	Alcohol/Sustancias explosivas en estado puro	7	16	0.001	100	0.98	0.95
3	Alcohol/Pólvora en estado puro	8	6	0.1	80	0.96	0.87
4	Alcohol/TNT en estado puro	8	4	0.1	50	1	0.88
5	Pólvora en estado puro/TNT en estado puro	9	200	0.1	20	1	0.45
8	Alcohol/Pólvora en estado puro/TNT en estado puro	8	60	0.1	20	0.99	0.71

En la Figura 49a, se presentan las curvas ROC y el AUC de dos de los modelos de la base de datos 1, en los cuales se puede apreciar que al igual que en los modelos PLS-DA, el desempeño de los modelos en los datos de entrenamiento es mucho mejor que el de los datos de prueba, ya que la curva ROC se encuentra mucho más cerca de la esquina superior izquierda ( $AUC \gg 0.5$ ) para los datos de entrenamiento representada por líneas continuas, y muy cerca de la región de indecisión ( $AUC \approx 0.5$ ) para los datos de prueba representada por líneas discontinuas. En el caso de la base de datos 2, cómo se observa en la Figura 49b, la diferencia del desempeño entre los datos de entrenamiento y prueba no es significativa, a excepción de los resultados obtenidos en la tarea 5 cuyo análisis se presenta posteriormente, por lo tanto, se puede asumir que los modelos no se sobreajustan a los datos de entrenamiento.



**Figura 49.** Curvas ROC para la descripción del desempeño de los modelos en datos de entrenamiento y prueba

Respecto al desempeño de los modelos de la base de datos 1 para clasificación de sustancias explosivas puras y con mezclas este es inferior al de los modelos de clasificación de únicamente sustancias explosivas puras, como se puede observar en los resultados del AUC y gráficamente en las curvas ROC de la Figura 50. Además, como se muestra en la Figura 50b y Figura 50c, el desempeño de los modelos de clasificación de TNT es superior a los de clasificación de pólvora. En el caso de la base de datos 2, como se muestra en la Figura 50d, el desempeño del modelo de clasificación de TNT es superior al de pólvora. Por lo tanto, se espera que bajo las condiciones actuales del prototipo y con concentraciones entre 3 y 5gr de sustancias explosivas este sea capaz de clasificar correctamente en mayor porcentaje observaciones de TNT que de pólvora con los modelos de redes neuronales.

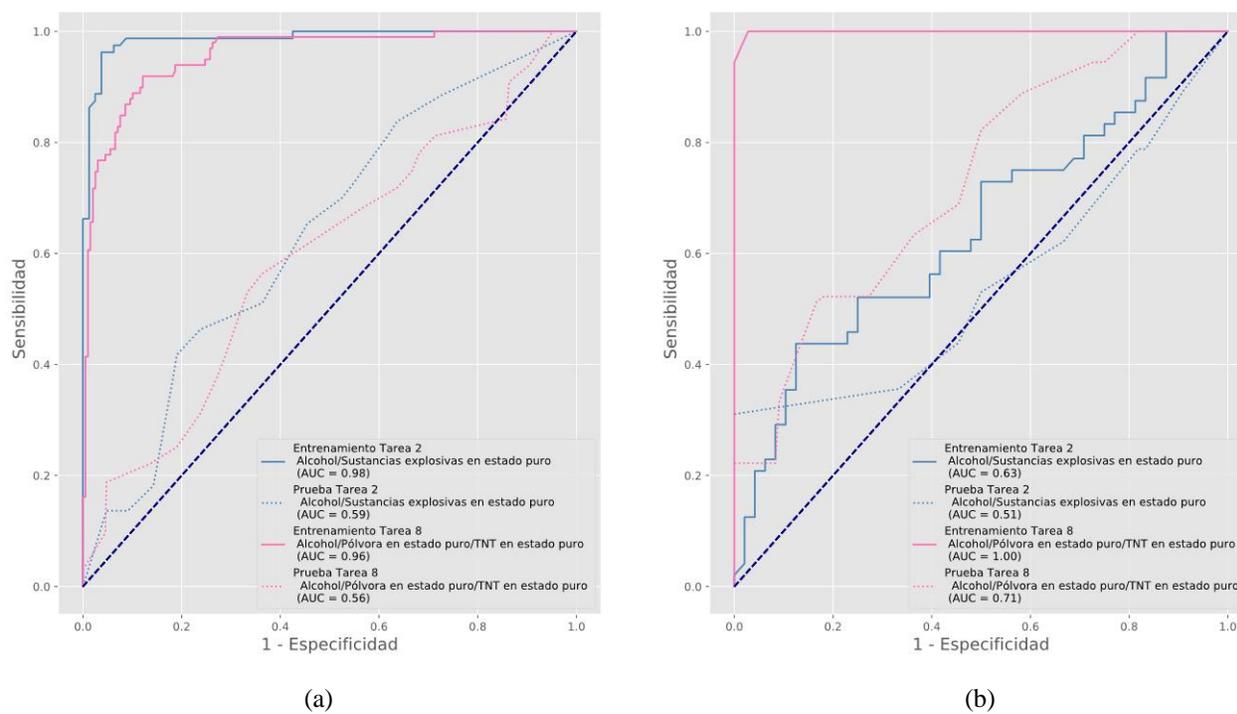


**Figura 50.** Curvas ROC para descripción del desempeño de los modelos para clasificación de sustancias explosivas puras y mezclas

En cuanto a la diferencia entre un modelo que determina si una observación es sustancia explosiva o no y uno que determina si una observación no es sustancia explosiva y si lo es el tipo entre TNT y pólvora. Se pudo observar, como se muestra en la Figura 51a para la base de datos 1, que el desempeño de un modelo que únicamente determina si una sustancia es o no explosiva es similar al modelo que determina además el tipo de sustancia explosiva que es. Por lo cual, si estos modelos son utilizados en las condiciones de los experimentos realizados para elaborar la base de datos 1, el prototipo tendría una capacidad de clasificación entre sustancias explosivas y no explosivas con el modelo de la Tarea 2 similar a la capacidad de clasificación entre TNT, pólvora y alcohol con el modelo de la Tarea 8.

En el caso de la base de datos 2, como se muestra en la Figura 51b, el desempeño de un modelo que únicamente determina si una sustancia es o no. Además, al comparar el desempeño

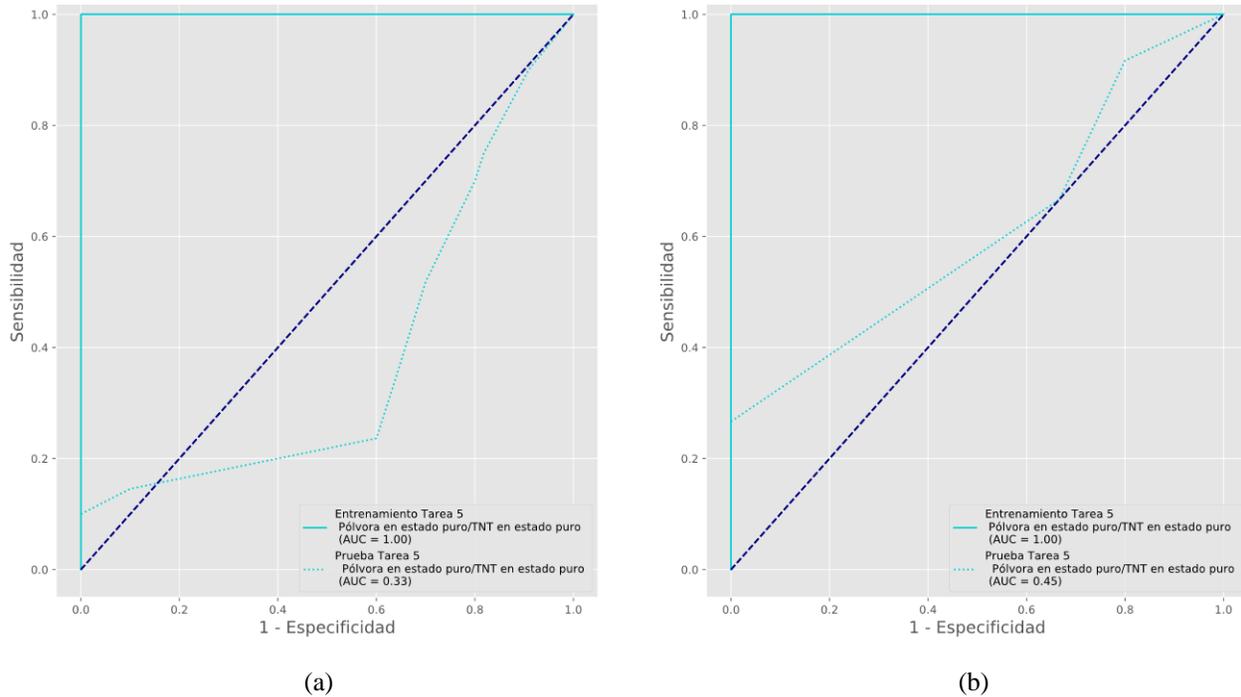
entre los modelos realizados con bajas concentraciones de sustancia explosiva (experimentos de la base de datos 1 con concentraciones entre 0.1 y 3gr de pólvora y TNT) y de los realizados con mayores concentraciones (experimentos de la base de datos 2 con concentraciones entre 3 y 5gr de pólvora y TNT), se pudo observar que se obtuvo mejores resultados al aumentar el nivel de concentración de sustancia.



**Figura 51.** Curvas ROC para descripción del desempeño de los modelos de clasificación entre dos y tres clases de sustancias

Finalmente, al analizar los resultados de los modelos de clasificación entre TNT y pólvora de la Figura 52a para la base de datos 1 y Figura 52b para la base de datos 2, se determinó que el prototipo e-nose tiene un mal desempeño ( $AUC < 0.5$ ) al clasificar pólvora y TNT, sin embargo, se obtuvo un mejor resultado con el modelo de la base de datos 2 para clasificación de sustancias con concentraciones entre 3 y 5gr. Además, en las figuras se identifica que el resultado de la

clasificación de los datos de entrenamiento es alto ( $AUC=1$ ) en comparación a los de prueba, es decir, los datos de entrenamiento se sobreajustan a los datos de prueba.



**Figura 52.** Curvas ROC para descripción del desempeño de los modelos de clasificación entre pólvora y TNT

### 6.3.2. Modelos de Regresión

En la Tabla 22 y **Tabla 23**, se presentan los resultados de los modelos de regresión para la base de datos 1 y 2, respectivamente, en los cuales se describe el número de variables latentes (LV) utilizado para la reducción de dimensionalidad, el número de neuronas de la red, la tasa de aprendizaje y el número de épocas, además del MSE y R2 de los datos de entrenamiento y prueba para la evaluación desempeño de los modelos. Estas métricas en conjunto con las gráficas de los valores reales vs los valores predichos por los modelos fueron utilizadas para el análisis que se realizó a continuación.

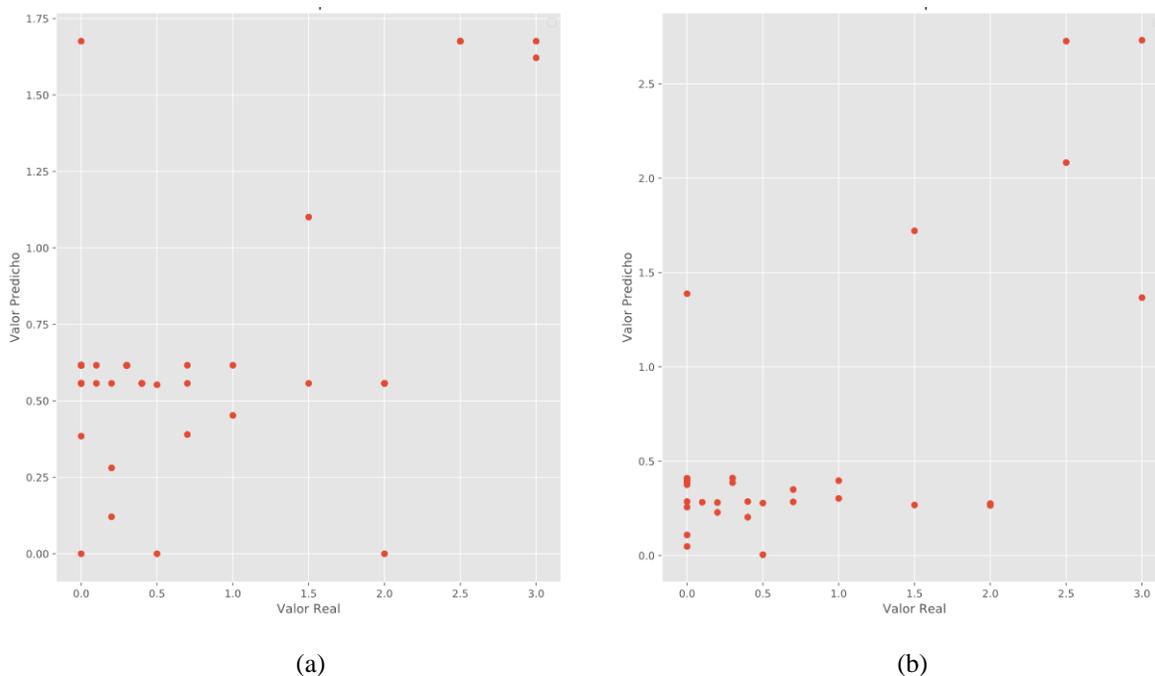
**Tabla 22***Base de datos 1- Resultados de modelos de regresión MLP*

N°	TAREA	#LV	#NEURONAS	TASA DE APRENDIZAJE	ÉPOCAS	MSE		R2	
						TRAIN	TEST	TRAIN	TEST
1	Alcohol/Sustancias explosivas en estado puro y mezclas	2	60	0.001	400	0.45	0.59	0.37	0.30
2	Alcohol/Sustancias explosivas en estado puro	2	400	0.001	400	0.39	0.47	0.51	0.47
3	Alcohol/Pólvora en estado puro	2	200	0.001	40	0.38	0.59	0.34	0.23
4	Alcohol/TNT en estado puro	2	200	0.001	200	0.37	0.80	0.46	-0.01
5	Pólvora en estado puro/TNT en estado puro	2	80	0.001	400	0.46	0.39	0.50	0.57
6	Alcohol/Pólvora en estado puro y mezcla de pólvora	2	60	0.01	80	0.52	0.67	0.02	0.05
7	Alcohol/TNT en estado puro y mezcla de TNT	2	100	0.001	200	0.51	0.61	0.21	0.22
8	Alcohol/Pólvora en estado puro/TNT en estado puro	2	1000	0.001	200	0.36	0.58	0.26	0.02

**Tabla 23***Base de datos 2- Resultados de modelos de regresión MLP*

N°	TAREA	#LV	#NEURONAS	TASA DE APRENDIZAJE	ÉPOCAS	MSE		R2	
						TRAIN	TEST	TRAIN	TEST
2	Alcohol/Pólvora en estado puro/TNT en estado puro	4	1000	0.01	20	2.06	3.71	0.47	0.10
3	Alcohol/Pólvora en estado puro	4	2	0.001	200	2.08	5.92	0.50	-0.37
4	Alcohol/TNT en estado puro	3	2	0.1	60	2.48	3.59	0.43	0.17
5	Pólvora en estado puro/TNT en estado puro	3	1000	0.001	400	0.44	1.19	0.29	-0.63
8	Alcohol/Pólvora en estado puro/TNT en estado puro	4	1000	0.01	20	2.06	3.71	0.47	0.10

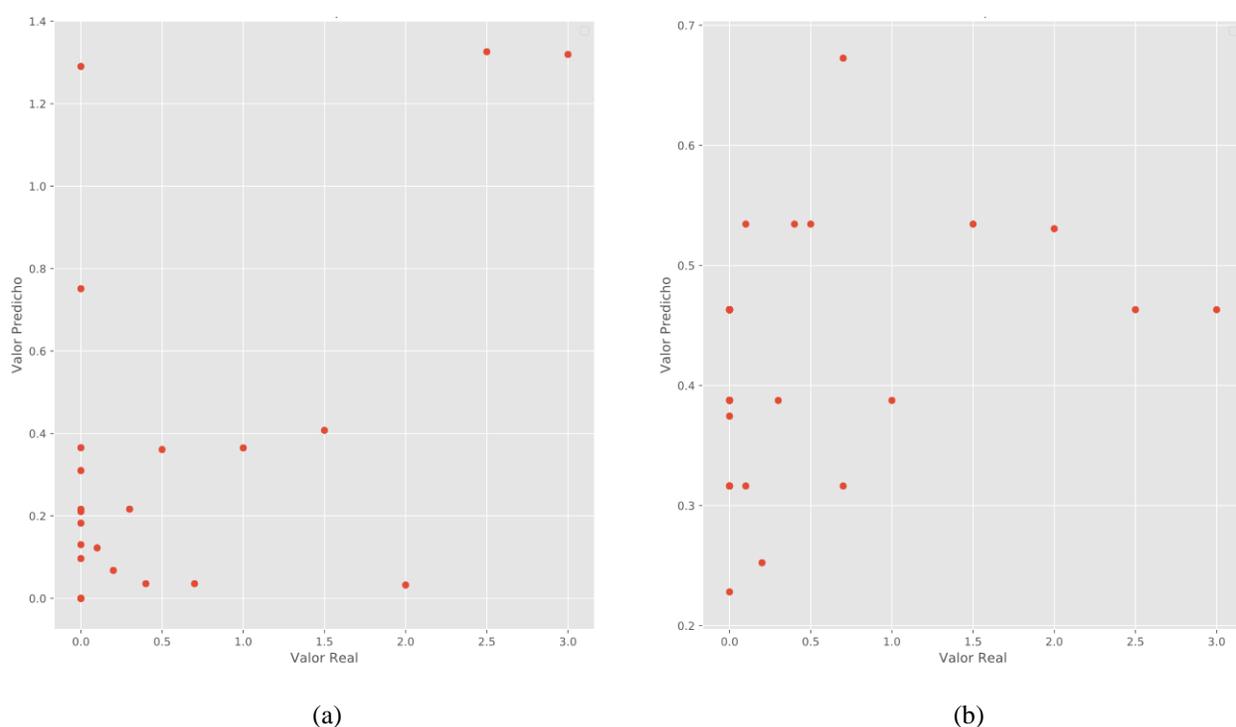
En base a los resultados obtenidos en los datos de entrenamiento y prueba, se determinó que los modelos de la base de datos 1 y 2, no se sobreajustan a los datos de entrenamiento ya que el resultado de R2 y MSE para el conjunto de entrenamiento es similar al obtenido en el conjunto de prueba. En la Figura 53a, se presenta la gráfica de los valores reales vs los valores predichos por el modelo de regresión de sustancias explosivas puras y mezclas, y en la Figura 53b del modelo de regresión de únicamente sustancias explosivas puras, ambos para la base de datos1. Estos modelos se encargan de predecir la concentración de sustancia explosiva sin importar el tipo de sustancia explosiva que es, por lo tanto, proporcionan una única salida (concentración predicha por el modelo) por cada observación. El modelo con el que se obtuvo mejores resultados fue aquel encargado de predecir la concentración de sustancias explosivas puras, esto puede verificar en las gráficas ya que los valores predichos se acercan más a los reales, sobre todo con las concentraciones más bajas y altas de sustancia.



**Figura 53.** Valores predichos por los modelos de regresión

(a) Regresión de sustancias explosivas puras y mezclas (b) Regresión de sustancias explosivas puras

En la Figura 54a, se presenta la gráfica de los valores reales y los valores predichos por el modelo de regresión de pólvora en estado puro, y en la Figura 54b del modelo de regresión de pólvora en estado puro y mezclas, ambos para la base de datos 1. El modelo con el que se obtuvo mejores resultados al igual que con los modelos analizados en la parte superior, fue aquel encargado de predecir la concentración de pólvora en estado puro, sin embargo, el desempeño no es bueno, ya que como se observa en las dos gráficas los valores predichos se alejan de los reales.

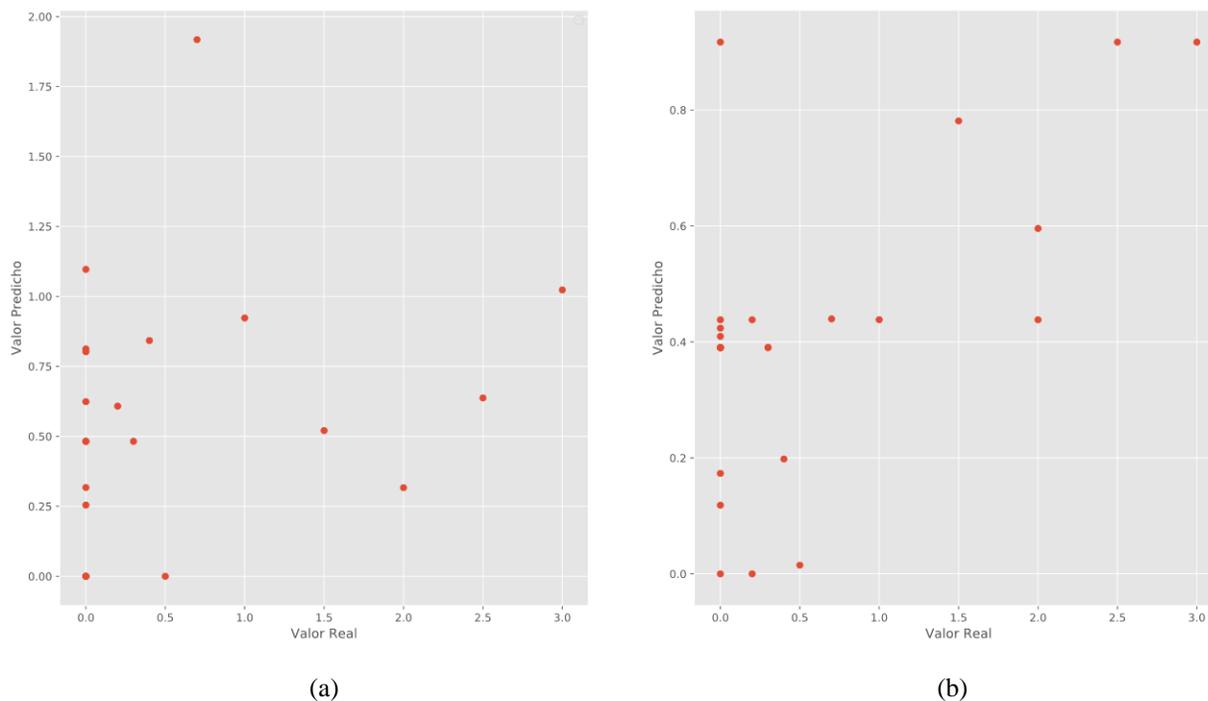


**Figura 54.** Valores predichos por los modelos de regresión

(a) Regresión de pólvora en estado puro (b) Regresión de pólvora en estado puro y mezclas

En la Figura 55a, se presenta la gráfica de los valores reales vs los valores predichos por el modelo de regresión de TNT en estado puro, y en la Figura 55b del modelo de regresión de TNT en estado puro y mezclas, ambos para la base de datos 1. El modelo con el que se obtuvo mejores resultados fue aquel encargado de predecir la concentración de TNT con mezclas. Sin embargo, su

desempeño no es bueno ( $R^2 \ll 0.4$  y  $MSE \gg 0$ ), ya que como se aprecia en las gráficas los valores de concentraciones predichas se alejan de los valores reales.

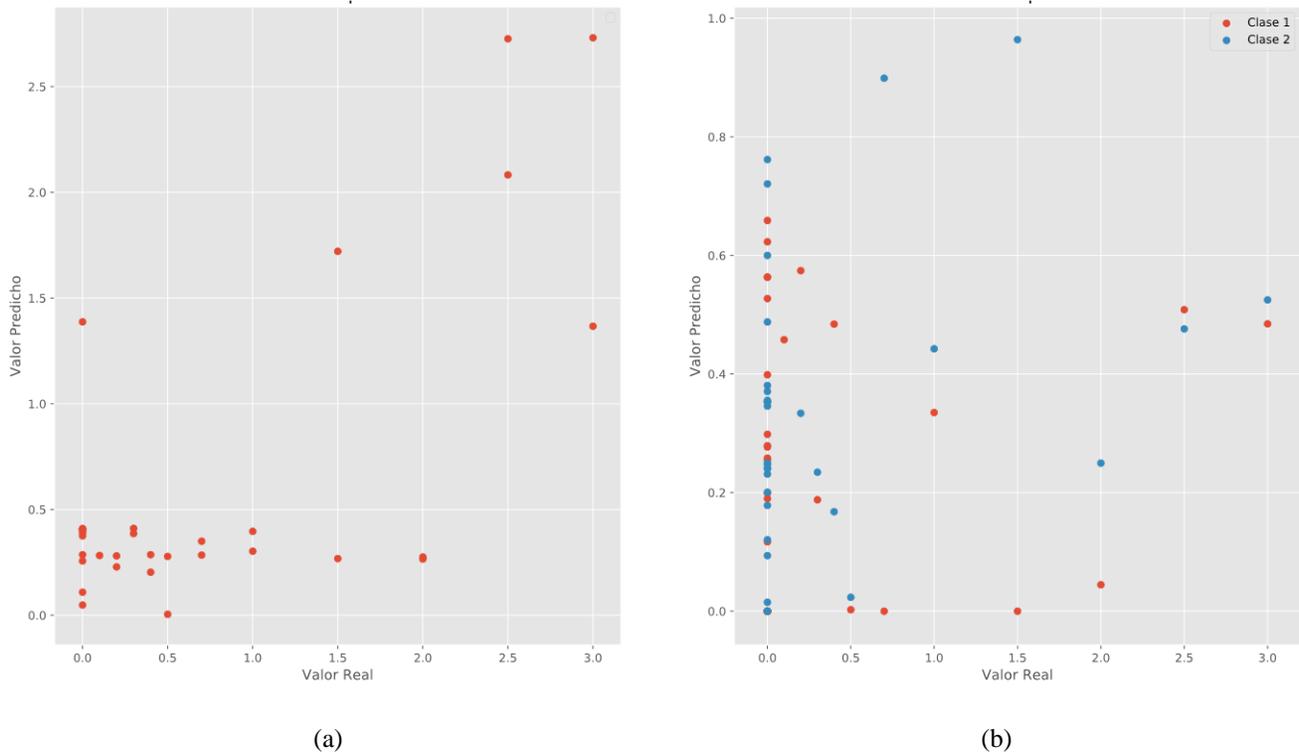


**Figura 55.** Valores predichos por los modelos de regresión

(a) Regresión de TNT en estado puro (b) Regresión de TNT en estado puro y mezclas

En la Figura 56a, se presenta la gráfica de los valores reales vs los valores predichos por el modelo de regresión de sustancias explosivas puras, y en la Figura 56b la gráfica del modelo de regresión de TNT y pólvora en estado puro, que a diferencia del primero este si cuenta con dos salidas: uno para la concentración predicha de la clase 1 correspondiente a pólvora y otro para la clase 2 correspondiente a TNT, ambos para la base de datos 1. De estos resultados se determinó que el prototipo bajo las condiciones iniciales en que se encontraba y con bajas concentraciones de sustancia explosiva, era capaz de predecir mejor la concentración de las sustancias con un modelo en el que no se toma en cuenta si la observación que ingresa al modelo pertenece a una clase de

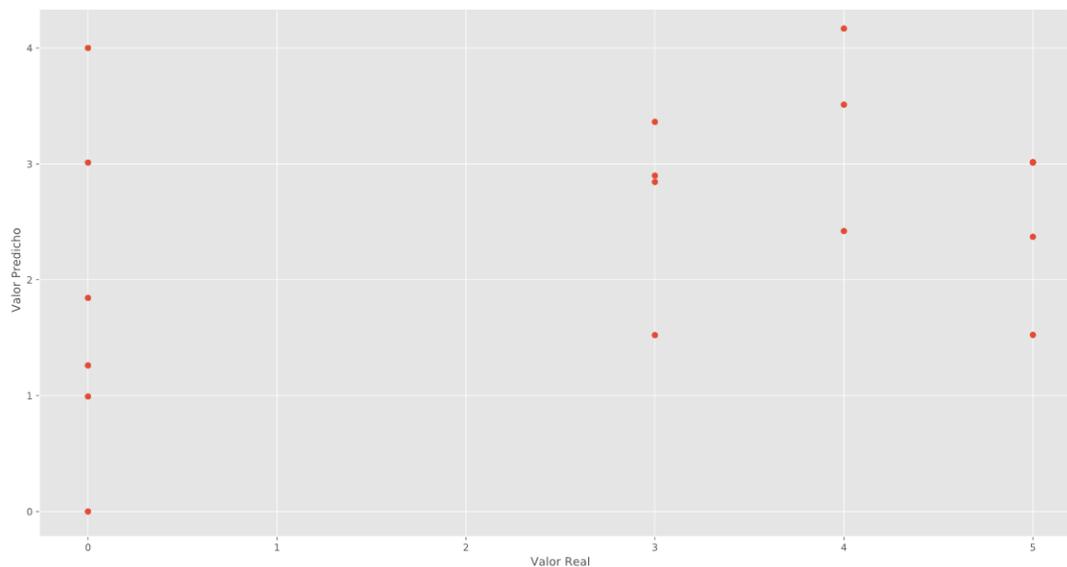
explosivo u otra que con un modelo de dos clases que predice la concentración por cada clase de sustancia.



**Figura 56.** Valores predichos por los modelos de regresión

(a) Regresión de sustancias explosivas puras (b) Regresión de TNT y pólvora en estado puro

Con relación a los modelos de la base de datos 2, el prototipo tuvo un mal desempeño en todos los casos ( $R^2 \ll 0.4$  y  $MSE \gg 0$ ), como se observa en la Figura 57 de los valores reales vs los valores predichos por el modelo de regresión de sustancias explosivas puras. Las observaciones sin sustancia explosiva son incorrectamente identificadas con valores de hasta 4 gr y en el caso de las sustancias con el más alto nivel de concentración: 5gr, los resultados predichos varían, pero no bajan de 1gr.



**Figura 57.** Valores predichos por el modelo de regresión de sustancias explosivas puras para la base de datos 2

En el presente capítulo se describieron los conceptos utilizados para el desarrollo de los modelos de clasificación y de regresión con redes neuronales MLP con una capa oculta. Posteriormente, se describió el proceso para la creación de los modelos en el que se definió el tipo de escalamiento usado y cuál fue el uso que se dio a la técnica de validación cruzada. Finalmente se presentó una descripción de los resultados obtenidos para cada una de las tareas.

# CAPÍTULO 7

## Deep Learning

En este capítulo se presenta una descripción de los conceptos utilizados en el desarrollo de los modelos de clasificación y regresión mediante la técnica de deep learning que consiste una red neuronal profunda Long Short-Term Memory (LSTM) importante para el aprendizaje profundo de series de tiempo. A continuación, se presenta el procedimiento para la elaboración de los modelos, en el que se describe el tipo de tratamiento previo al ingreso de los datos a la red neuronal profunda y los hiperparámetros seleccionados durante la validación cruzada. Por último, se muestran los resultados obtenidos con los modelos desarrollados en base al desempeño evaluado con las métricas mencionadas en el Capítulo 2.

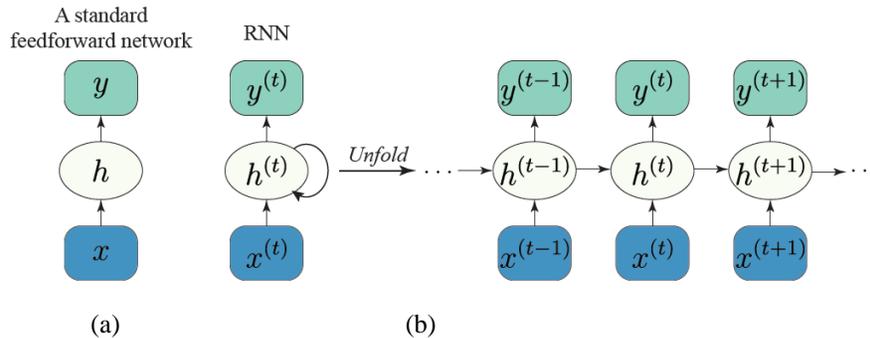
### 7.1. Conceptos Básicos

Las redes neuronales clásicas denominadas perceptrón multicapa (MLP), presentadas en el Capítulo 6, son un buen punto de inicio para modelar problemas de clasificación y regresión de series temporales, sin embargo, cuentan con ciertas limitaciones ya que los pasos de tiempo se modelan como características de los datos de entrada, por lo que la red no tiene una comprensión explícita del orden o estructura temporal entre muestras (Jason Brownlee, 2017).

#### 7.1.1. Red Neuronal Recurrente (RNN)

Una RNN, se diferencia de una red neuronal clásica por el tipo de neuronas que posee en las capas ocultas. En este tipo de red las entradas de sus neuronas en las capas ocultas no son únicamente

las salidas de la capa anterior, sino a su vez, la salida de la neurona para un tiempo anterior, como se muestra en la Figura 58a, con ello se dice que la red obtiene una memoria de los eventos pasados (Raschka & Mirjalili, 2017). A estas neuronas se las conoce con el nombre de neuronas recurrentes y pueden representarse de forma desplegada a través del tiempo como se observa en la Figura 58b.



**Figura 58.** Arquitectura de una red neuronal clásica y de una red neuronal recurrente

Fuente: (Raschka & Mirjalili, 2017)

Cada neurona recurrente tiene un conjunto de pesos para las entradas  $x$  y otro para las salidas de un tiempo anterior  $y_{(t-1)}$  (Geron, 2017). Por lo cual, la salida de una neurona recurrente puede representarse mediante la suma ponderada de la entrada y de la salida para un tiempo anterior, como se muestra en la Ecuación 55.

$$y = \phi(x_{(t)}w_x + y_{(t-1)}w_y + b) \quad ( 55 )$$

Donde  $\phi$  representa la función de activación de la neurona,  $b$  el término de bias,  $x_{(t)}$  la entrada actual de la neurona y  $w_x$  el peso correspondiente a la entrada,  $y_{(t-1)}$  la salida de la neurona para un tiempo anterior y  $w_y$  su peso respectivo.

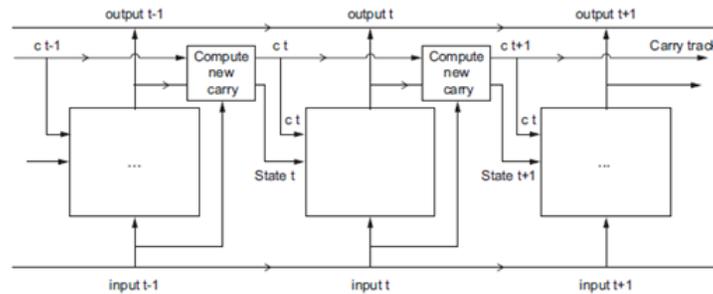
El entrenamiento de las RNN se realiza mediante retropropagación a través del tiempo (BPTT Backpropagation through time), el cual consiste en desplegar las neuronas a través del tiempo, como se mostró en la **Figura 58b** y luego utilizar la retropropagación regular definida en el

Capítulo 6. Al igual que en la retropropagación regular, la función de costo y los gradientes de la función se propagan hacia atrás a través de la red desplegada, y por último los pesos del modelo se actualizan. Sin embargo, existen ciertos problemas que aparecen durante la etapa de entrenamiento, uno de ellos es que las transformaciones por las que pasan los datos al pasar por una RNN provocan que se pierda cierta información después de cada paso de tiempo, otro de los problemas que surgen con este tipo de neuronas es que al desplegar las neuronas a través del tiempo para series temporales de larga duración se obtendrá una red muy profunda que puede sufrir de problemas de desvanecimiento o explosión de gradientes provocando que el entrenamiento sea lento (Geron, 2017). Para resolver estos problemas se han creado otro tipo de neuronas recurrentes que cuentan con memoria a largo plazo, como lo es la memoria larga a corto plazo (LSTM).

### **7.1.2. Red Neuronal de Memoria Larga a Corto Plazo (LSTM)**

La red neuronal de memoria a largo plazo (LSTM), es un tipo de red neuronal recurrente (RNN) que gracias a su diseño evita el problema de desvanecimiento/explosión de gradiente que se presenta en una RNN simple. Una neurona LSTM, como se muestra en la Figura 59, cuenta con una entrada y una salida por cada paso de tiempo y dos estados: un estado a largo plazo  $c$  y un estado actual o a corto plazo  $h$ . Internamente cuenta con tres compuertas: una de entrada que se encarga de decidir qué valores de entrada pueden ingresar al estado a largo plazo  $c$ , una de salida que decide la salida de la neurona y por ende el estado a corto plazo en base a la entrada y al estado a largo plazo, y una de olvido que decide qué información se debe descartar del estado a largo plazo (Jason Brownlee, 2017). Con estas compuertas la neurona podrá aprender a reconocer entradas importantes, almacenarlas en el estado a largo plazo, preservarlas durante el tiempo necesario y extraerlas cuando corresponda (Geron, 2017). El entrenamiento de una red LSTM consiste en la aplicación del método de retropropagación a través del tiempo (BPTT), el cual como se mencionó

anteriormente consiste en desplegar las neuronas a través del tiempo y luego utilizar el método de retropropagación regular.



**Figura 59.** Estructura de una neurona LSTM

Fuente: (Chollet, 2018)

Para la construcción de los modelos de aprendizaje profundo LSTM se definió una única etapa principal: la construcción del modelo de predicción, ya que este tipo de red podrá aprender de los datos en bruto sin requerir realizar manualmente de una extracción de las características de las observaciones. Las etapas para su construcción se muestran en la Figura 60 y su descripción en las subsecciones siguientes.



**Figura 60.** Proceso para la construcción de modelo de deep learning LSTM

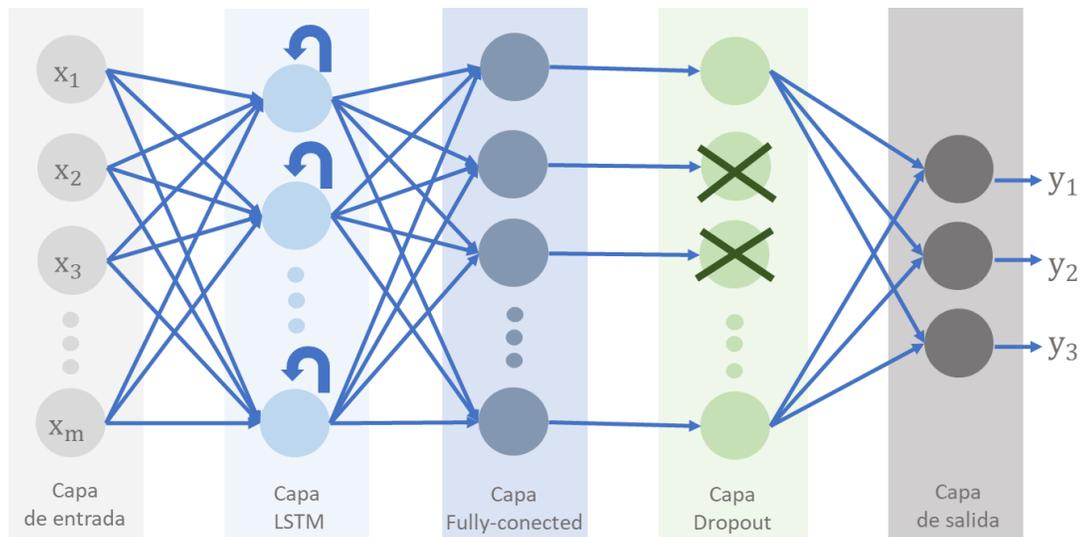
El proceso de construcción del modelo de deep learning consta de tres etapas: El ingreso de los datos de entrada y salida, la construcción del modelo de predicción y la obtención de las salidas predichas.

### 7.1.3. Construcción del Modelo de Deep Learning LSTM

Para la construcción de la red neuronal profunda LSTM se deberá definir la arquitectura de la red, las funciones de activación de las neuronas de las capas ocultas y de salida, la función objetivo que va a ser minimizada durante el proceso de aprendizaje y la función de optimización que minimizará la función objetivo, cada uno de estos parámetros se describen a continuación.

#### 7.1.3.1. Arquitectura de la Red

Para el modelo de deep learning se seleccionaron dos capas ocultas: una capa LSTM que extraerá las características de las series temporales, una capa fully-connected o completamente conectada encargada de interpretar las características extraídas por la capa LSTM y finalmente, una capa dropout o de abandono encargada de eliminar temporalmente un porcentaje de neuronas de la capa anterior de forma randómica para prevenir el sobreajuste del modelo a los datos de entrenamiento, cada una de las partes que conforman la red se muestran en la Figura 61.



**Figura 61.** Arquitectura de la red neuronal de deep learning LSTM

La grafica muestra la red de deep learning LSTM con dos capas oculta seleccionada para el desarrollo de los modelos de clasificación y cuantificación.

### 7.1.3.2. Funciones de Activación

Como se mencionó en el Capítulo 6, existen diversos tipos de funciones de activación para las capas ocultas, sin embargo, se decido usar la función relu de la Ecuación 56 en las capas ocultas ya que es la más usada actualmente por su bajo costo computacional y a que no existen problemas de desvanecimiento de gradiente para valores positivos extremos de entrada (Raschka & Mirjalili, 2017).

#### Función ReLU

$$\phi(z) = \begin{cases} z, & z \geq 0 \\ 0, & z < 0 \end{cases} \quad ( 56 )$$

Para la capa de salida del modelo de clasificación se utilizó la función de activación softmax, definida mediante la Ecuación 57, que calcula la probabilidad de que la entrada  $z$  pertenezca la clase  $i$ , donde  $k$  representa el número de clases. Para el modelo de regresión la función de activación elegida fue la función identidad de la Ecuación 58.

#### Función softmax

$$\phi(z) = \frac{e^{z_i}}{\sum_{j=1}^p e^{z_j}} \quad ( 57 )$$

#### Función identidad

$$\phi(z) = z \quad ( 58 )$$

### 7.1.3.3. Función Objetivo

Como se mencionó en el Capítulo 6, es necesario definir una función objetivo que mida el error entre las predicciones realizadas con la red y el objetivo verdadero para utilizarla como señal de retroalimentación para el ajuste del valor de los pesos durante la etapa de entrenamiento (Chollet,

2018). Ya que la selección de la función objetivo depende de la tarea del modelo, para los modelos de regresión se seleccionó el error cuadrático medio (MSE) presentado en la Ecuación 59, el cual calcula el promedio del cuadrado de la diferencia entre la salida predicha  $\hat{y}$  y el objetivo  $y$  (Rao & McMahan, 2019).

$$J_{MSE}(y, \hat{y}) = \frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2 \quad ( 59 )$$

Para los modelos de clasificación la entropía categórica cruzada presentada en la Ecuación 60, con la cual se calcula el promedio del producto del logaritmo negativo de la salida predicha  $\hat{y}$  con el objetivo  $y$  (Rao & McMahan, 2019).

$$J(\theta) = -\frac{1}{n} \sum_{i=0}^n \sum_{j=0}^m y_{(i,j)} \times \log(\hat{y}_{(i,j)}) \quad ( 60 )$$

#### 7.1.3.4. Función de Optimización

La función de optimización es aquella que implementa el algoritmo de retropropagación para ajustar los pesos del modelo según la información proporcionada por la función objetivo (Chollet, 2018). Como se definió en el Capítulo 6, existen diversos tipos de funciones de optimización para redes neuronales, como lo son el descenso de gradiente (GD) y el descenso de gradiente estocástico (SGD), siendo el último el seleccionado como función de optimización. El proceso para el desarrollo de los modelos se detalla en la siguiente sección.

## 7.2. Generación de Modelos de Deep Learning LSTM

El proceso de desarrollo de los modelos de deep learning, como se mencionó en la sección anterior consiste en la construcción del modelo de predicción para lo cual se empleó la librería

*Keras*, la cual permite definir las capas de la red neuronal en forma secuencial. Finalmente, con la arquitectura de la red definida se procedió a entrenar los modelos de clasificación y regresión con los hiperparámetros seleccionados durante la validación cruzada.

### **7.2.1. Escalamiento**

Al igual que en los capítulos anteriores se seleccionó los métodos de escalamiento empleados en el Capítulo 5. Los cuales centran las características de los datos en la media y los escala con una varianza igual a 1 para tareas de clasificación y únicamente los centra en la media para tareas de regresión.

### **7.2.2. Balanceo de Clases**

Al tener un mayor número de observaciones de una clase que de otra, el algoritmo de aprendizaje aprenderá implícitamente un modelo que optimiza las predicciones basadas en la clase más abundante en el conjunto de datos (Raschka & Mirjalili, 2017). Una forma de solucionar este problema es balancear la cantidad de observaciones en cada clase ya sea mediante una técnica de sobre muestro o un submuestreo de los datos. En este caso se utilizó el algoritmo *SMOTE*, el cual sobre muestrea la clase minoritaria creando observaciones sintéticas (Chawla et al., 2002).

### **7.2.3. Validación Cruzada**

Se utilizó la validación cruzada de  $k=10$  iteraciones para seleccionar el número de neuronas en cada una de las capas ocultas, la tasa de aprendizaje para la función de optimización y el número de épocas para el entrenamiento del modelo. Con los hiperparámetros seleccionados se entrenó y evaluó los modelos con los datos de entrenamiento y validación, respectivamente. Los resultados

del desempeño de los modelos en los datos de entrenamiento y prueba se presentan en la siguiente sección.

## 7.3. Resultados del Desempeño de los Modelos de Deep Learning LSTM

### 7.3.1. Modelos de Clasificación

En la **Tabla 24** y **Tabla 25**, se presentan los resultados obtenidos con los modelos de clasificación para la base de datos 1 y 2, en los cuales se describe el número de neuronas de la red en las capas ocultas, la tasa de aprendizaje y el número de épocas, además del AUC de los datos de entrenamiento y prueba para la evaluación desempeño de los modelos.

**Tabla 24**

*Base de datos 1- Resultados de los modelos de clasificación de deep learning*

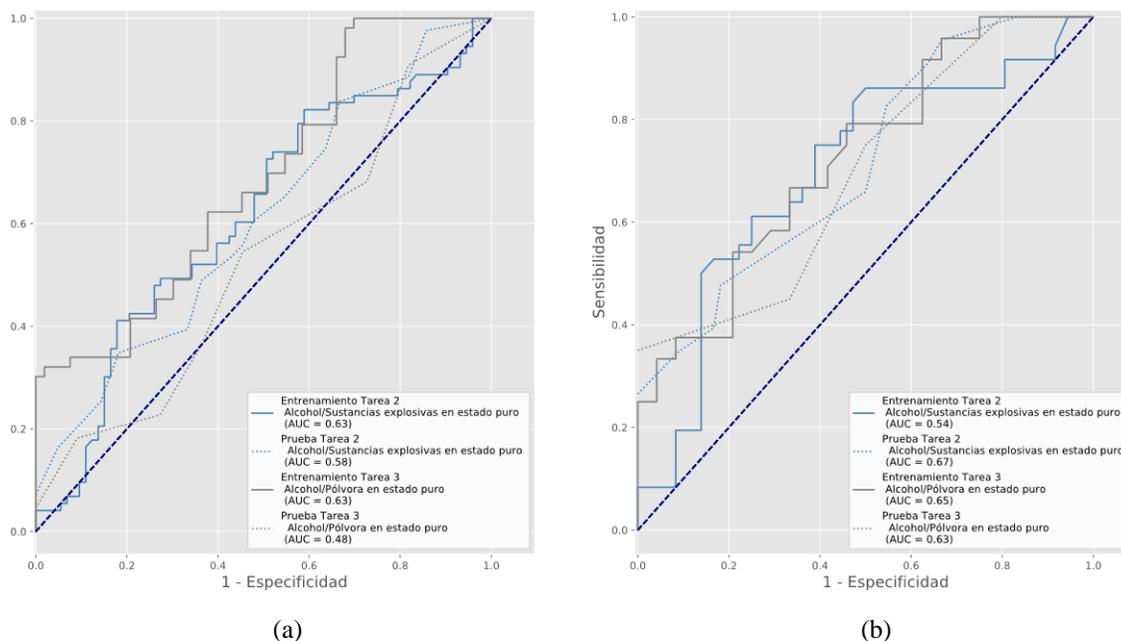
N°	TAREA	#NEURONAS		TASA DE APRENDIZAJE	ÉPOCAS	AUC	
		CAPA LSTM	CAPA FULLY CONNECTED			TRAIN	TEST
1	Alcohol/Sustancias explosivas en estado puro y mezclas	100	100	0.01	20	0.66	0.63
2	Alcohol/Sustancias explosivas en estado puro	100	100	0.01	20	0.63	0.58
3	Alcohol/Pólvora en estado puro	10	10	0.1	20	0.63	0.48
4	Alcohol/TNT en estado puro	10	10	0.1	20	0.6	0.68
5	Pólvora en estado puro/TNT en estado puro	10	10	0.1	20	0.51	0.67
6	Alcohol/Pólvora en estado puro y mezcla de pólvora	10	10	0.1	20	0.59	0.55
7	Alcohol/TNT en estado puro y mezcla de TNT	10	10	0.1	20	0.66	0.73
8	Alcohol/Pólvora en estado puro/TNT en estado puro	100	100	0.01	10	0.53	0.52

**Tabla 25**

Base de datos 2- Resultados de los modelos de clasificación de deep learning

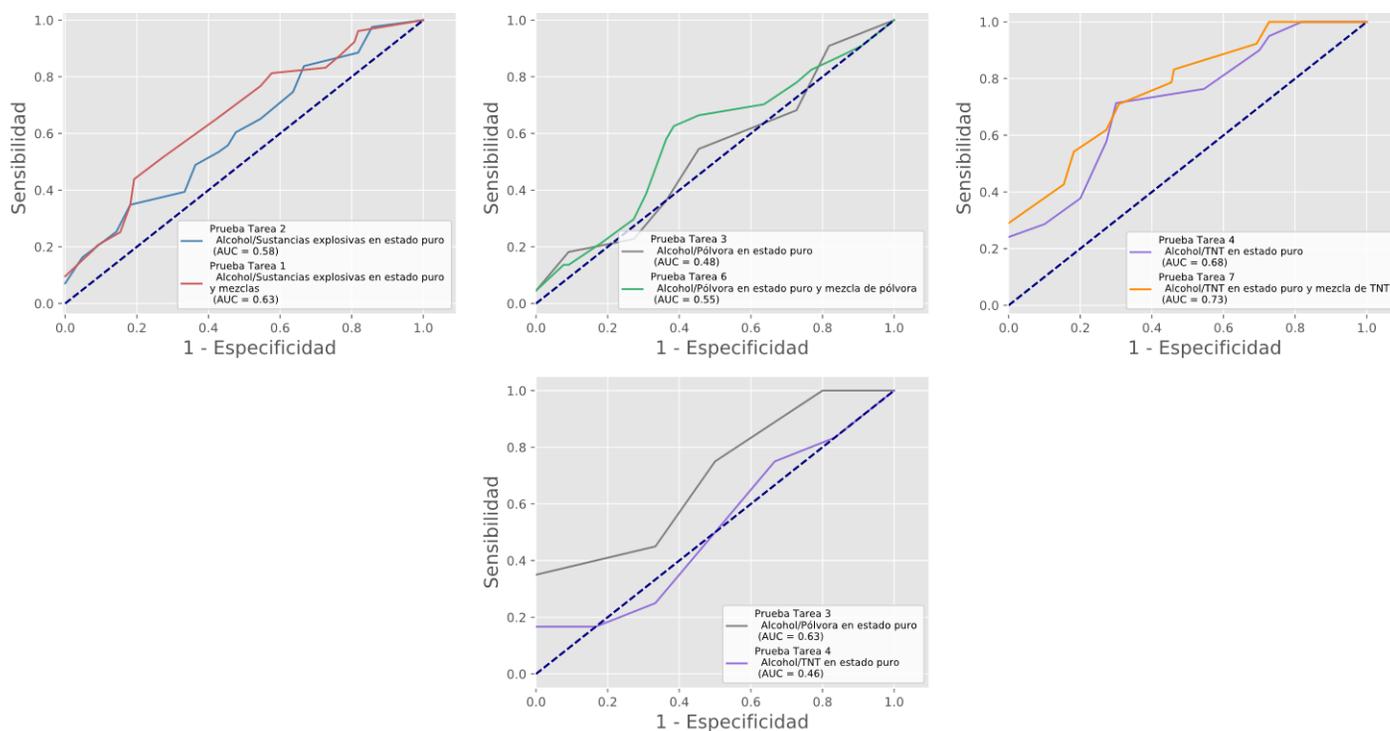
N°	TAREA	#NEURONAS		TASA DE APRENDIZAJE	ÉPOCAS	AUC	
		CAPA LSTM	CAPA FULLY CONNECTED			TRAIN	TEST
2	Alcohol/Sustancias explosivas en estado puro	100	100	0.01	10	0.54	0.67
3	Alcohol/Pólvora en estado puro	50	50	0.1	10	0.65	0.63
4	Alcohol/TNT en estado puro	10	10	0.1	20	0.67	0.46
5	Pólvora en estado puro/TNT en estado puro	10	10	0.1	20	0.72	0.47
8	Alcohol/Pólvora en estado puro/TNT en estado puro	10	10	0.1	20	0.64	0.53

En la Figura 62, se presentan las curvas ROC y el AUC de dos de los modelos de la base de datos 1 y 2, en los cuales se puede apreciar que contrario a los modelos generados con los algoritmos de capítulos anteriores, el desempeño de los modelos en los datos de entrenamiento es similar al de los datos de prueba, por lo tanto, se puede asumir que los modelos no se sobreajustan a los datos de entrenamiento.



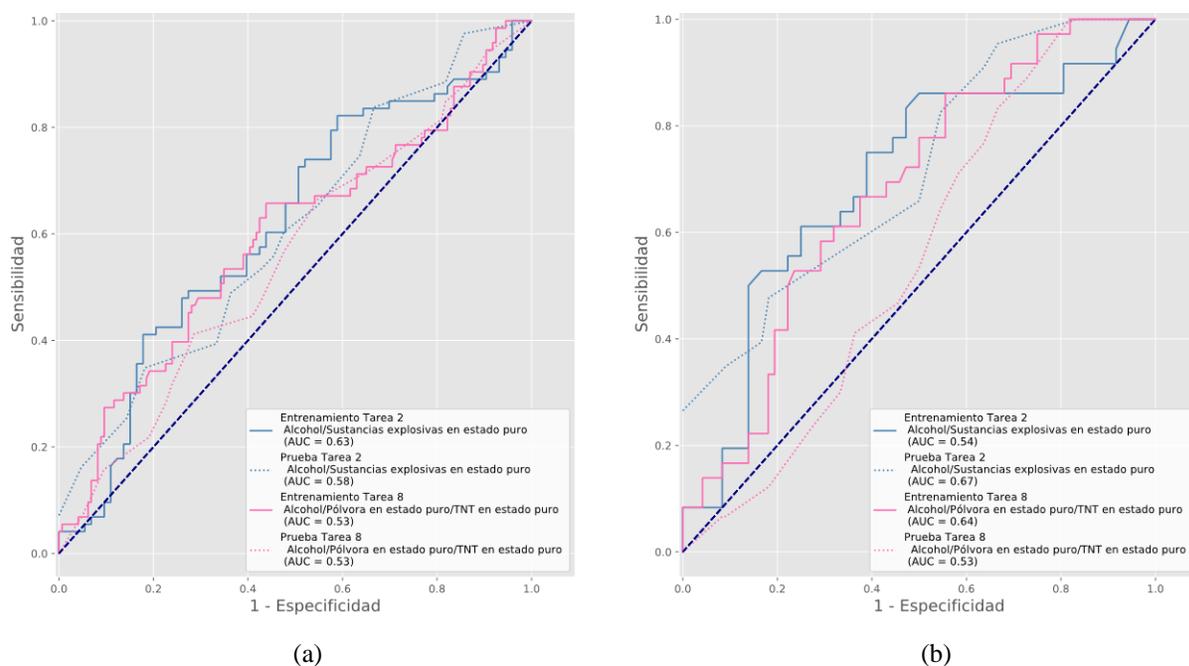
**Figura 62.** Curvas ROC para la descripción del desempeño de los modelos en datos de entrenamiento y prueba

Respecto al desempeño de los modelos de la base de datos 1 para clasificación de sustancias explosivas puras y con mezclas este es superior al de los modelos de clasificación de únicamente sustancias explosivas puras, como se puede observar en los resultados del AUC y gráficamente en las curvas ROC de la Figura 63. Además, como se muestra en la Figura 63b y Figura 63c, el desempeño de los modelos de clasificación de TNT es superior a los de clasificación de pólvora. En el caso de la base de datos 2, como se muestra en la Figura 63d, el desempeño del modelo de clasificación de pólvora es superior al de TNT. Por lo tanto, se espera que bajo las condiciones actuales del prototipo y con concentraciones entre 3 y 5gr de sustancias explosivas este sea capaz de clasificar correctamente en mayor porcentaje observaciones de pólvora que de TNT con los modelos de deep learning.



**Figura 63.** Curvas ROC para descripción del desempeño de los modelos para clasificación de sustancias explosivas puras y mezclas

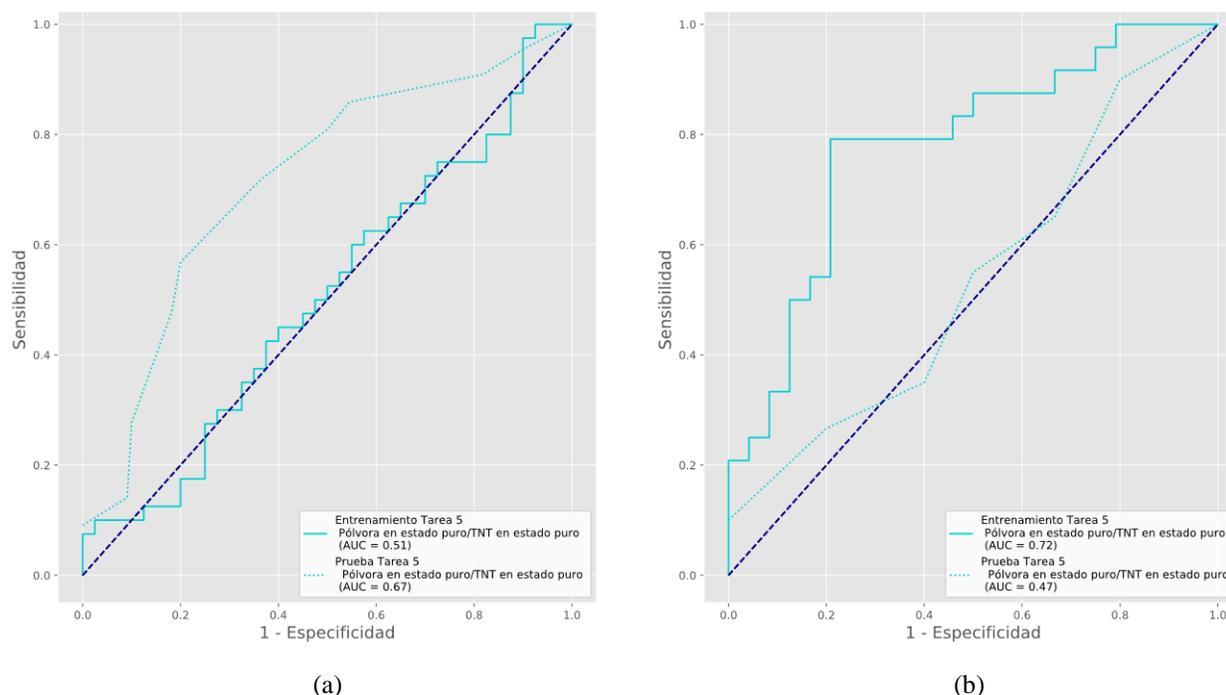
En cuanto a la diferencia entre un modelo que determina si una observación es sustancia explosiva o no y uno que determina si una observación no es sustancia explosiva y si lo es el tipo entre TNT y pólvora. Se pudo observar, como se muestra en la Figura 64 para la base de datos 1 y 2, que el desempeño de un modelo que únicamente determina si una sustancia es o no fue superior al modelo que identifica además el tipo de sustancia explosiva. Además, al comparar el desempeño entre los modelos realizados con bajas concentraciones de sustancia explosiva (experimentos de la base de datos 1 con concentraciones entre 0.1 y 3gr de pólvora y TNT) y de los realizados con mayores concentraciones (experimentos de la base de datos 2 con concentraciones entre 3 y 5gr de pólvora y TNT), se pudo observar que se obtuvo mejores resultados al aumentar el nivel de concentración de sustancia.



**Figura 64.** Curvas ROC para descripción del desempeño de los modelos de clasificación entre dos y tres clases de sustancias

Finalmente, al analizar los resultados de los modelos de clasificación entre TNT y pólvora de la Figura 65a para la base de datos 1 y Figura 65b para la base de datos 2, se determinó que el

prototipo e-nose es capaz de clasificar pólvora y TNT ( $AUC > 0.5$ ), contrario al modelo de la base de datos 2 para clasificación de sustancias con concentraciones entre 3 y 5gr cuyo  $AUC < 0.5$ . Además, en las figuras se identifica que el resultado de la clasificación de los datos de entrenamiento es bajo en comparación a los de prueba, es decir, los modelos se subajustan a los datos de entrenamiento.



**Figura 65.** Curvas ROC para descripción del desempeño de los modelos de clasificación entre pólvora y TNT

### 7.3.2. Modelos de Regresión

En la **Tabla 26** y **Tabla 27**, se presentan los resultados de los modelos de regresión para la base de datos 1 y 2, respectivamente, en los cuales se describe el número de neuronas de la red, la tasa de aprendizaje y el número de épocas, además del MSE y R2 de los datos de entrenamiento y prueba para la evaluación desempeño de los modelos. Con estos resultados se analizó el desempeño de los modelos de regresión con redes neuronales profundas LSTM.

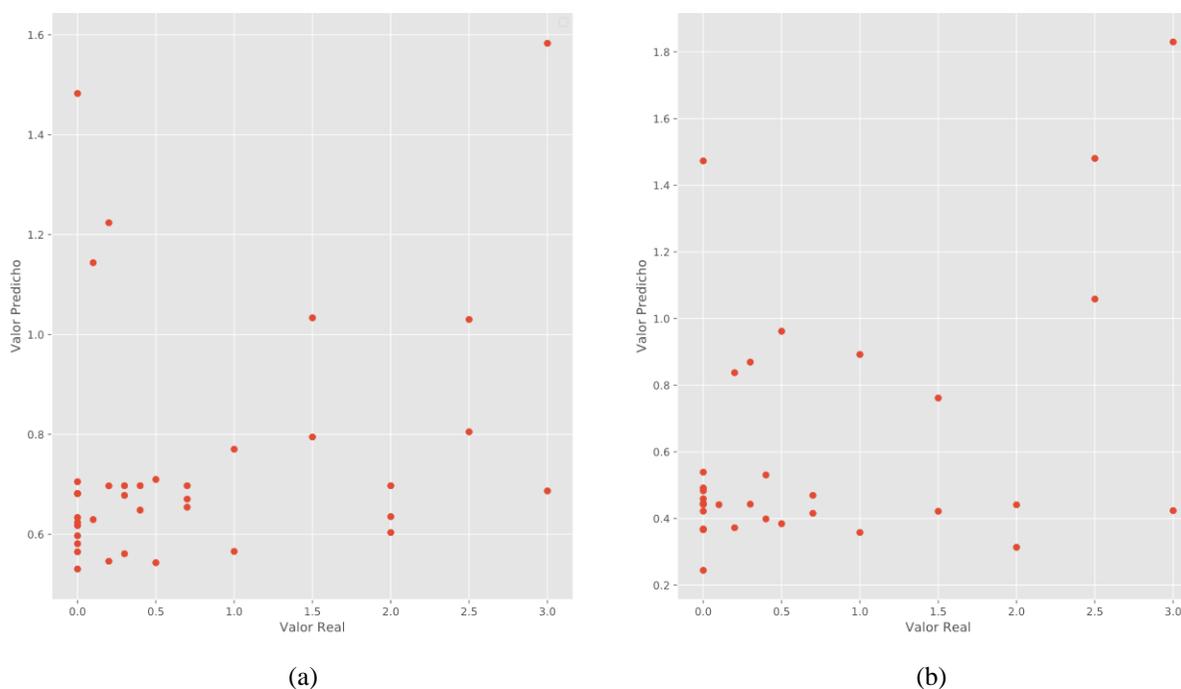
**Tabla 26***Base de datos 1- Resultados de los modelos de regresión de deep learning*

N°	TAREA	#Neuronas		Tasa de aprendizaje	Epocas	MSE		R2	
		Capa LSTM	Capa Fully Connected			TRAIN	TEST	TRAIN	TEST
1	Alcohol/Sustancias explosivas en estado puro y mezclas	10	10	0.1	20	0.67	0.77	0.10	0.10
2	Alcohol/Sustancias explosivas en estado puro	100	100	0.01	50	0.64	0.75	0.21	0.17
3	Alcohol/Pólvora en estado puro	10	10	0.1	20	0.57	0.75	0.01	0.01
4	Alcohol/TNT en estado puro	10	10	0.1	20	0.65	0.71	0.04	0.10
5	Pólvora en estado puro/TNT en estado puro	10	10	0.1	20	0.92	0.87	0.01	0.03
6	Alcohol/Pólvora en estado puro y mezcla de pólvora	10	10	0.1	20	0.51	0.77	0.05	-0.09
7	Alcohol/TNT en estado puro y mezcla de TNT	10	10	0.1	20	0.64	0.77	0.01	0.01
8	Alcohol/Pólvora en estado puro/TNT en estado puro	10	10	0.1	20	0.48	0.62	0.01	-0.04

**Tabla 27***Base de datos 2- Resultados de los modelos de regresión de deep learning*

N°	TAREA	#Neuronas		Tasa de aprendizaje	Epocas	MSE		R2	
		Capa LSTM	Capa Fully Connected			TRAIN	TEST	TRAIN	TEST
2	Alcohol/Pólvora en estado puro/TNT en estado puro	10	10	0.1	20	4.09	4.36	-0.05	-0.06
3	Alcohol/Pólvora en estado puro	10	10	0.1	20	5.52	5.41	-0.34	-0.25
4	Alcohol/TNT en estado puro	10	10	0.1	20	4.69	4.69	-0.08	-0.08
5	Pólvora en estado puro/TNT en estado puro	10	10	0.1	20	0.64	0.74	-0.03	-0.02
8	Alcohol/Pólvora en estado puro/TNT en estado puro	10	10	0.1	20	5.28	5.27	-0.43	-0.41

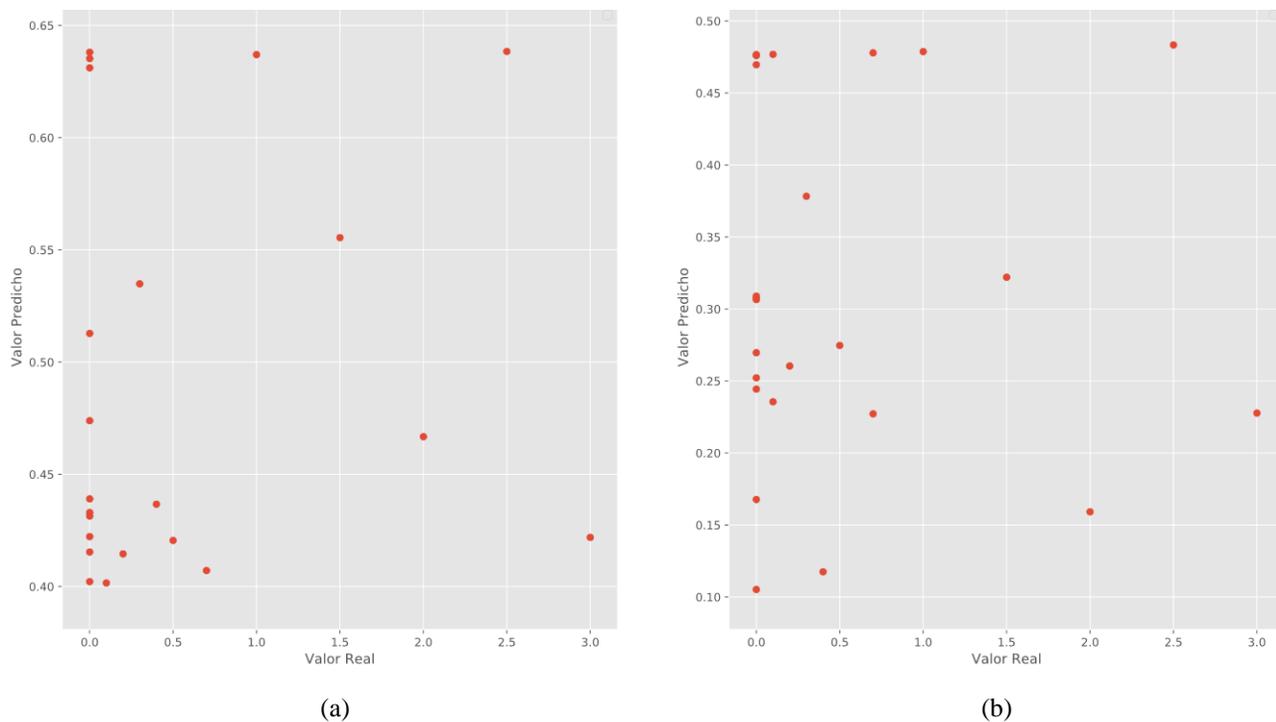
En base a los resultados obtenidos en los datos de entrenamiento y prueba, se determinó que los modelos de la base de datos 1 y 2, no se sobreajustan a los datos de entrenamiento ya que el resultado de R2 y MSE para el conjunto de entrenamiento es similar al obtenido en el conjunto de prueba. En la Figura 66a, se presenta la gráfica de los valores reales vs los valores predichos por el modelo de regresión de sustancias explosivas puras y mezclas, y en la Figura 66b del modelo de regresión de únicamente sustancias explosivas puras, ambos para la base de datos 1. Estos modelos se encargan de predecir la concentración de sustancia explosiva sin importar el tipo de sustancia explosiva que es, por lo tanto, proporcionan una única salida (concentración predicha por el modelo) por cada observación. El modelo con el que se obtuvo mejores resultados fue aquel encargado de predecir la concentración de sustancias explosivas puras, esto puede verificar en las gráficas ya que los valores predichos se acercan más a los reales, sobre todo con las concentraciones más bajas de sustancia.



**Figura 66.** Valores predichos por los modelos de regresión

(a) Regresión de sustancias explosivas puras y mezclas (b) Regresión de sustancias explosivas puras

En la Figura 67a, se presenta la gráfica de los valores reales vs los valores predichos por el modelo de regresión de pólvora en estado puro, y en la Figura 67b del modelo de regresión de pólvora en estado puro y mezclas, ambos para la base de datos 1. El modelo con el que se obtuvo mejores resultados fue aquel encargado de predecir la concentración de pólvora en estado puro, sin embargo, el desempeño no es bueno, ya que como se observa en las dos gráficas los valores predichos se alejan de los reales.

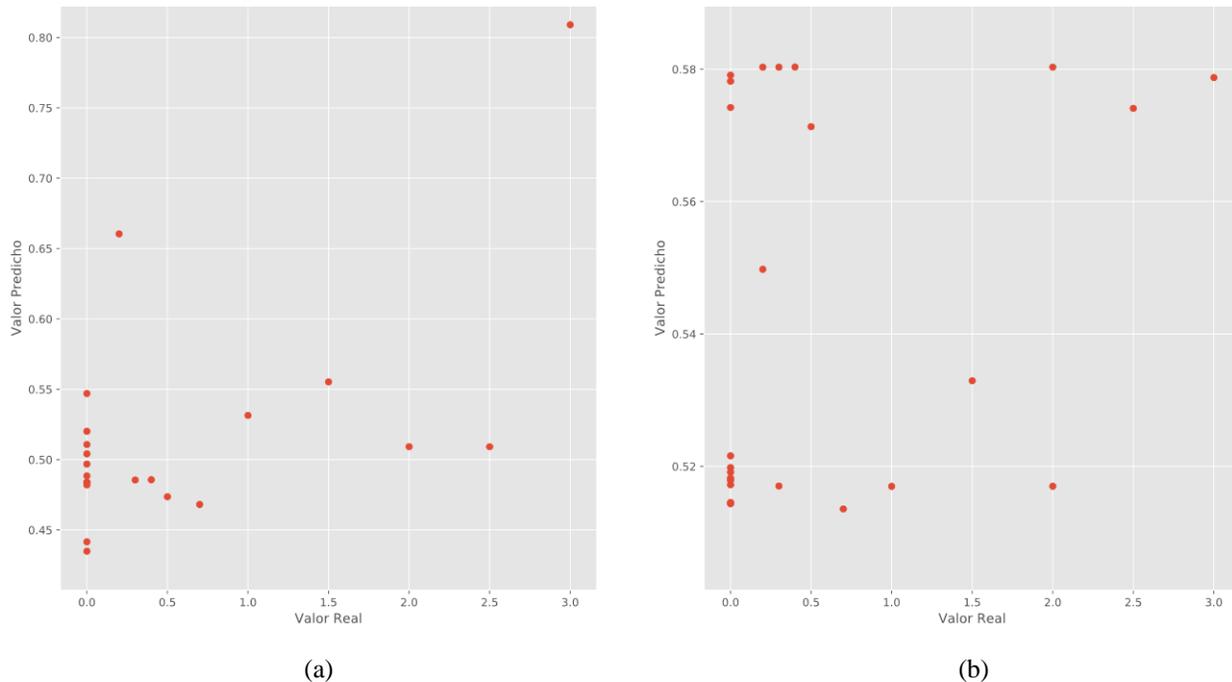


**Figura 67.** Valores predichos por los modelos de regresión

(a) Regresión de pólvora en estado puro (b) Regresión de pólvora en estado puro y mezclas

En la Figura 68a, se presenta la gráfica de los valores reales vs los valores predichos por el modelo de regresión de TNT en estado puro, y en la Figura 68b del modelo de regresión de TNT en estado puro y mezclas, ambos para la base de datos 1. El modelo con el que se obtuvo mejores resultados fue aquel encargado de predecir la concentración de TNT en estado puro. Sin embargo,

su desempeño no es bueno ( $R^2 \ll 0.4$  y  $MSE \gg 0$ ), ya que como se aprecia en las gráficas los valores de concentraciones predichas se alejan de los valores reales.

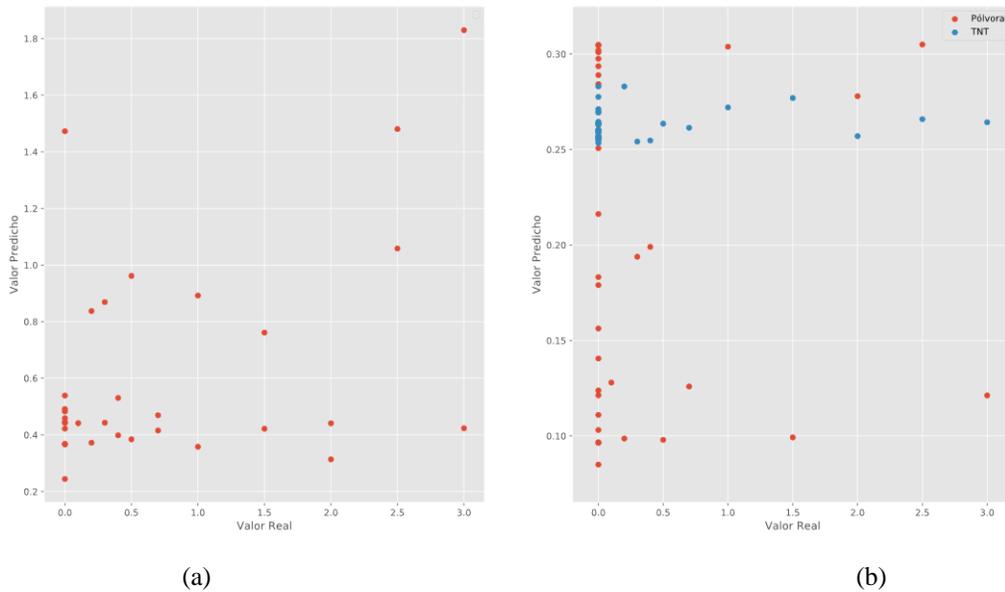


**Figura 68.** Valores predichos por los modelos de regresión

(a) Regresión de TNT en estado puro (b) Regresión de TNT en estado puro y mezclas

En la Figura 69a, se presenta la gráfica de los valores reales vs los valores predichos por el modelo de regresión de sustancias explosivas puras, y en la Figura 69b la gráfica del modelo de regresión de TNT y pólvora en estado puro, que a diferencia del primero este si cuenta con dos salidas: uno para la concentración predicha de la clase 1 correspondiente a pólvora y otro para la clase 2 correspondiente a TNT, ambos para la base de datos 1. De estos resultados se determinó que el prototipo bajo las condiciones iniciales en que se encontraba y con bajas concentraciones de sustancia explosiva, era capaz de predecir mejor la concentración de las sustancias con un modelo en el que no se toma en cuenta si la observación que ingresa al modelo pertenece a una clase de

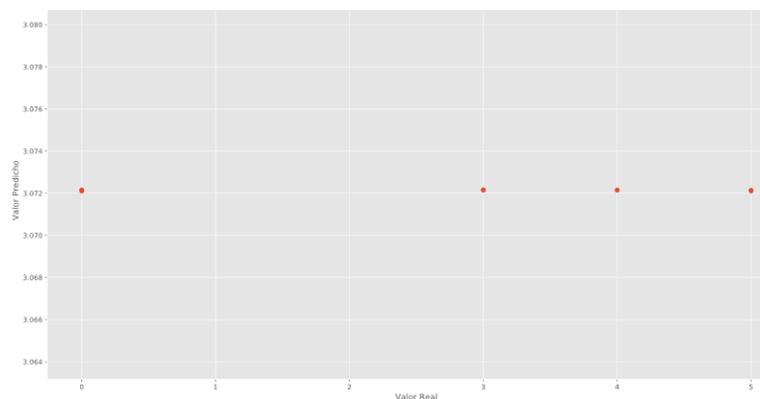
explosivo u otra que con un modelo de dos clases que predice la concentración por cada clase de sustancia.



**Figura 69.** Valores predichos por los modelos de regresión

(a) Regresión de sustancias explosivas puras (b) Regresión de TNT y pólvora en estado puro

Con relación a los modelos de la base de datos 2, el prototipo tuvo un mal desempeño en todos los casos ( $R^2 < 0.4$  y  $MSE > 0$ ), como se observa en la Figura 70 de los valores reales vs los valores predichos por el modelo de regresión de sustancias explosivas puras. Las observaciones sin sustancia explosiva son incorrectamente identificadas con valores de 3.072 gr.



**Figura 70.** Valores predichos por el modelo de regresión de sustancias explosivas puras para la base de datos 2

En el presente capítulo se describieron los conceptos utilizados para el desarrollo de los modelos de clasificación y de regresión con redes neuronales profundas LSTM. Posteriormente, se describió el proceso para la creación de los modelos en el que se definió el tipo de escalamiento usado y cuál fue el uso que se dio a la técnica de validación cruzada. Finalmente se presentó una descripción de los resultados obtenidos para cada una de las tareas.

# CAPÍTULO 8

## Análisis de Resultados y Pruebas

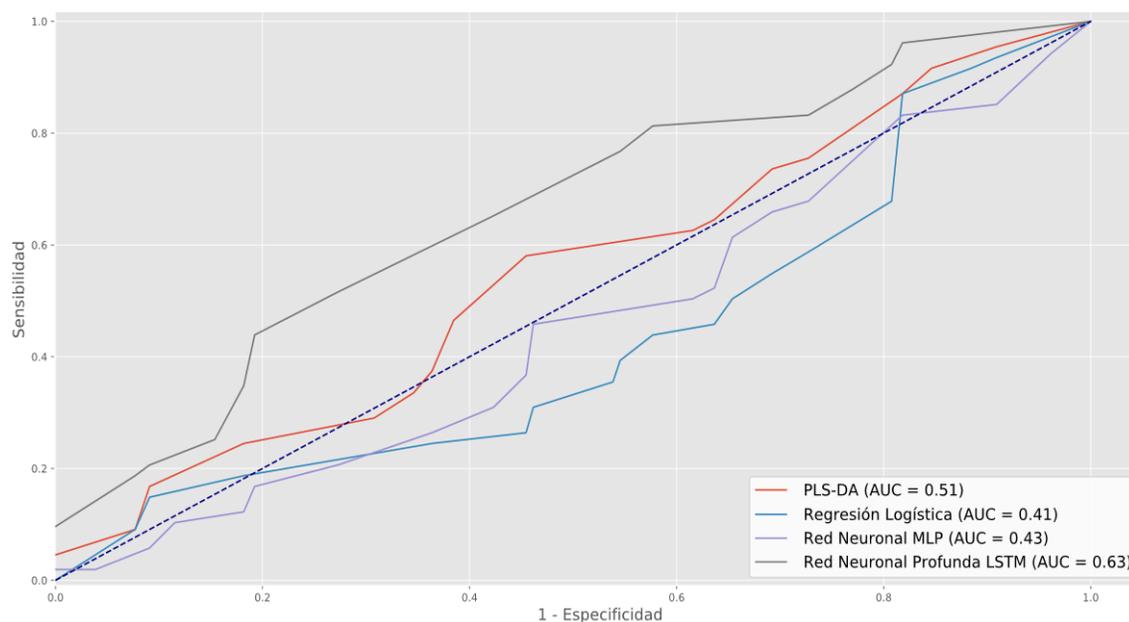
En este capítulo se presenta un análisis comparativo del desempeño de los modelos de clasificación y cuantificación de sustancias explosivas presentados en los capítulos previos en base al área bajo la curva (AUC) y del R<sup>2</sup>, respectivamente, además, se presentan los resultados reales y predichos con el modelo con mejor desempeño para cada una de las tareas, además, para el caso de los modelos de clasificación se definió un umbral para determinar la pertenencia de los experimentos a una clase u otra. A continuación, se describe el desempeño de la interfaz gráfica de usuario para el entrenamiento y prueba de los modelos. Finalmente, se discuten los resultados en base al desempeño de los modelos y el tiempo de entrenamiento de los mismos

### 8.1. Análisis Comparativo de los Modelos de Clasificación

Para realizar el análisis comparativo de los modelos de clasificación se elaboró una curva ROC por cada tarea, la cual evalúa el desempeño de los modelos generados con la base de datos inicial generada en (Salazar, 2018) y con la base de datos actual del prototipo optimizado en (Jacome, 2019), en base a la sensibilidad y tasa de falsos positivos (1-especificidad) con diferentes valores de umbral.

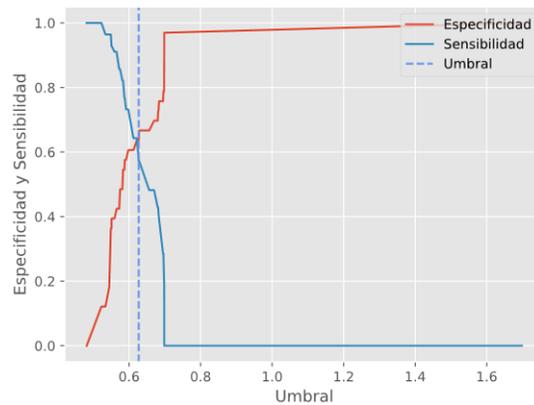
### 8.1.1. Modelos de Tarea 1: Clasificación de Alcohol/Sustancias Explosivas en Estado Puro y Mezclas

Los modelos generados para la Tarea 1 permiten identificar si el experimento analizado es o no sustancia explosiva (pura o mezclada con jabón o pasta dental) pero no su tipo. En la Figura 71, se observa la curva ROC de los modelos de esta tarea: PLS-DA, Regresión Logística, red neuronal MLP y red neuronal profunda LSTM, con las cuales se determina que el modelo con mejor desempeño es la red neuronal profunda LSTM con un AUC de 0.63 para los datos de prueba.



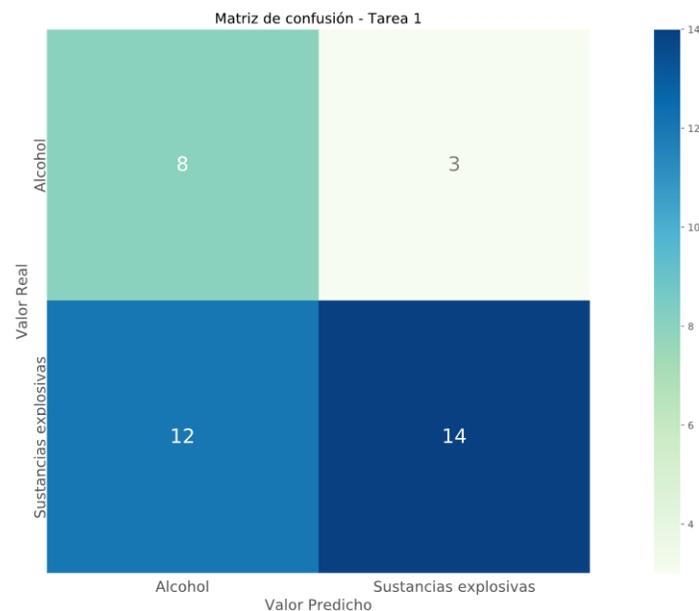
**Figura 71.** Curva ROC de modelos de clasificación de alcohol/sustancias explosivas en estado puro y mezclas

Por lo tanto, se utilizó este modelo para evaluar la sensibilidad y especificidad del modelo con un umbral fijo de 0.63. Como se observa en la Figura 72, este umbral se seleccionó en el punto en que la sensibilidad y especificidad tienen el mismo valor, es decir, en el cual la tasa de falsos positivos y falsos negativos es baja.



**Figura 72.** Umbral de modelo de clasificación de alcohol/sustancias explosivas en estado puro y mezclas LSTM

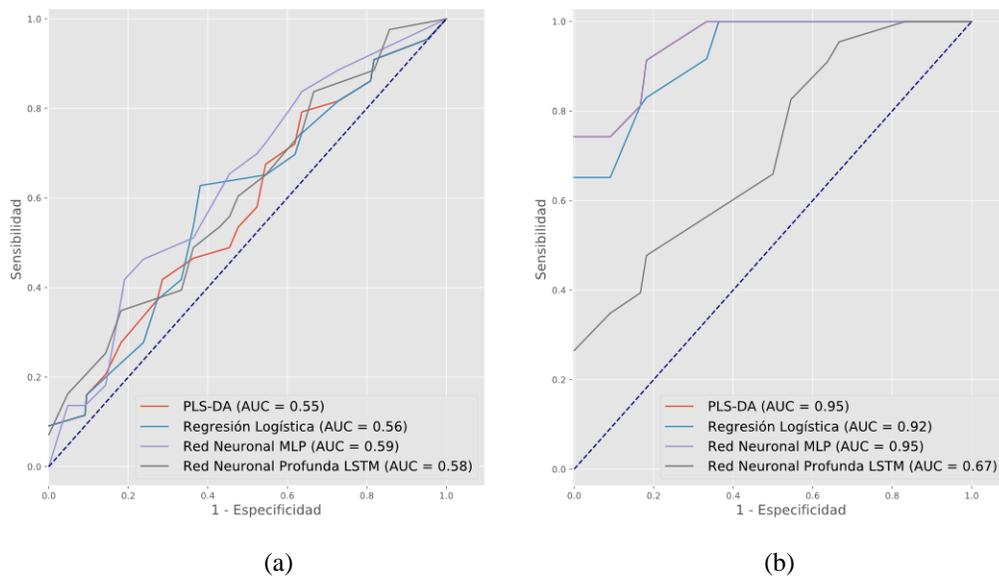
Los resultados de la clasificación de los datos de prueba para el modelo de deep learning LSTM se muestra en la matriz de confusión de la Figura 73, de la cual se determina que de 37 observaciones, 14 son clasificadas correctamente como sustancias explosivas y 8 como alcohol (sustancias no explosivas), por lo tanto, la sensibilidad de este modelo en la detección de sustancias explosivas puras o con mezclas de pólvora y jabón con concentraciones entre 0.1 y 3 gr es de 0.54, la especificidad de 0.73 y la precisión de 0.82.



**Figura 73.** Matriz de confusión del modelo de clasificación de alcohol/sustancias explosivas en estado puro y mezclas

### 8.1.2. Modelos de Tarea 2: Clasificación de Alcohol/Sustancias Explosivas en Estado Puro

Los modelos generados para la Tarea 2 permiten identificar si el experimento analizado es o no sustancia explosiva pura pero no su tipo. En la Figura 74, se observan las curvas ROC de los modelos de esta tarea: PLS-DA, regresión logística, red neuronal MLP y red neuronal profunda LSTM. En la Figura 74a para la base de datos 1, se determina que el modelo con mejor desempeño es la red neuronal MLP con un AUC de 0.59 para los datos de prueba y en la figura Figura 74b para la base de datos 2 la red neuronal MLP y PLS-DA ambos con un AUC de 0.95 para los datos de prueba.

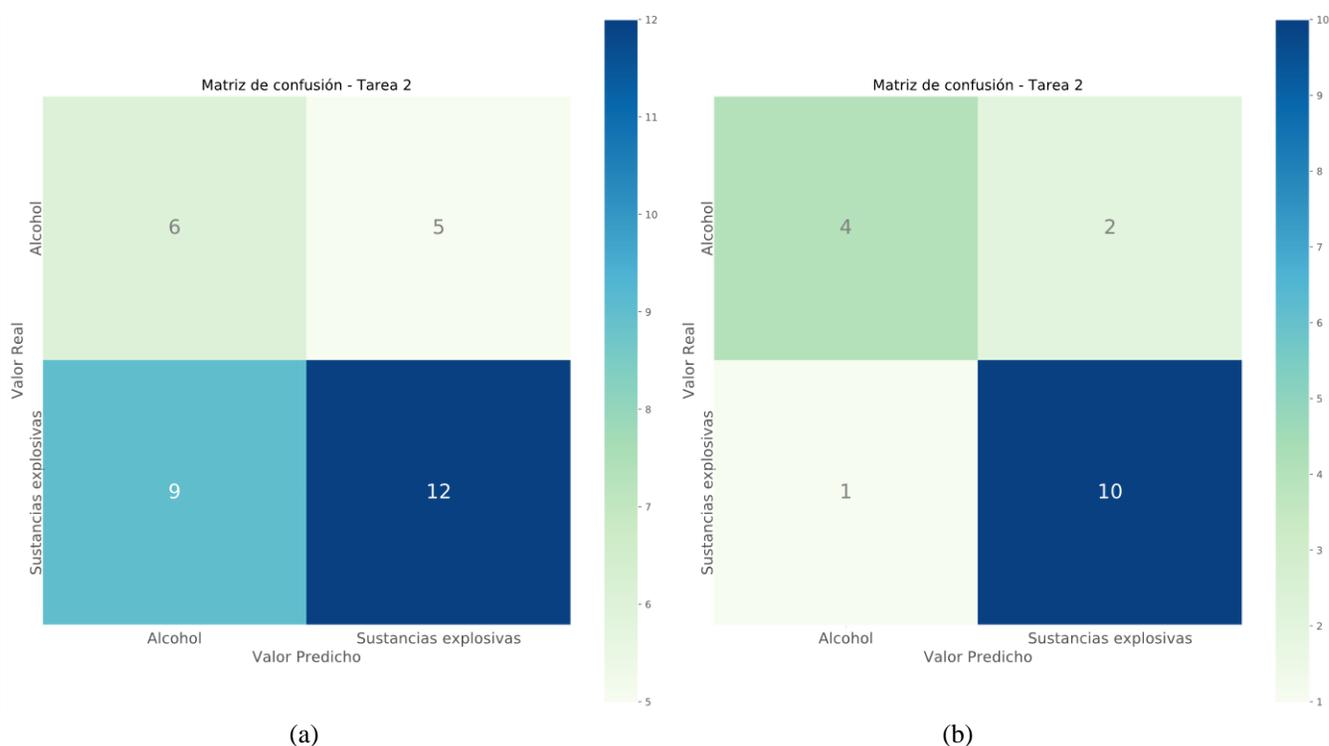


**Figura 74.** Curva ROC de modelos de clasificación de alcohol/sustancias explosivas en estado puro

(a) Base de datos 1, (b) Base de datos 2

Por lo tanto, se utilizó los modelos generados con la red neuronal MLP para evaluar su sensibilidad y especificidad con un umbral fijo de 0.5. Este umbral es el mismo para ambas bases de datos y se seleccionó en el punto en que la sensibilidad y especificidad tienen el mismo valor, es decir, en el cual la tasa de falsos positivos y falsos negativos es baja.

Los resultados para la base de datos 1 de la clasificación de los datos de prueba con el modelo de redes neuronales MLP se muestran en la matriz de confusión de la Figura 75a, de la cual se determina que de 32 observaciones, 12 son clasificadas correctamente como sustancias explosivas y 6 como alcohol (sustancias no explosivas), por lo tanto, la sensibilidad de este modelo en la detección de sustancias explosivas puras con concentraciones entre 0.1 y 3 gr es de 0.52, la especificidad de 0.55 y la precisión de 0.69. Los resultados para la base de datos 2 se muestran en la matriz de confusión de la Figura 75b, de la cual se determina que de 17 observaciones, 10 son clasificadas correctamente como sustancias explosivas y 4 como alcohol (sustancias no explosivas), por lo tanto, la sensibilidad de este modelo en la detección de sustancias explosivas puras con concentraciones entre 3 y 5 gr es de 0.91, la especificidad de 0.67 y la precisión de 0.83.

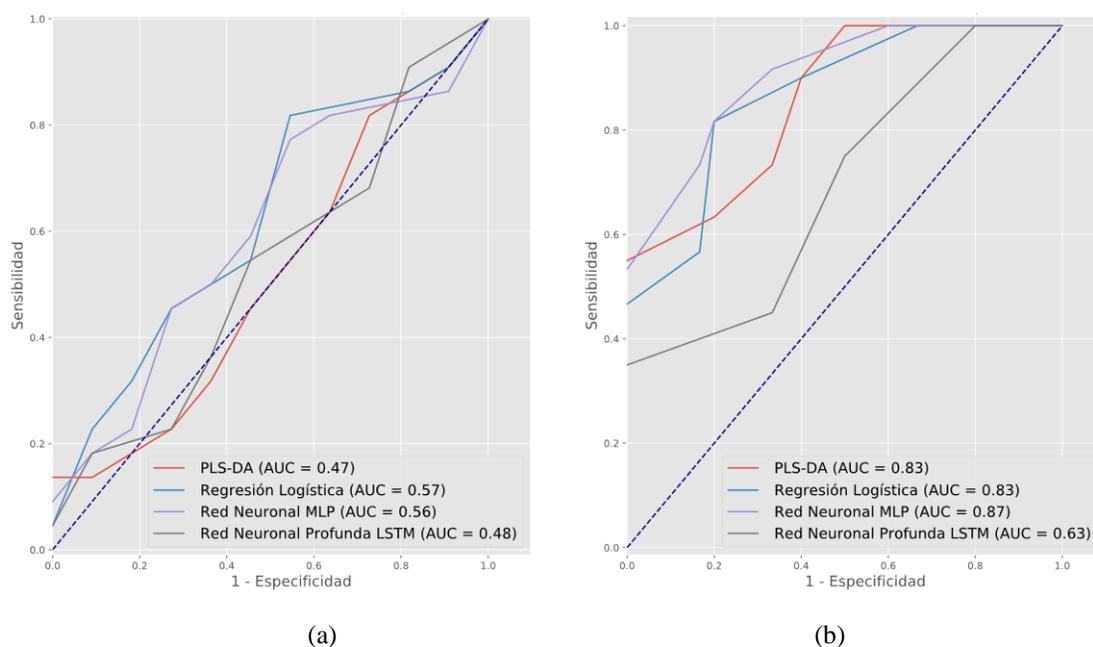


**Figura 75.** Matriz de confusión del modelo de clasificación de alcohol/sustancias explosivas en estado puro

(a) Base de datos 1, (b) Base de datos 2

### 8.1.3. Modelos de Tarea 3: Clasificación de Alcohol/Pólvora en Estado Puro

Los modelos generados para la Tarea 3 permiten identificar si el experimento analizado es pólvora en estado puro o alcohol. En la Figura 76, se observan las curvas ROC de los modelos de esta tarea: PLS-DA, Regresión Logística, red neuronal MLP y red neuronal profunda LSTM. En la Figura 76a para la base de datos 1, se determina que el modelo de regresión logística es aquel con mejor desempeño con un AUC de 0.57 para los datos de prueba y en la figura Figura 76b para la base de datos 2 la red neuronal MLP con un AUC de 0.87 para los datos de prueba.

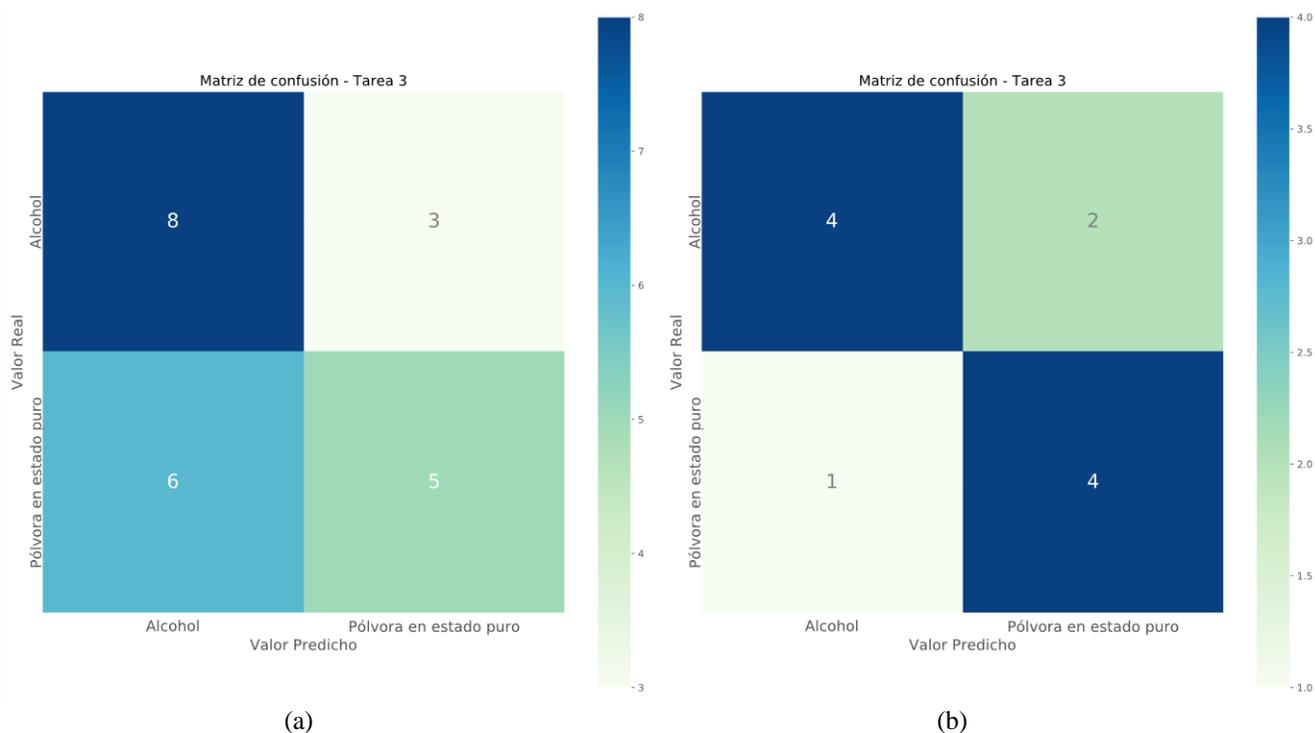


**Figura 76.** Curva ROC de modelos de clasificación de alcohol/pólvora en estado puro

(a) Base de datos 1, (b) Base de datos 2

Por lo tanto, se utilizó los modelos generados con el método de regresión logística para la base de datos 1 y con la red neuronal MLP para la base de datos 2, con las cuales se evalúa su sensibilidad y especificidad con un umbral fijo de 0.5. Este umbral es el mismo para ambas bases de datos y se seleccionó en el punto en que la sensibilidad y especificidad tienen el mismo valor, es decir, en el cual la tasa de falsos positivos y falsos negativos es baja.

Los resultados para la base de datos 1 de la clasificación de los datos de prueba con el modelo de regresión logística se muestran en la matriz de confusión de la Figura 77a, de la cual se determina que de 22 observaciones, 5 son clasificadas correctamente como pólvora en estado puro (sustancia explosiva) y 8 como alcohol (sustancia no explosiva), por lo tanto, la sensibilidad de este modelo en la detección de sustancias explosivas puras con concentraciones entre 0.1 y 3 gr es de 0.45, la especificidad de 0.73 y la precisión de 0.62. Los resultados para la base de datos 2 con la red neuronal MLP se muestran en la matriz de confusión de la Figura 77b, de la cual se determina que de 11 observaciones, 4 son clasificadas correctamente como pólvora en estado puro (sustancia explosiva) y 4 como alcohol (sustancia no explosiva), por lo tanto, la sensibilidad de este modelo en la detección de sustancias explosivas puras con concentraciones entre 3 y 5 gr es de 0.80, la especificidad de 0.67 y la precisión de 0.67.

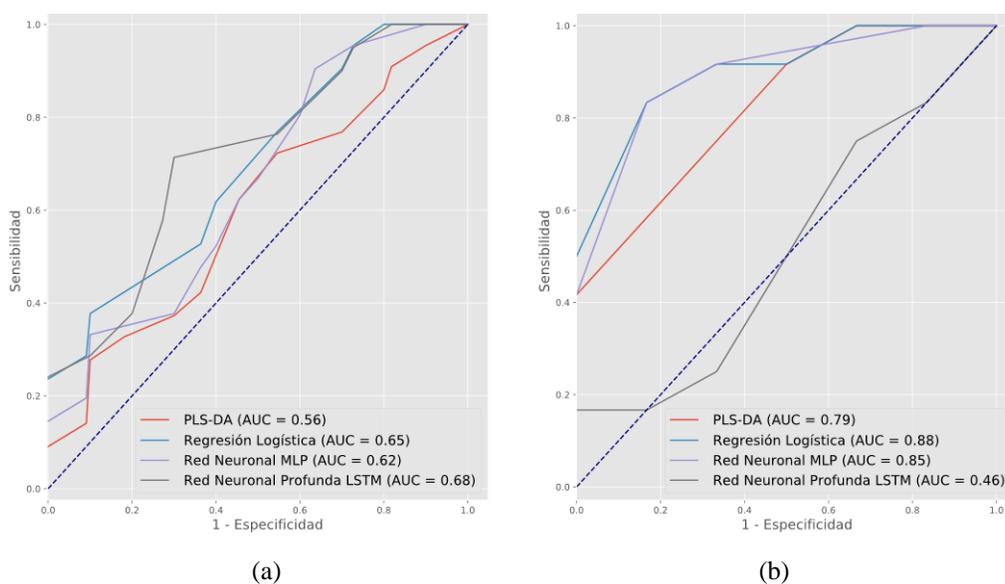


**Figura 77.** Matriz de confusión del modelo de clasificación de alcohol/pólvora en estado puro

(a) Base de datos 1, (b) Base de datos 2

### 8.1.4. Modelos de Tarea 4: Clasificación de Alcohol/TNT en Estado Puro

Los modelos generados para la Tarea 4 permiten identificar si el experimento analizado es TNT en estado puro o alcohol. En la Figura 78, se observan las curvas ROC de los modelos de esta tarea: PLS-DA, Regresión Logística, red neuronal MLP y red neuronal profunda LSTM. En la Figura 78a para la base de datos 1, se determina que el modelo generado con la red neuronal profunda LSTM es aquel con mejor desempeño con un AUC de 0.68 para los datos de prueba y en la Figura 78b para la base de datos 2 el modelo de regresión logística con un AUC de 0.88 para los datos de prueba.

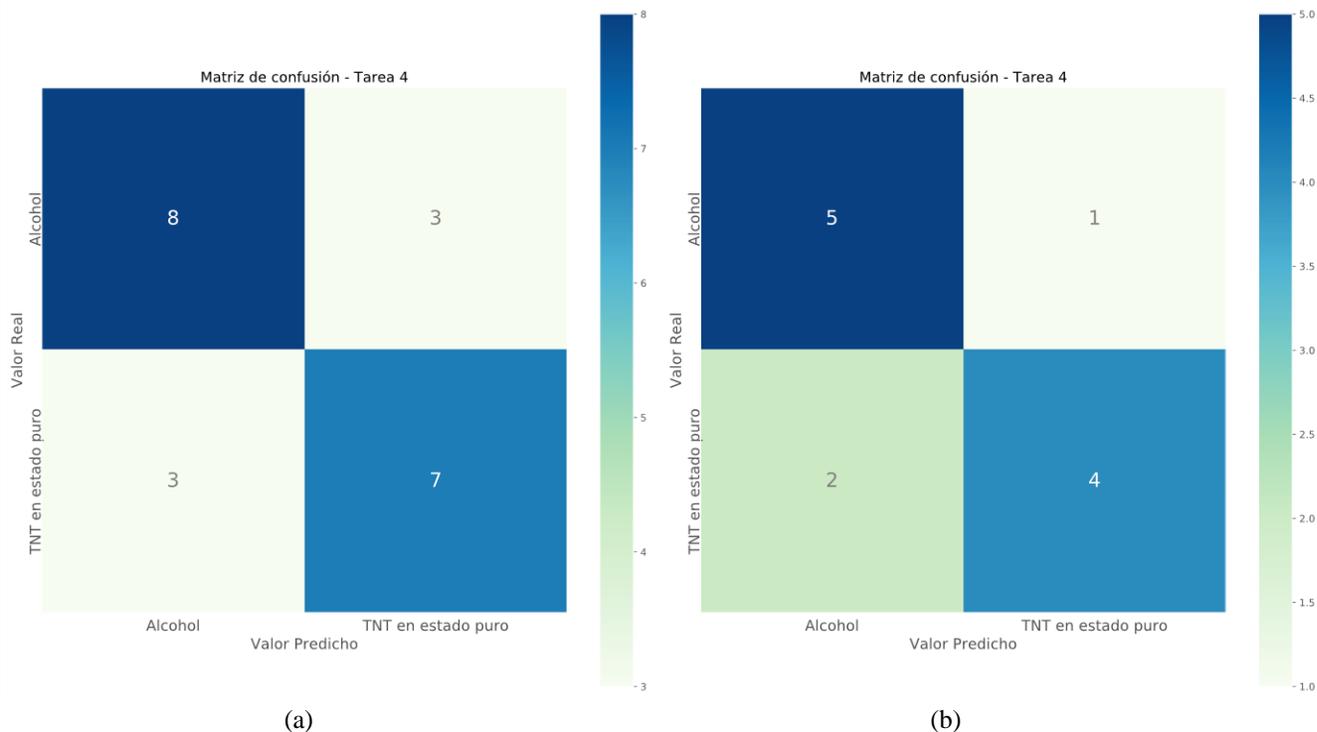


**Figura 78.** Curva ROC de modelos de clasificación de alcohol/TNT en estado puro

(a) Base de datos 1, (b) Base de datos 2

Por lo tanto, se utilizó los modelos generados con la red LSTM para la base de datos 1 y con el método de regresión logística para la base de datos 2, con las cuales se evalúa su sensibilidad y especificidad con un umbral fijo de 0.37 para el modelo de la base de datos 1 y de 0.5 para el de la base de datos 2. Este umbral se seleccionó en el punto en que la sensibilidad y especificidad tienen el mismo valor, es decir, en el cual la tasa de falsos positivos y falsos negativos es baja.

Los resultados para la base de datos 1 de la clasificación de los datos de prueba con el modelo de redes neuronales profundas LSTM se muestran en la matriz de confusión de la Figura 79a, de la cual se determina que de 21 observaciones, 7 son clasificadas correctamente como TNT en estado puro (sustancia explosiva) y 8 como alcohol (sustancia no explosiva), por lo tanto, la sensibilidad de este modelo en la detección de TNT en estado puro con concentraciones entre 0.1 y 3 gr es de 0.70, la especificidad de 0.73 y la precisión de 0.70. Los resultados para la base de datos 2 con el método de regresión logística se muestran en la matriz de confusión de la Figura 79b, de la cual se determina que de 12 observaciones, 4 son clasificadas correctamente como TNT en estado puro (sustancia explosiva) y 5 como alcohol (sustancia no explosiva), por lo tanto, la sensibilidad de este modelo en la detección de TNT en estado puro con concentraciones entre 3 y 5 gr es de 0.67, la especificidad de 0.83 y la precisión de 0.80.

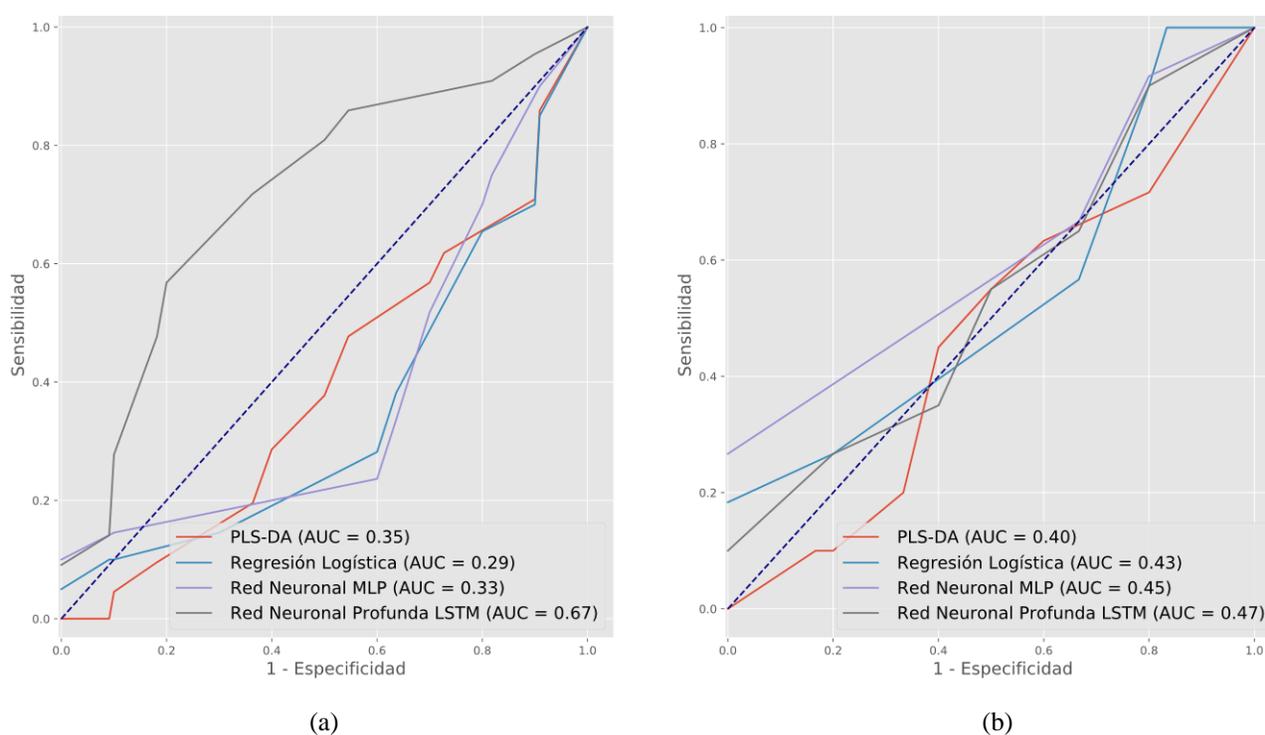


**Figura 79.** Matriz de confusión del modelo de clasificación de alcohol/TNT en estado puro

(a) Base de datos 1, (b) Base de datos 2

### 8.1.5. Modelos de Tarea 5: Clasificación de Pólvora en Estado Puro/TNT en Estado Puro

Los modelos generados para la Tarea 5 permiten identificar si el experimento analizado es pólvora o TNT en estado puro. En la Figura 80, se observan las curvas ROC de los modelos de esta tarea: PLS-DA, Regresión Logística, red neuronal MLP y red neuronal profunda LSTM. En la Figura 80a para la base de datos 1, se determina que el modelo generado con la red neuronal profunda LSTM es aquel con mejor desempeño con un AUC de 0.67 para los datos de prueba y en la Figura 80b para la base de datos 2 de igual forma el modelo generado con la red neuronal profunda LSTM con un AUC de 0.47 para los datos de prueba.



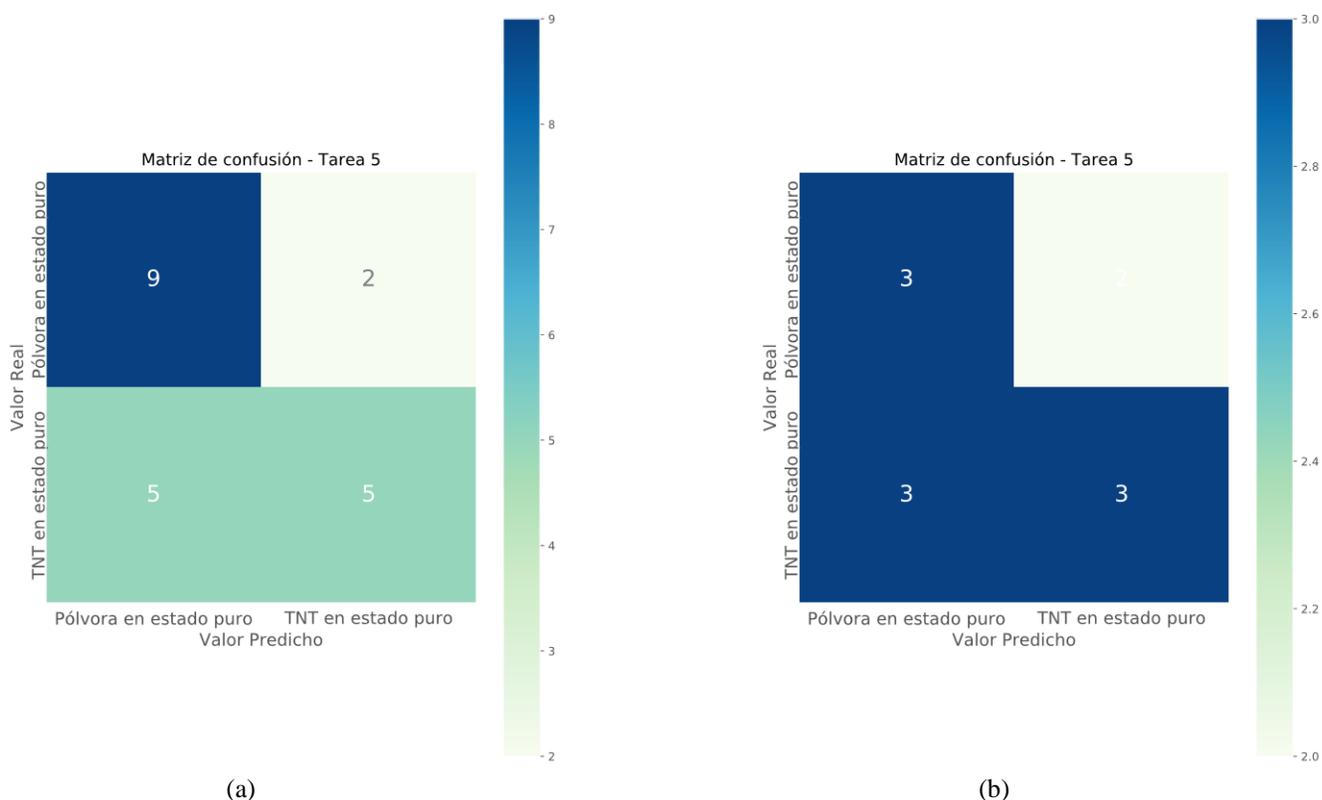
**Figura 80.** Curva ROC de modelos de clasificación de alcohol/TNT en estado puro

(a) Base de datos 1, (b) Base de datos 2

Por lo tanto, se utilizó los modelos generados con la red LSTM para la base de datos 1 y 2, con las cuales se evalúa su sensibilidad y especificidad con un umbral fijo de 0.48 para el modelo de la base de datos 1 y de 0.37 para el de la base de datos 2. Este umbral se seleccionó en el punto en que la sensibilidad y especificidad tienen el mismo valor, es decir, en el cual la tasa de falsos positivos y falsos negativos es baja.

Los resultados para la base de datos 1 de la clasificación de los datos de prueba con el modelo de redes neuronales profundas LSTM se muestran en la matriz de confusión de la Figura 81a, de la cual se determina que de 21 observaciones, 9 son clasificadas correctamente como pólvora en estado puro y 5 como TNT en estado puro, por lo tanto, la sensibilidad de este modelo en la detección de pólvora en estado puro con concentraciones entre 0.1 y 3 gr es de 0.82, la especificidad de 0.50 y la precisión de 0.64, además, la sensibilidad del modelo en la detección de TNT en estado puro con concentraciones entre 0.1 y 3 gr es de 0.50, la especificidad de 0.82 y la precisión de 0.71.

Los resultados para la base de datos 2 se muestran en la matriz de confusión de la Figura 81b, de la cual se determina que de 11 observaciones, 3 son clasificadas correctamente como pólvora en estado puro y 3 como TNT en estado puro, por lo tanto, la sensibilidad de este modelo en la detección de pólvora en estado puro con concentraciones entre 3 y 5 gr es de 0.60, la especificidad de 0.50 y la precisión de 0.50, además, la sensibilidad del modelo en la detección de TNT en estado puro con concentraciones entre 3 y 5 gr es de 0.50, la especificidad de 0.60 y la precisión de 0.60.

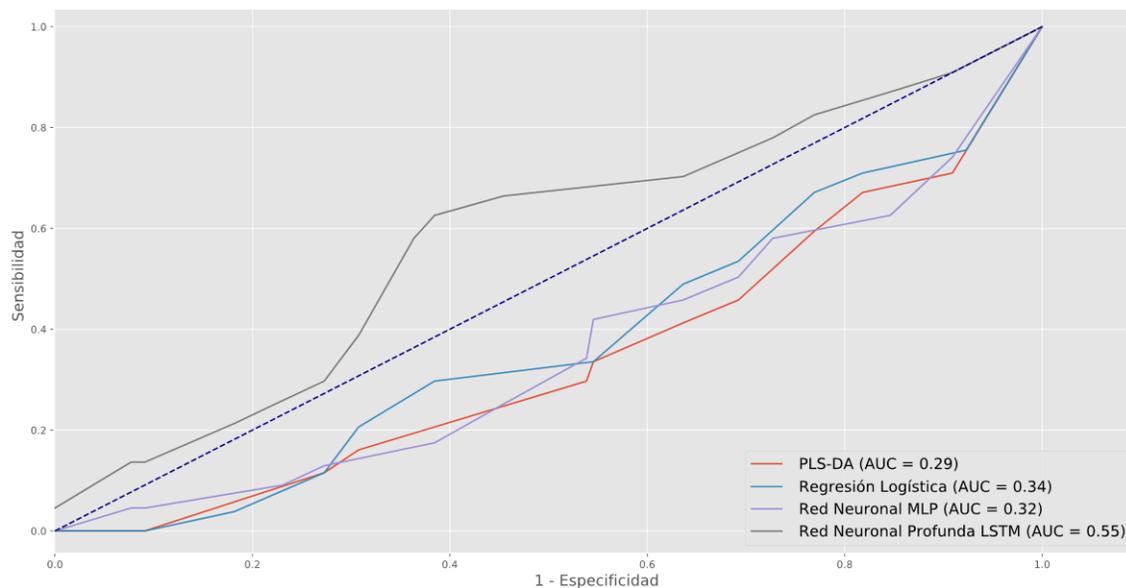


**Figura 81.** Matriz de confusión del modelo de clasificación de pólvora en estado puro /TNT en estado puro

(a) Base de datos 1, (b) Base de datos 2

### 8.1.6. Modelos de Tarea 6: Clasificación de Alcohol/Pólvora en Estado Puro y Mezcla de Pólvora

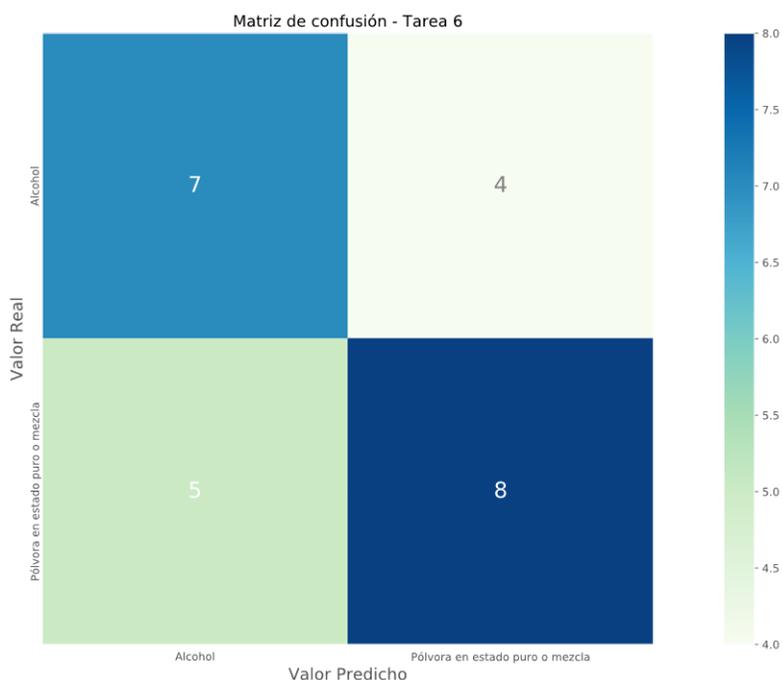
Los modelos generados para la Tarea 6 permiten identificar si el experimento analizado es alcohol o pólvora en estado puro o mezclado con jabón o pasta dental. En la Figura 82, se observan las curvas ROC de los modelos para la base de datos 1 de esta tarea: PLS-DA, Regresión Logística, red neuronal MLP y red neuronal profunda LSTM. De las cuales, se determina que el modelo generado con la red neuronal profunda LSTM es aquel con mejor desempeño con un AUC de 0.55 para los datos de prueba.



**Figura 82.** Curva ROC de modelos de clasificación de Alcohol/Pólvora en estado puro y mezcla de pólvora

Por lo tanto, se utilizó el modelo generado con la red LSTM, con la cual se evalúa su sensibilidad y especificidad con un umbral fijo de 0.42. Este umbral se seleccionó en el punto en que la sensibilidad y especificidad tienen el mismo valor, es decir, en el cual la tasa de falsos positivos y falsos negativos es baja.

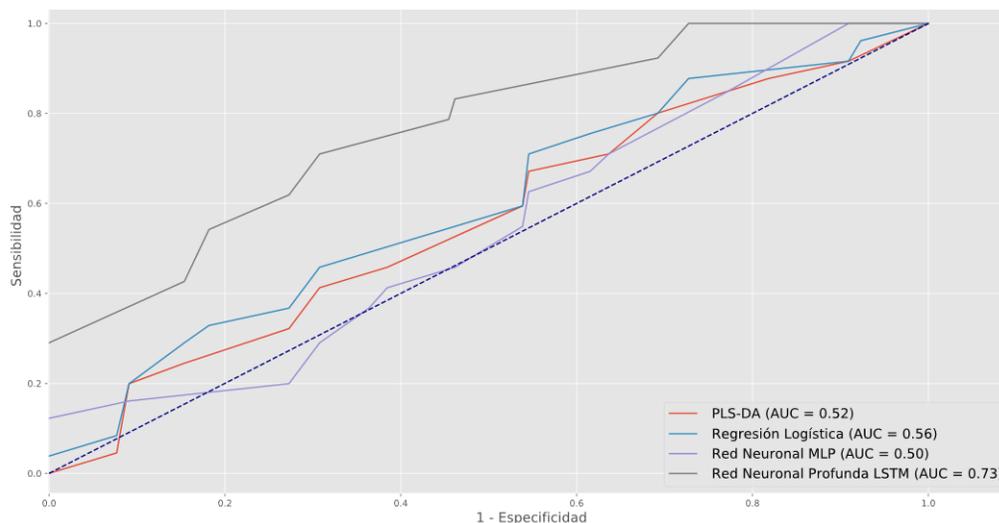
Los resultados de la clasificación de los datos de prueba con el modelo de redes neuronales profundas LSTM se muestran en la matriz de confusión de la Figura 83, de la cual se determina que de 24 observaciones, 7 son clasificadas correctamente como alcohol y 8 como pólvora en estado puro o mezclado con jabón o pasta dental, por lo tanto, la sensibilidad de este modelo en la detección de pólvora con concentraciones entre 0.1 y 3 gr es de 0.62, la especificidad de 0.64 y la precisión de 0.67.



**Figura 83.** Matriz de confusión del modelo de clasificación de Alcohol/Pólvora en estado puro y mezcla de pólvora

### 8.1.7. Modelos de Tarea 7: Clasificación de Alcohol/TNT en Estado Puro y Mezcla de TNT

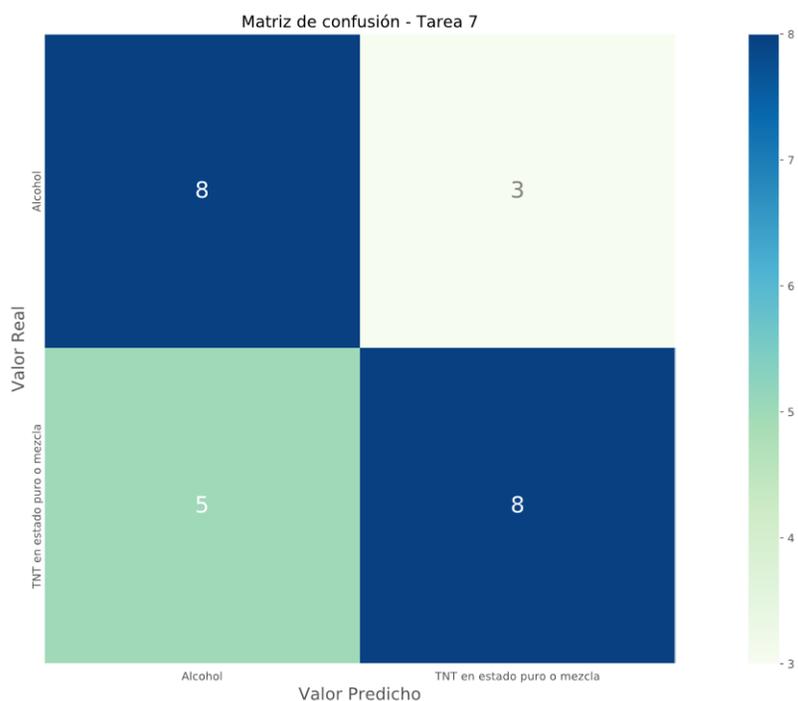
Los modelos generados para la Tarea 7 permiten identificar si el experimento analizado es alcohol o TNT en estado puro o mezclado con jabón o pasta dental. En la Figura 84, se observan las curvas ROC de los modelos para la base de datos 1 de esta tarea: PLS-DA, Regresión Logística, red neuronal MLP y red neuronal profunda LSTM. De las cuales, se determina que el modelo generado con la red neuronal profunda LSTM es aquel con mejor desempeño con un AUC de 0.73 para los datos de prueba.



**Figura 84.** Curva ROC de modelos de clasificación de Alcohol/TNT en estado puro y mezcla de TNT

Por lo tanto, se utilizó el modelo generado con la red LSTM, con la cual se evalúa su sensibilidad y especificidad con un umbral fijo de 0.48. Este umbral se seleccionó en el punto en que la sensibilidad y especificidad tienen el mismo valor, es decir, en el cual la tasa de falsos positivos y falsos negativos es baja.

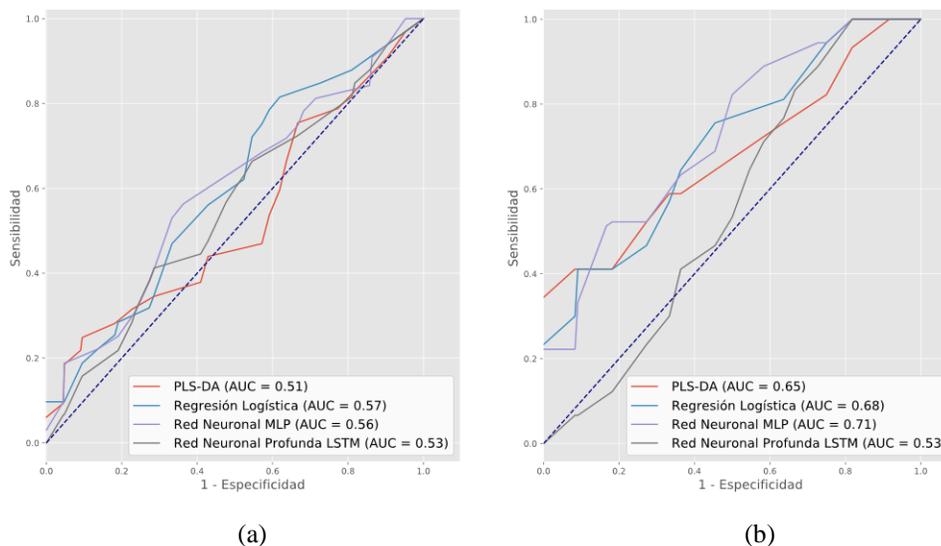
Los resultados de la clasificación de los datos de prueba con el modelo de redes neuronales profundas LSTM se muestran en la matriz de confusión de la Figura 85, de la cual se determina que de 24 observaciones, 8 son clasificadas correctamente como alcohol y 8 como TNT en estado puro o mezclado con pasta dental o jabón, por lo tanto, la sensibilidad de este modelo en la detección de TNT con concentraciones entre 0.1 y 3 gr es de 0.62, la especificidad de 0.73 y la precisión de 0.73.



**Figura 85.** Matriz de confusión del modelo de clasificación de Alcohol/TNT en estado puro y mezcla de TNT

### 8.1.8. Modelos de Tarea 8: Clasificación de Alcohol/Pólvora en Estado Puro/TNT en Estado Puro

Los modelos generados para la Tarea 8 permiten identificar si el experimento analizado es sustancia explosiva o no y su tipo: pólvora o TNT en estado puro. En la Figura 86, se observan las curvas ROC de los modelos de esta tarea: PLS-DA, Regresión Logística, red neuronal MLP y red neuronal profunda LSTM. En la Figura 86a para la base de datos 1, se determina que el modelo generado con el método de regresión logística es aquel con mejor desempeño con un AUC de 0.57 para los datos de prueba y en la Figura 86b para la base de datos 2 la red neuronal MLP con un AUC de 0.71 para los datos de prueba.



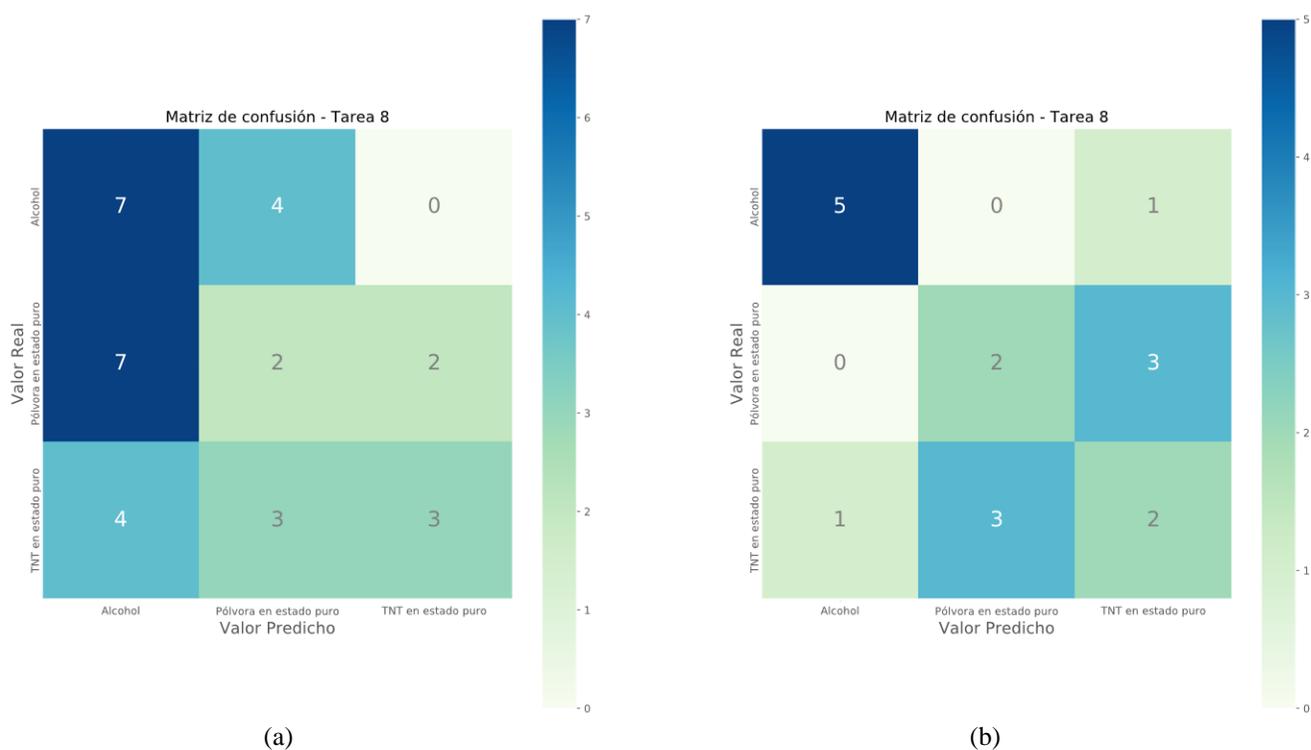
**Figura 86.** Curva ROC de modelos de clasificación de Alcohol/Pólvora en estado puro/TNT en estado puro

(a) Base de datos 1, (b) Base de datos 2

Por lo tanto, se utilizó los modelos generados con el método de regresión logística para la base de datos 1 y con la red neuronal MLP para la base de datos 2, con las cuales se evalúa su sensibilidad y especificidad con un umbral fijo de 0.35 para la clase alcohol, 0.42 para la clase pólvora y de 0.35 para la clase TNT, para de la base de datos 1 y de 0.5 para la clase alcohol, 0.5 para la clase pólvora y de 0.5 para la clase TNT, para la base de datos 2. Estos umbrales se seleccionaron en el punto en que la sensibilidad y especificidad tienen el mismo valor, es decir, en el cual la tasa de falsos positivos y falsos negativos es baja.

Los resultados para la base de datos 1 de la clasificación de los datos de prueba con el modelo de regresión logística se muestran en la matriz de confusión de la Figura 87a, de la cual se determina que de 32 observaciones, 7 son clasificadas correctamente como alcohol (sustancia no explosiva), 2 como pólvora en estado puro y 3 como TNT en estado puro, por lo tanto, la sensibilidad de este modelo en la detección de alcohol es de 0.64, la especificidad de 0.31 y la precisión de 0.39, en la detección de pólvora en estado puro con concentraciones entre 0.1 y 3 gr es de 0.18, la especificidad

de 0.59 y la precisión de 0.22 y en la detección de TNT en estado puro con concentraciones entre 0.1 y 3 gr la sensibilidad es de 0.30, la especificidad de 0.82 y la precisión de 0.60. Los resultados para la base de datos 2 de la clasificación de los datos de prueba con el modelo de redes neuronales MLP, se muestran en la matriz de confusión de la Figura 87b, de la cual se determina que de 17 observaciones, 5 son clasificadas correctamente como alcohol (sustancia no explosiva), 2 como pólvora en estado puro y 2 como TNT en estado puro, por lo tanto, la sensibilidad de este modelo en la detección de alcohol es de 0.83, la especificidad de 0.80 y la precisión de 0.83, en la detección de pólvora en estado puro con concentraciones entre 3 y 5 gr la sensibilidad es de 0.40, la especificidad de 0.70 y la precisión de 0.40, y en la detección de TNT en estado puro con concentraciones entre 3 y 5 gr la sensibilidad es de 0.33, la especificidad de 0.64 y la precisión de 0.33.

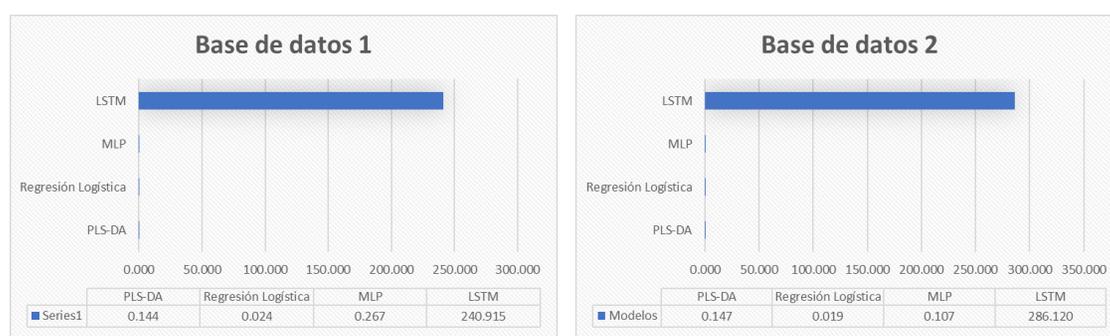


**Figura 87.** Matriz de confusión del modelo de clasificación de Alcohol/Pólvora en estado puro/TNT en estado puro.

(a) Base de datos 1, (b) Base de datos 2

### 8.1.9. Tiempo de Entrenamiento de Modelos de Clasificación

En la Figura 88, se presenta el tiempo total de entrenamiento para los 8 modelos de clasificación con cada uno de los algoritmos para la base de datos 1 (Figura 88a) y 2 (Figura 88b), con el cual se determina que el tiempo de entrenamiento de los modelos con el método de regresión logística es inferior al tiempo empleado para el entrenamiento de los modelos generados el método de deep learning mediante la red neuronal profunda LSTM .



(a)

(b)

**Figura 88.** Tiempo de entrenamiento de modelos de clasificación. (a) Base de datos 1, (b) Base de datos 2

Estos resultados se obtuvieron en un computador con un procesador Inter(R) Core (TM) i7-8550U, 12 GB RAM, ejecutando Spyder 3.3.6 en Windows 10.

## 8.2. Análisis Comparativo de los Modelos de Cuantificación

Para realizar el análisis comparativo de los modelos de clasificación se compara el error cuadrático medio (MSE) y el valor R2 para los datos de prueba de cada uno de los modelos y se selecciona el modelo con mejor desempeño para la representación gráfica del valor predicho y el real de los experimentos.

### 8.2.1. Modelos de Tarea 1: Regresión de Alcohol/Sustancias Explosivas en Estado Puro y Mezclas

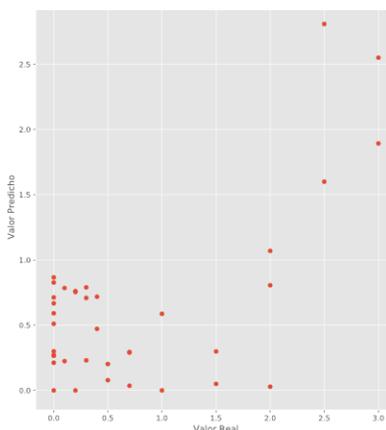
Los modelos generados para la Tarea 1 permiten identificar la cantidad sustancia explosiva (pura o mezclada con jabón o pasta dental) del experimento analizado, sin importar su clase. En la **Tabla 28**, se observa el MSE y R2 de esta tarea para los modelos de regresión: PLS-R, red neuronal MLP y red neuronal profunda LSTM, con las cuales se determina que aquel con mejor desempeño es el modelo PLS-R con un MSE de 0.51 y un R2 de 0.38 para los datos de prueba.

**Tabla 28**

*MSE y R2 de modelos de regresión de alcohol/sustancias explosivas en estado puro y mezclas*

TAREA	PLS-R		MLP		LSTM	
	MSE	R2	MSE	R2	MSE	R2
Alcohol/Sustancias explosivas en estado puro y mezclas	0.51	0.38	0.59	0.3	0.77	0.10

Los resultados de la predicción de la cantidad de sustancia explosiva (pura o mezclada con jabón o pasta dental) entre 0 y 3gr con los datos de prueba para el modelo PLS-R se muestra en la Figura 89, de la cual se determina que el valor predicho por el modelo se aleja del valor real en experimentos con concentraciones intermedias (entre 0.5 y 2.5 gr).



**Figura 89.** Valores predichos por el modelo de regresión de alcohol/sustancias explosivas en estado puro y mezclas

### 8.2.2. Modelos de Tarea 2: Regresión de Alcohol/Sustancias Explosivas en Estado Puro

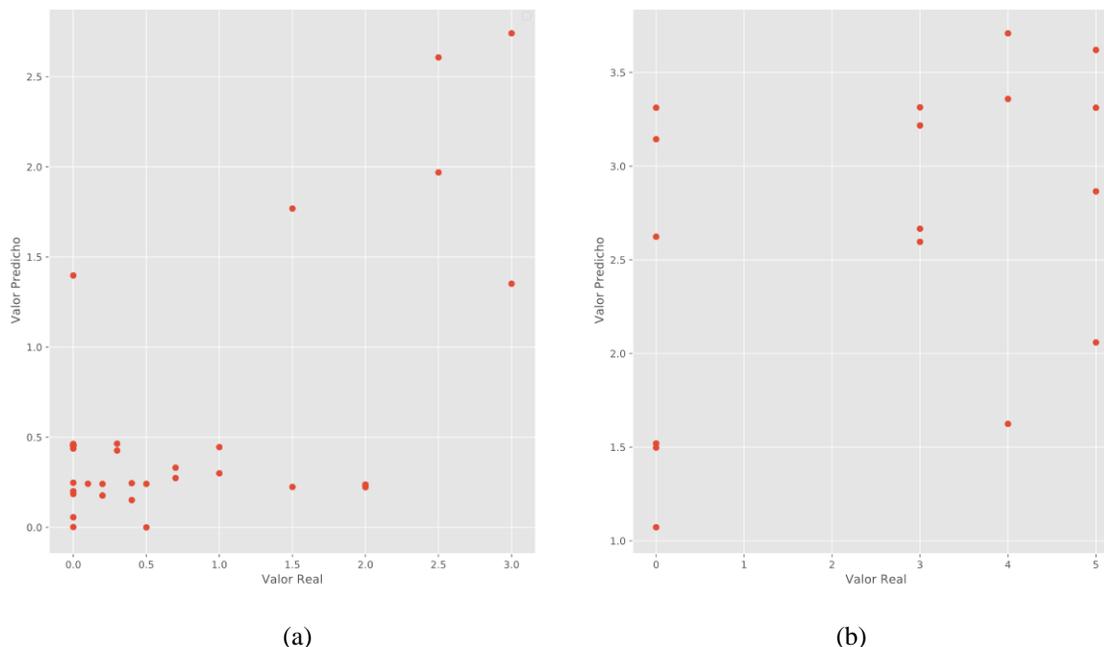
Los modelos generados para la Tarea 2 permiten identificar la cantidad sustancia explosiva pura del experimento analizado, sin importar su clase. En la **Tabla 29**, se observa el MSE y R2 de esta tarea para los modelos de regresión: PLS-R, red neuronal MLP y red neuronal profunda LSTM, con las cuales se determina que el modelo con mejor desempeño para la base de daros 1 es la red neuronal MLP con un MSE de 0.47 y un R2 de 0.47 para los datos de prueba y para la base de daros 2 el modelo PLS-R con un MSE de 3.40 y un R2 de 0.17 para los datos de prueba.

**Tabla 29**

*MSE y R2 de modelos de regresión de alcohol/sustancias explosivas en estado puro*

TAREA	BASE DE DATOS 1						BASE DE DATOS 2					
	PLS-R		MLP		LSTM		PLS-R		MLP		LSTM	
	MSE	R2	MSE	R2	MSE	R2	MSE	R2	MSE	R2	MSE	R2
Alcohol/Sustancias explosivas en estado puro	0.47	0.46	0.47	0.47	0.75	0.17	3.40	0.17	3.71	0.10	4.36	-0.06

Los resultados de la predicción de la cantidad de sustancia explosiva pura entre 0 y 3 gr con los datos de prueba para el modelo MLP se muestra en la Figura 90a, de la cual se determina que el valor predicho por el modelo se aleja del valor real en experimentos con concentraciones intermedias (entre 1 y 2 gr). En el caso del modelo para la predicción de la cantidad de sustancia explosiva pura entre 3 y 5 gr con los datos de prueba para el modelo PLS-R se muestra en la Figura 90b, de la cual se determina que el valor predicho por el modelo se aleja del valor real en experimentos con concentraciones de 0 y 5 gr de sustancia explosiva pura.



**Figura 90.** Valores predichos por el modelo de regresión de alcohol/sustancias explosivas en estado puro

(a) Base de datos 1, (b) Base de datos 2

### 8.2.3. Modelos de Tarea 3: Regresión de Alcohol/Pólvora en Estado Puro

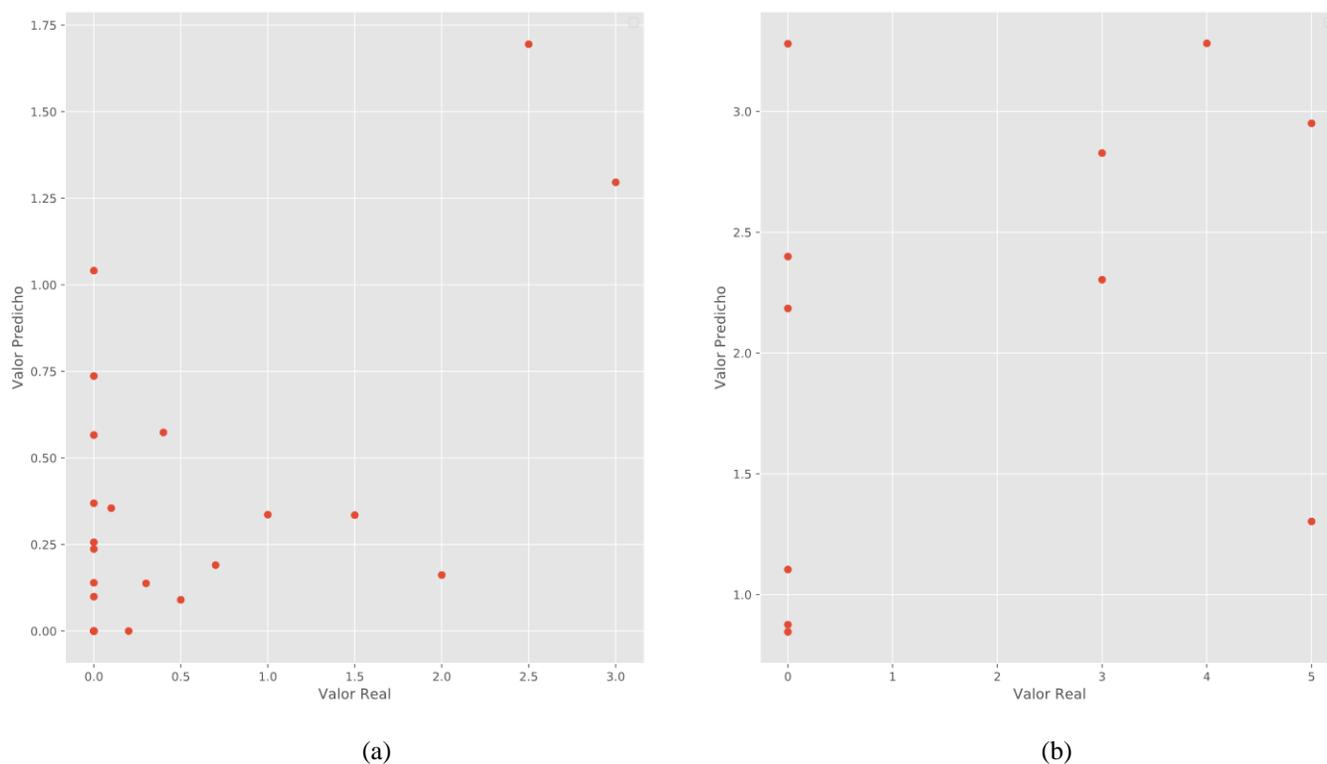
Los modelos generados para la Tarea 3 permiten identificar la cantidad pólvora en estado puro del experimento analizado. En la **Tabla 30**, se observa el MSE y R2 de esta tarea para los modelos de regresión: PLS-R, red neuronal MLP y red neuronal profunda LSTM, con las cuales se determina que el modelo con mejor desempeño para la base de datos 1 y 2 es el modelo PLS-R con un MSE de 0.52 y un R2 de 0.31 para los datos de prueba de la base de datos 1 y un MSE de 3.89 y un R2 de 0.09 para los datos de prueba de la base de datos 2.

**Tabla 30**

*MSE y R2 de modelos de regresión de alcohol/pólvora en estado puro*

TAREA	BASE DE DATOS 1						BASE DE DATOS 2					
	PLS-R		MLP		LSTM		PLS-R		MLP		LSTM	
	MSE	R2	MSE	R2	MSE	R2	MSE	R2	MSE	R2	MSE	R2
Alcohol/Pólvora en estado puro	0.52	0.31	0.59	0.23	0.75	0.01	3.89	0.09	5.92	-0.37	5.41	-0.25

Los resultados de la predicción de la cantidad de pólvora en estado puro entre 0 y 3 gr con los datos de prueba para el modelo PLS-R de la base de datos 1 se muestra en la Figura 91a, de la cual se determina que el valor predicho por el modelo se aleja del valor real, sin embargo, la concentración predicha por el modelo si aumenta a medida que la concentración real lo hace. En el caso del modelo para la predicción de la cantidad de sustancia explosiva pura entre 3 y 5 gr con los datos de prueba para el modelo PLS-R de la base de datos 2, se muestra en la Figura 91b, de la cual se determina que el valor predicho por el modelo se aleja del valor real en experimentos con concentraciones de 0 y 5 gr de pólvora en estado puro.



**Figura 91.** Valores predichos por el modelo de regresión de alcohol/pólvora en estado puro

(a) Base de datos 1, (b) Base de datos 2

#### 8.2.4. Modelos de Tarea 4: Regresión de Alcohol/TNT en Estado Puro

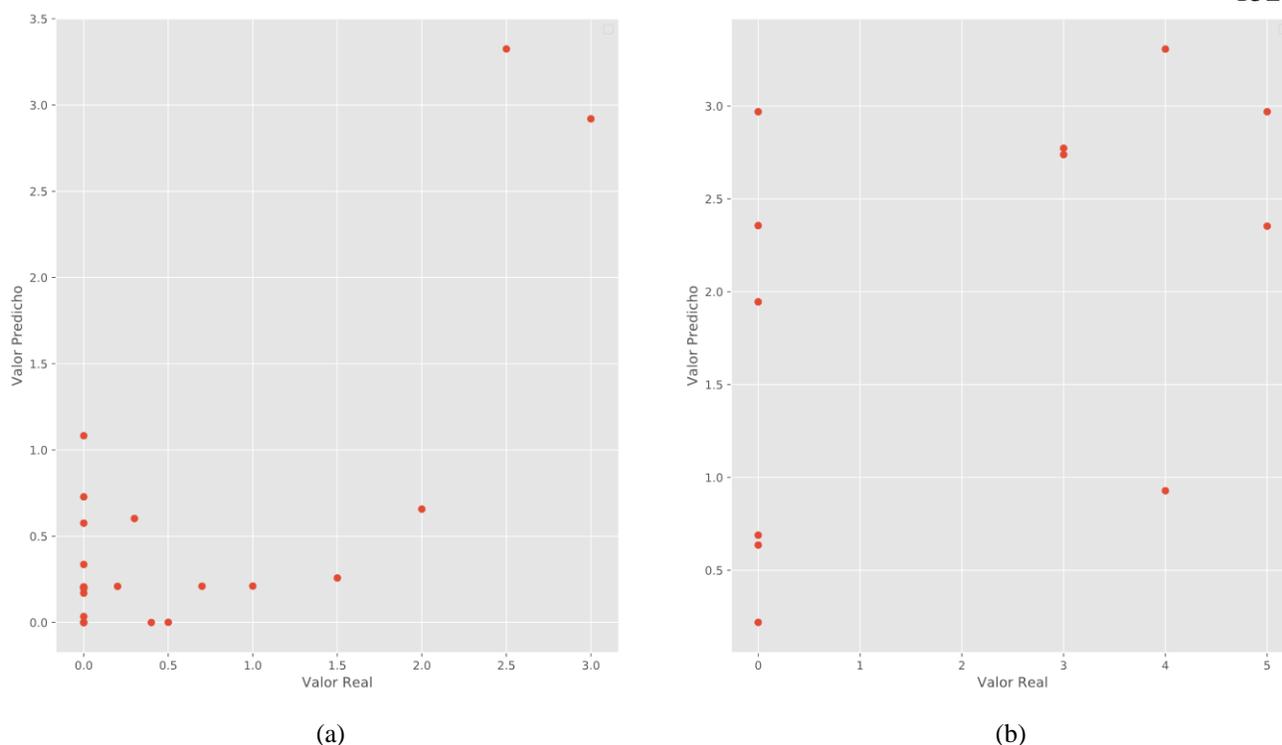
Los modelos generados para la Tarea 4 permiten identificar la cantidad TNT en estado puro del experimento analizado. En la **Tabla 31**, se observa el MSE y R2 de esta tarea para los modelos de regresión: PLS-R, red neuronal MLP y red neuronal profunda LSTM, con las cuales se determina que el modelo con mejor desempeño para la base de datos 1 y 2 es el modelo PLS-R con un MSE de 0.36 y un R2 de 0.53 para los datos de prueba de la base de datos 1 y un MSE de 3.35 y un R2 de 0.22 para los datos de prueba de la base de datos 2.

**Tabla 31**

*MSE y R2 de modelos de regresión de alcohol/pólvora en estado puro*

TAREA	BASE DE DATOS 1						BASE DE DATOS 2					
	PLS-R		MLP		LSTM		PLS-R		MLP		LSTM	
	MSE	R2	MSE	R2	MSE	R2	MSE	R2	MSE	R2	MSE	R2
Alcohol/TNT en estado puro	0.36	0.53	0.8	-0.01	0.71	0.10	3.35	0.22	3.59	0.17	4.69	-0.08

Los resultados de la predicción de la cantidad de TNT en estado puro entre 0 y 3 gr con los datos de prueba para el modelo PLS-R de la base de datos 1 se muestra en la Figura 92a, de la cual se determina que el valor predicho por el modelo se acerca al valor real, sobre todo con concentraciones de 0gr y 3gr de TNT, además, la concentración predicha por el modelo si aumenta a medida que la concentración real lo hace. En el caso del modelo para la predicción de la cantidad de sustancia explosiva pura entre 3 y 5 gr con los datos de prueba para el modelo PLS-R de la base de datos 2, se muestra en la Figura 92b, de la cual se determina que el valor predicho por el modelo se aleja del valor real en experimentos con concentraciones de 0 y 5 gr de sustancia explosiva pura.



**Figura 92.** Valores predichos por el modelo de regresión de alcohol/TNT en estado puro

(a) Base de datos 1, (b) Base de datos 2

### 8.2.5. Modelos de Tarea 5: Regresión de Pólvora en Estado Puro/TNT en Estado Puro

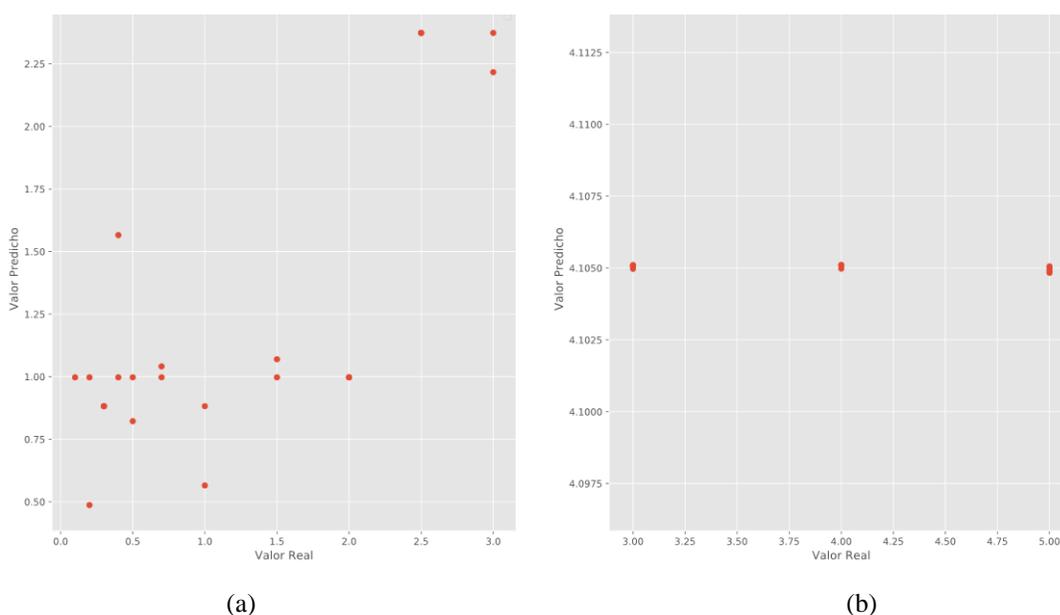
Los modelos generados para la Tarea 5 permiten identificar la cantidad de pólvora o TNT en estado puro del experimento analizado, sin importar su clase, su diferencia de los modelos generados para la Tarea 2 es que este modelo no se entrenó con sustancias no explosivas (0 gr de concentración). En la **Tabla 32**, se observa el MSE y R2 de esta tarea para los modelos de regresión: PLS-R, red neuronal MLP y red neuronal profunda LSTM, con las cuales se determina que el modelo con mejor desempeño para la base de datos 1 es la red neuronal MLP con un MSE de 0.39 y un R2 de 0.57, y la red neuronal profunda LSTM para la base de datos 2 con un MSE de 0.74 y un R2 de -0.02 para los datos de prueba de la base de datos 2.

**Tabla 32**

*MSE y R2 de modelos de pólvora en estado puro/TNT en estado puro*

TAREA	BASE DE DATOS 1						BASE DE DATOS 2					
	PLS-R		MLP		LSTM		PLS-R		MLP		LSTM	
	MSE	R2	MSE	R2	MSE	R2	MSE	R2	MSE	R2	MSE	R2
Pólvora en estado puro/TNT en estado puro	0.42	0.52	0.39	0.57	0.87	0.03	0.98	-0.34	1.19	-0.63	0.74	-0.02

Los resultados de la predicción de la cantidad de TNT en estado puro entre 0 y 3 gr con los datos de prueba para el modelo PLS-R de la base de datos 1 se muestra en la Figura 93a, de la cual se determina que el valor predicho por el modelo se acerca al valor real, sobre todo con concentraciones de 0gr y 3gr de TNT, además, la concentración predicha por el modelo si aumenta a medida que la concentración real lo hace. En el caso del modelo para la predicción de la cantidad de sustancia explosiva pura entre 3 y 5 gr con los datos de prueba para el modelo PLS-R de la base de datos 2, se muestra en la Figura 93b, de la cual se determina que el valor predicho por el modelo se aleja del valor real en experimentos con concentraciones de 0 y 5 gr de sustancia explosiva pura.



**Figura 93.** Valores predichos por el modelo de regresión de alcohol/TNT en estado puro

(a) Base de datos 1, (b) Base de datos 2

### 8.2.6. Modelos de Tarea 6: Regresión de Alcohol/Pólvora en Estado Puro y Mezcla de Pólvora

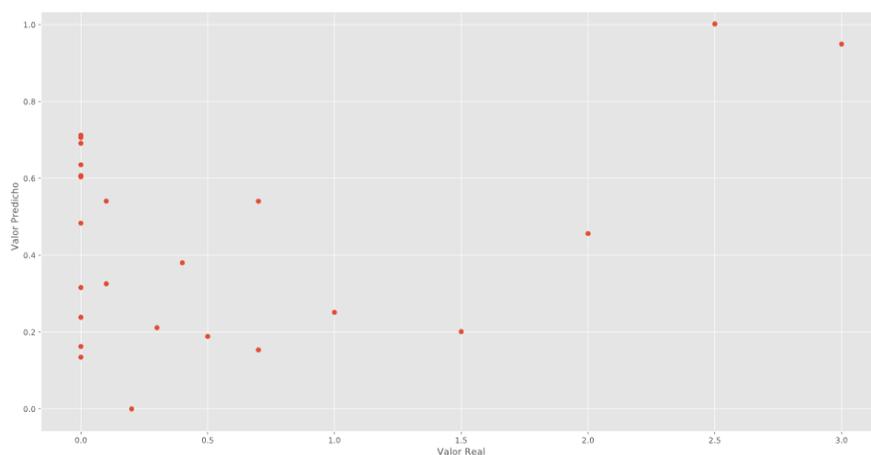
Los modelos generados para la Tarea 6 permiten identificar la cantidad pólvora (en estado puro o mezclado con jabón o pasta dental) del experimento analizado. En la **Tabla 33**, se observa el MSE y R2 de esta tarea para los modelos de regresión: PLS-R, red neuronal MLP y red neuronal profunda LSTM, con las cuales se determina que el modelo con mejor desempeño es el modelo PLS-R con un MSE de 0.61 y un R2 de 0.12 para los datos de prueba de la base de datos 1.

**Tabla 33**

*MSE y R2 de modelos de regresión de alcohol/pólvora en estado puro y mezcla de pólvora*

TAREA	BASE DE DATOS 1					
	PLS-R		MLP		LSTM	
	MSE	R2	MSE	R2	MSE	R2
Alcohol/Pólvora en estado puro y mezcla de pólvora	0.61	0.12	0.67	0.05	0.77	-0.09

Los resultados de la predicción de la cantidad de pólvora entre 0 y 3 gr con los datos de prueba para el modelo PLS-R de la base de datos 1 se muestra en la Figura 94, de la cual se determina que el valor predicho por el modelo se aleja del valor real, sin embargo, la concentración predicha por el modelo si aumenta a medida que la concentración real lo hace.



**Figura 94.** Valores predichos por el modelo de regresión de alcohol/pólvora en estado puro y mezcla de pólvora

### 8.2.7. Modelos de Tarea 7: Regresión de Alcohol/TNT en Estado Puro y Mezcla de TNT

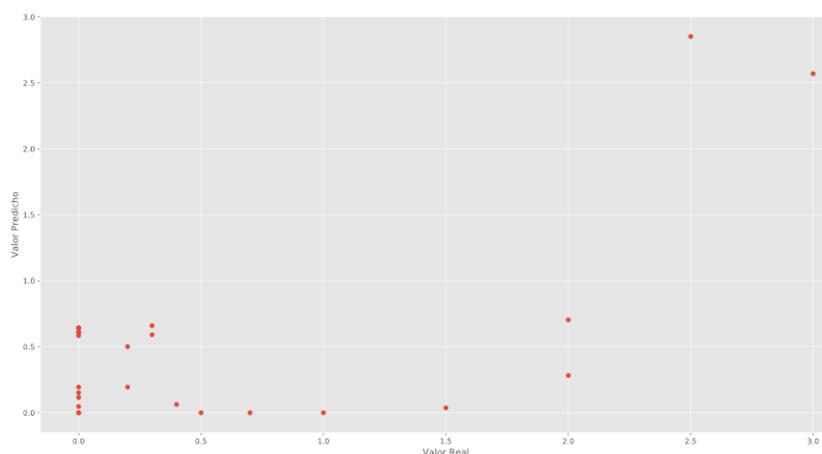
Los modelos generados para la Tarea 6 permiten identificar la cantidad TNT (en estado puro o mezclado con jabón o pasta dental) del experimento analizado. En la **Tabla 34**, se observa el MSE y R2 de esta tarea para los modelos de regresión: PLS-R, red neuronal MLP y red neuronal profunda LSTM, con las cuales se determina que el modelo con mejor desempeño es el modelo PLS-R con un MSE de 0.46 y un R2 de 0.40 para los datos de prueba de la base de datos 1.

**Tabla 34**

*MSE y R2 de modelos de regresión de alcohol/TNT en estado puro y mezcla de TNT*

TAREA	BASE DE DATOS 1					
	PLS-R		MLP		LSTM	
	MSE	R2	MSE	R2	MSE	R2
Alcohol/TNT en estado puro y mezcla de TNT	0.46	0.40	0.61	0.22	0.77	0.01

Los resultados de la predicción de la cantidad de TNT entre 0 y 3 gr con los datos de prueba para el modelo PLS-R de la base de datos 1 se muestra en la Figura 95, de la cual se determina que el valor predicho por el modelo se aleja del valor real en valores de concentración entre 0.5 y 2 gr, sin embargo, la concentración predicha por el modelo si aumenta a medida que la concentración real lo hace.



**Figura 95.** Valores predichos por el modelo de regresión de alcohol/TNT en estado puro y mezcla de TNT

### 8.2.8. Modelos de Tarea 8: Regresión de Alcohol/Pólvora en Estado Puro/TNT en Estado Puro

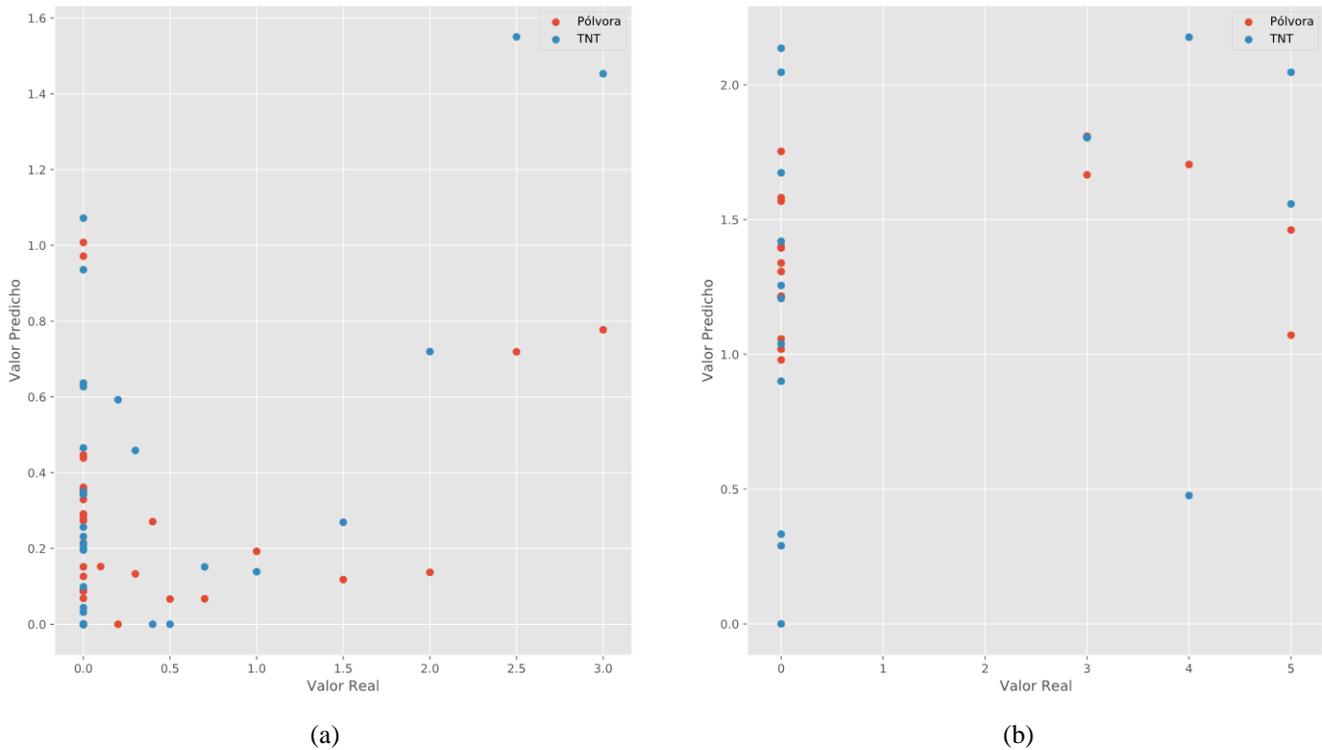
Los modelos generados para la Tarea 8 permiten identificar la cantidad de pólvora o TNT en estado puro del experimento analizado, su diferencia de los modelos generados para la Tarea 2 es que es un modelo multiclase. En la **Tabla 35**, se observa el MSE y R2 de esta tarea para los modelos de regresión: PLS-R, red neuronal MLP y red neuronal profunda LSTM, con las cuales se determina que el modelo con mejor desempeño para la base de datos 1 y 2 es el modelo PLS-R con un MSE de 0.46 y un R2 de 0.20, y con un MSE de 3.71 y un R2 de 0.10 para los datos de prueba de la base de datos 2.

**Tabla 35**

*MSE y R2 de modelos de regresión de alcohol/pólvora en estado puro/TNT en estado puro*

TAREA	BASE DE DATOS 1						BASE DE DATOS 2					
	PLS-R		MLP		LSTM		PLS-R		MLP		LSTM	
	MSE	R2	MSE	R2	MSE	R2	MSE	R2	MSE	R2	MSE	R2
Alcohol/Pólvora en estado puro/TNT en estado puro	0.47	0.20	0.58	0.02	0.62	-0.04	3.40	0.08	4.16	-0.12	5.27	-0.41

Los resultados de la predicción de la cantidad de pólvora o TNT en estado puro entre 0 y 3 gr con los datos de prueba para el modelo PLS-R de la base de datos 1 se muestra en la Figura 96a, de la cual se determina que el valor predicho por el modelo se aleja del valor real, además, la concentración predicha por el modelo si aumenta a medida que la concentración real lo hace. En el caso del modelo para la predicción de la cantidad de pólvora o TNT en estado puro entre 3 y 5 gr con los datos de prueba para el modelo PLS-R de la base de datos 2, se muestra en la Figura 96b, de la cual se determina que el valor predicho por el modelo también se aleja del valor real.

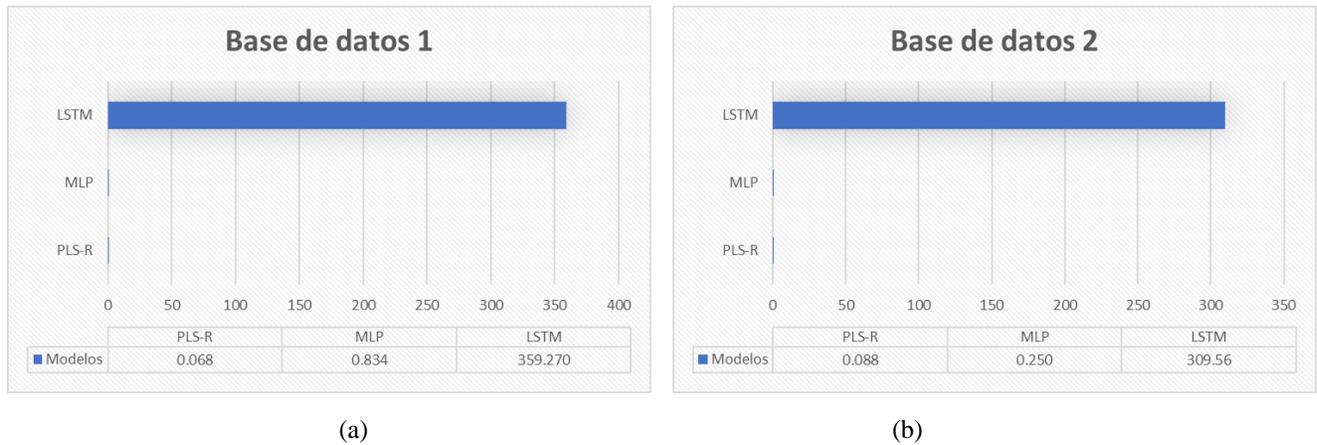


**Figura 96.** Valores predichos por el modelo de regresión de Alcohol/Pólvora en estado puro/TNT en estado puro

(a) Base de datos 1, (b) Base de datos 2

### 8.2.9. Tiempo de Entrenamiento de los Modelos de Regresión

En la Figura 97, se presenta el tiempo total de entrenamiento para los 8 modelos de regresión con cada uno de los algoritmos para la base de datos 1 (Figura 97a) y 2 (Figura 97b), con el cual se determina que el tiempo de entrenamiento de los modelos con el método PLS-R es inferior al tiempo empleado para el entrenamiento de los modelos generados el método de deep learning mediante la red neuronal profunda LSTM .



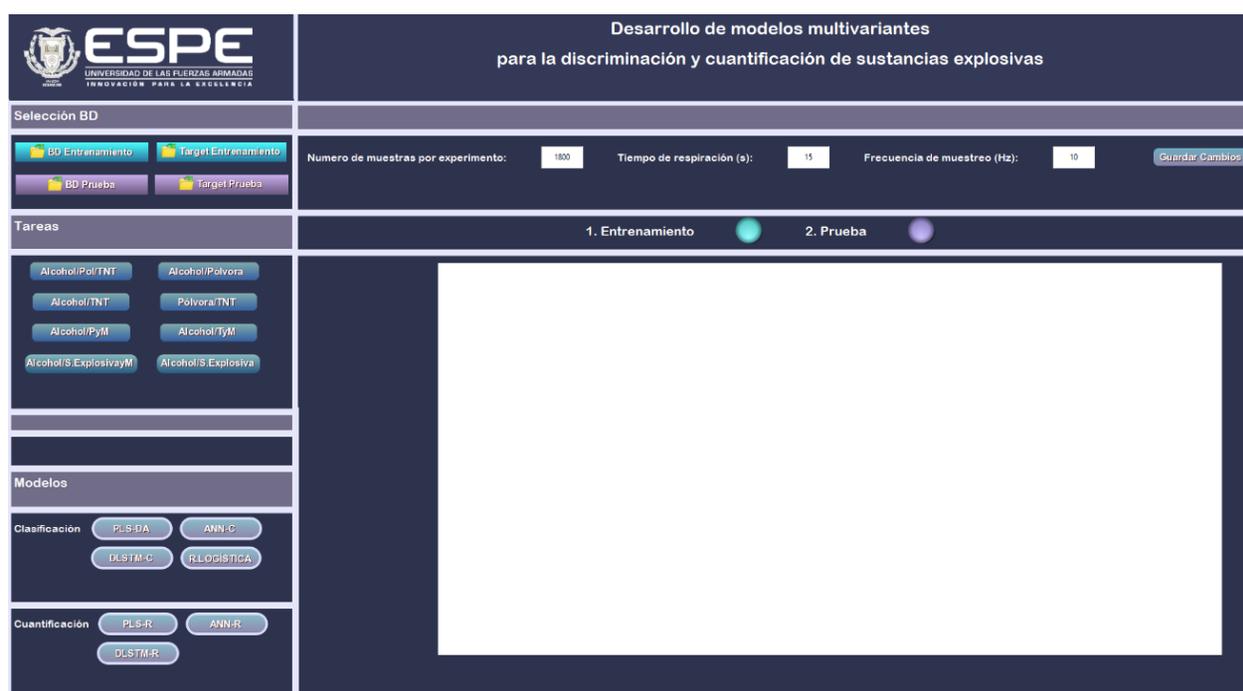
(a) (b)  
**Figura 97.** Tiempo de entrenamiento de modelos de regresión

(a) Base de datos 1, (b) Base de datos 2

Estos resultados se obtuvieron en un computador con un procesador Inter(R) Core (TM) i7-8550U, 12 GB RAM, ejecutando Spyder 3.3.6 en Windows 10.

## 8.3. Exploración del Desempeño de la Interfaz Gráfica de Usuario

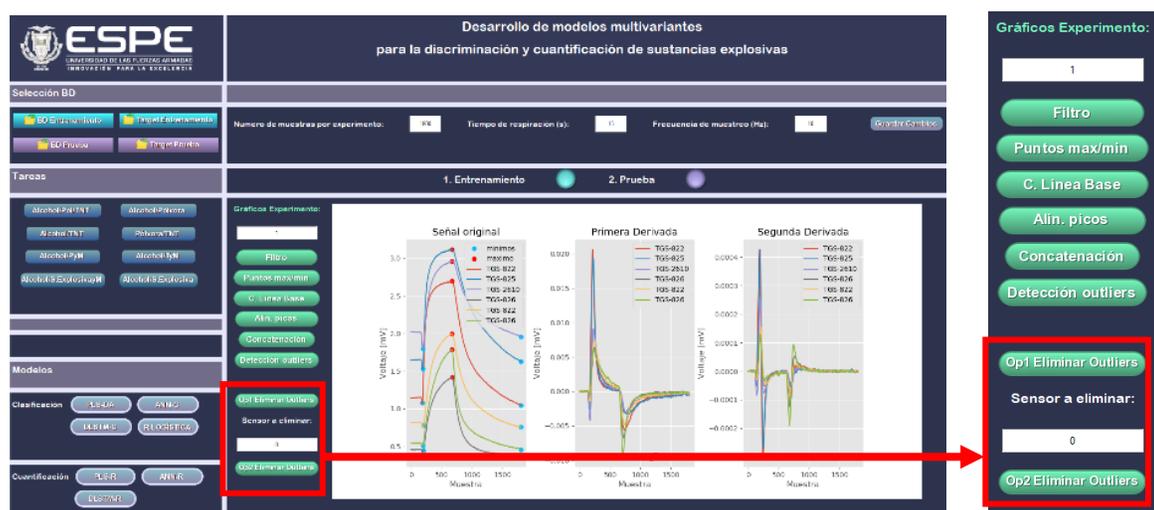
En base a la estructura presentada en el Capítulo 2, la interfaz resultante es presentada en la Figura 98. Esta interfaz permite utilizar los modelos generados para evaluarlos con nuevos datos de prueba o reentrenar los modelos con nuevos datos de entrenamiento. Su desempeño se indica a continuación.



**Figura 98.** Interfaz gráfica de usuario del proyecto de investigación

Para realizar el entrenamiento de un modelo, es necesario seleccionar la base de datos de entrenamiento, su target, la tarea que el modelo va a realizar, el algoritmo de regresión o clasificación, el número de muestras por experimento, el tiempo de respiración del prototipo e-nose y la frecuencia de muestreo de los datos. Con esa información y después de haber pulsado el botón

de entrenamiento, se realiza de forma automática el preprocesamiento de los datos, además que en la interfaz se puede observar el resultado de cada una de las etapas de preprocesamiento, como se observa en la Figura 99, mediante la selección del experimento que se desea visualizar y la etapa de la cual se quiere observar el resultado. Después de la detección de outliers, en la pantalla se presentan dos botones para la selección del método para eliminar las señales atípicas, mediante la opción 1, se eliminan únicamente estas señales y mediante la opción dos se elimina el sensor que se haya identificado como causante de la mayoría de los experimentos con outliers.



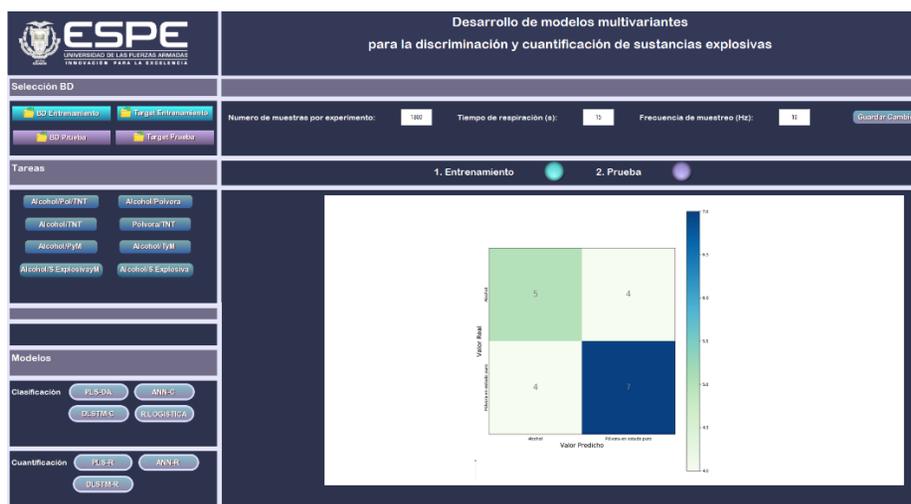
**Figura 99.** Preprocesamiento en la interfaz gráfica de usuario del proyecto de investigación

Al eliminarse las señales atípicas, se inicia el entrenamiento del modelo, finalizada esta etapa aparecerá una ventana emergente con los resultados del entrenamiento en función al tiempo que tomo este, el AUC de entrenamiento y la clase predicha por el modelo. El ejemplo presentado en la Figura 100, presenta el resultado de un modelo de clasificación de pólvora en estado puro, en la cual 0 representa un experimento no explosivo (alcohol) y 1 pólvora. Además, en la ventana principal se indicará la matriz de confusión resultante, como se observa en la Figura 101 para los modelos de clasificación y el valor predicho vs el real para los modelos de regresión. En el caso de

seleccionar no el entrenamiento sino la prueba de los modelos ya generados, se realizará automáticamente el preprocesamiento de las señales y se evaluará el desempeño de estos modelos con los datos de prueba, tal como se indicó en las figuras.

n de sustancias explosivas			
Tiempo de entrenamiento	AUC entrenamiento PLS-DA	y clasif real	y clasif predicho
3.5861005783081055	0.6717171717171717	[1 0 1 1 1 1 0 1 1 0 1 0 1 1 0 1 0 0 0 0]	[0 0 1 1 1 0 0 1 1 0 1 1 0 1 1 0 0 1 0 1]

**Figura 100.** Ventana emergente de resultado de entrenamiento de modelo PLS-DA



**Figura 101.** Matriz de confusión resultante del entrenamiento de modelo PLS-DA

## 8.4. Discusión de Resultados

En cuanto a los modelos de clasificación, el desempeño de los modelos de machine learning fue superior a los de deep learning para la base de datos actual de concentraciones de sustancia explosiva entre 3 y 5 gr, y en el caso de los modelos de la base de datos inicial, para concentraciones de sustancias explosiva entre 0 y 3gr los modelos de deep learning.

Para la mayoría de las tareas, los modelos de regresión con la técnica PLS-R ofrecieron mejores resultados. Sin embargo, el desempeño fue bajo con un  $R^2$  que no supero el 0.57 en el caso de la base de datos inicial y un  $R^2$  de 0.22 para la base de datos actual.

El desempeño de los modelos para la base de datos inicial con sustancias explosivas mezcladas con pasta dental o jabón fue inferior a los modelos generados con sustancias explosivas puras, razón por la cual para la elaboración de la base de datos actual no se tomaron en cuenta este tipo de experimentos.

Con los resultados de los modelos en función al tiempo de entrenamiento, se determinó que los modelos de deep learning tanto para tareas de clasificación como de regresión fueron los menos eficientes.

La interfaz gráfica de usuario permitió ejecutar en segundo plano la aplicación que integra tanto el preprocesamiento de los datos como los modelos. Además, se cuenta con la opción de reentrenar los modelos, por lo cual, es una herramienta útil para evaluar y seleccionar el modelo adecuado en función a la tarea que se desea que el prototipo realice y las características del prototipo.

En los modelos de machine learning el desempeño de los modelos en los datos de entrenamiento fue mucho mejor que el de los datos de prueba, ya que el  $AUC \gg 0.5$  para los datos de entrenamiento y  $AUC \approx 0.5$  para los datos de prueba. Esto puede deberse a que los datos con los que fueron entrenados los modelos de esta base de datos no son lo suficientemente representativos a los datos de prueba, lo cual se podría solucionar al adquirir más datos de entrenamiento que proporcionen más información al algoritmo para que pueda ser utilizado durante el proceso de entrenamiento y permita obtener modelos que no se sobreajusten a estos datos.

# CAPÍTULO 9

## Conclusiones y Recomendaciones

En este capítulo se presentan las conclusiones y recomendaciones del proyecto de investigación que aportaran al desarrollo de investigaciones futuras.

### 9.1. Conclusiones

Se generaron ocho modelos por cada base de datos para tareas de clasificación y cuantificación, de los cuales se evaluó su desempeño con sustancias explosivas puras y mezclas con pasta dental o jabón para la base de datos 1 y únicamente con sustancias explosivas puras para la base de datos 2. Los modelos con mejor desempeño para la clasificación de sustancias explosivas (puras y mezclas con pasta dental o jabón) y alcohol para la base de datos 1, fueron los generados con la red neuronal profunda LSTM, con una sensibilidad de 0.54 y especificidad de 0.73 en la clasificación entre sustancias explosivas y alcohol, una sensibilidad de 0.62 y especificidad de 0.64 en la clasificación entre pólvora y alcohol, y una sensibilidad de 0.62 y especificidad de 0.73 en la clasificación entre TNT y alcohol. Los modelos con mejor desempeño para la clasificación de sustancias explosivas puras y alcohol para la base de datos 1 y 2, fueron los generados con la red neuronal MLP, con una sensibilidad de 0.52 y especificidad de 0.55 para el primer caso, y una sensibilidad de 0.91 y especificidad de 0.67 para el segundo. En la clasificación de pólvora en estado puro y alcohol el modelo de regresión logística fue aquel con mejor desempeño para la base de datos 1, con una sensibilidad de 0.45 y especificidad de 0.73, y la red neuronal MLP para la base de datos 2 con una sensibilidad de 0.80 y especificidad y precisión de 0.67. Para la clasificación de TNT en

estado puro y alcohol, el modelo con mejor desempeño fue la red neuronal profunda LTMS para la base de datos 1 con una sensibilidad de 0.70 y especificidad de 0.73, y para la base de datos 2 el modelo de regresión logística con una sensibilidad de 0.67 y especificidad de 0.83. En la clasificación de pólvora y TNT en estado puro el modelo con mejor desempeño fue la red neuronal profunda LSTM para la base de datos 1 y 2, cuya sensibilidad en la detección de pólvora fue de 0.82 y especificidad de 0.50, en la detección de TNT una sensibilidad de 0.50 y especificidad de 0.82 para la base de datos 1, y para la base de datos 2 una sensibilidad en la detección de pólvora de 0.60 y especificidad de 0.50 y en la detección de TNT una sensibilidad de 0.50 y especificidad de 0.60. Para la clasificación entre alcohol, TNT y pólvora en estado puro el modelo con mejor desempeño fue el método de regresión logística para la base de datos 1 con una sensibilidad en la detección de alcohol de 0.64 y especificidad de 0.31, en la detección de pólvora una sensibilidad de 0.18 y especificidad de 0.59 y en la detección de TNT una sensibilidad de 0.30 y especificidad de 0.8 y para la base de datos 2 el modelo de redes neuronal MLP, con una sensibilidad en la detección de alcohol de 0.83 y especificidad de 0.80, en la detección de pólvora una sensibilidad de 0.40 y especificidad de 0.70 y en la detección de TNT una sensibilidad de 0.33 y especificidad de 0.64. Los modelos con mejor desempeño para la cuantificación de sustancias explosivas (puras y mezclas con pasta dental o jabón) y alcohol para la base de datos 1, fueron los generados con PLS-R, con un  $R^2$  de 0.38 para la cuantificación de sustancias explosivas y alcohol,  $R^2$  de 0.12 para la cuantificación de pólvora, y  $R^2$  de 0.40 para cuantificación de TNT. Los modelos con mejor desempeño para la cuantificación de sustancias explosivas puras y alcohol para la base de datos 1 y 2, fueron los generados con la red neuronal MLP para la base de datos 1 con un  $R^2$  de 0.47, y el generado con PLS-R para la base de datos 2 con un  $R^2$  de 0.17. Para la cuantificación de pólvora el modelo con mejor desempeño fue el modelo PLS-R con un  $R^2$  de 0.31 para la base de datos 1 y de

0.09 para la base de datos 2. Para la cuantificación de pólvora el modelo con mejor desempeño fue el modelo PLS-R con un R2 de 0.31 para la base de datos 1 y de 0.09 para la base de datos 2, para la cuantificación de TNT el modelo con mejor desempeño fue el modelo PLS-R con un R2 de 0.53 para la base de datos 1 y de 0.22 para la base de datos 2. Los modelos con mejor desempeño para la cuantificación de pólvora y TNT para la base de datos 1 y 2, fueron los generados con la red neuronal MLP para la base de datos 1 con un R2 de 0.57, y el generado con la red neuronal profunda LSTM para la base de datos 2 con un R2 de -0.02. En el caso del modelo multiclase para la cuantificación de TNT y polvora el modelo con mejor desempeño fue el modelo PLS-R con un R2 de 0.20 para la base de datos 1 y de 0.08 para la base de datos 2.

Se analizaron las bases de datos generadas con el prototipo e-nose desarrollado en (Vallejo & Zurita, 2017) y optimizado en (Jacome, 2019), de las cuales se identificaron 6 señales pertenecientes a la respuesta de sensores químicos TGS, con 650 muestras para el primer caso y 1800 muestras para el segundo. Estas señales contaban con ruido, principalmente en el sensor 5 para la base de datos inicial, sus líneas base no se encontraban en un mismo punto y los picos de las señales estaban desalineados. Además, se evaluaron las bases de datos mediante el análisis de componentes principales (PCA), con el cual se realizó un análisis exploratorio y con ayuda del valor  $T^2$  de Hotelling, se determinó que el sensor 5 en la primera base de datos genero la mayoría de las señales atípicas, contrario a la base de datos actual, en la cual no se identificó un sensor específico causante de estas señales.

Los modelos de clasificación y cuantificación de sustancias explosivas se generaron mediante las técnicas lineales de mínimos cuadrados parciales y regresión logística, y mediante dos

técnicas no lineales: una red neuronal artificial clásica perceptrón multicapa con una capa oculta y una red neuronal profunda LSTM.

Mediante el análisis comparativo de los modelos se determinó que los modelos de machine learning generados con la base de datos actual con 36 experimentos de entrenamiento de entre 3 y 5gr de sustancias explosivas puras, tuvieron un mejor desempeño que los modelos de deep learning, con un AUC entre 0.47 y 0.95 y un tiempo de entrenamiento entre 0.019 y 0.147 segundos. Sin embargo, para la base de datos inicial con 89 experimentos de entrenamiento de entre 0.1 y 3gr de sustancias explosivas puras y mezcladas con pasta dental o jabón, los modelos con mejor desempeño en clasificación fueron los modelos de deep learning, con un AUC entre 0.55 y 0.73, y un tiempo de entrenamiento de 240.915 segundos. Por lo cual se concluye que la cantidad de datos y el tiempo disponible para la elaboración de los modelos es importante a la hora de seleccionar el tipo de técnica a utilizar.

Se determinó que el prototipo tiene problemas al cuantificar la concentración de la sustancia explosiva de interés ya que el desempeño de los modelos de regresión en función del  $R^2$  no supero el 0.57 para la base de datos inicial y el 0.22 para la actual. Además, los modelos con mejor desempeño fueron en su mayoría aquellos generados mediante la técnica de mínimos cuadrados parciales.

El desempeño de los modelos de clasificación y cuantificación de TNT fue superior a los de pólvora. Por consiguiente, se espera que bajo las condiciones actuales del prototipo y con concentraciones entre 3 y 5gr de sustancias explosivas este sea capaz de clasificar correctamente en mayor porcentaje observaciones de TNT que de pólvora.

Al analizar el poder predictivo del prototipo e-nose con los experimentos de la base de datos 2 se determinó que el desempeño de los modelos de clasificación fue superior al generado con la base de datos 1, a pesar de que la cantidad de sustancia dopante en el prototipo disminuyó en un 50%.

Al comparar el desempeño entre los modelos realizados con bajas concentraciones de sustancia explosiva (experimentos de la base de datos 1 con concentraciones entre 0.1 y 3gr de pólvora y TNT) de los realizados con mayores concentraciones (experimentos de la base de datos 2 con concentraciones entre 3 y 5gr de pólvora y TNT), se pudo observar que se obtuvo mejores resultados al aumentar el nivel de concentración de sustancia. Por lo tanto, mientras mayor sea la concentración de sustancia explosiva que se desee clasificar el prototipo tendrá mayores posibilidades de realizar una clasificación correcta.

La aplicación desarrollada cuenta con una interfaz gráfica de usuario que permite ejecutar en segundo plano el preprocesamiento de las señales obtenidas con el prototipo e-nose, detectar outliers y utilizar los modelos generados en el presente proyecto de investigación para tareas de clasificación o cuantificación de sustancias explosivas. Además, permite reentrenar los modelos en base a las condiciones actuales del prototipo, es decir que, si existen modificaciones en el prototipo que cambien las condiciones con las que fueron entrenados los modelos inicialmente se podrán reentrenar con las condiciones actuales y evaluarlos, siendo un importante aporte para proyectos futuros en los cuales se podrá utilizar esta aplicación para seleccionar el modelo con mejor desempeño acorde a la tarea que se desea realizar y las condiciones actuales del prototipo.

## 9.2. Recomendaciones

Se recomienda generar una base de datos más grande para evaluar el desempeño de los modelos de deep learning, ya que se aconseja utilizar este tipo de técnicas cuando existe una cantidad considerable de datos para su entrenamiento.

La etapa de reducción de dimensionalidad demostró disminuir el tiempo de entrenamiento y la información redundante existente en los datos de entrenamiento, por lo tanto, se recomienda explorar otras técnicas de reducción de dimensionalidad y seleccionar la óptima para cada uno de los algoritmos utilizados.

Con una base de datos de cientos o miles de experimentos, sería importante probar otro tipo de arquitecturas de redes neuronales profundas recomendadas para series temporales, como la combinación de capas LSTM para la predicción de la serie temporal y capas convolucionales para la extracción de las características de entrada.

La temperatura de la cámara de sensores, al tener una alta incidencia en la respuesta de los sensores, debe ser controlada de forma adecuada, por lo tanto, se recomienda verificar que las condiciones de temperatura del prototipo e-nose sean similares tanto en la base de datos de entrenamiento como de prueba y en caso de ser necesario incluir la temperatura como característica de entrada para futuros modelos.

# REFERENCIAS

- Aggarwal, C. C. (2013). Probabilistic and Statistical Models for Outlier Detection. En *Outlier Analysis* (Springer, pp. 41–74). New York. [https://doi.org/https://doi.org/10.1007/978-1-4614-6396-2\\_2](https://doi.org/10.1007/978-1-4614-6396-2_2)
- Alpaydm, E. (2012). *Introduction to Machine Learning, 2nd ed. Natural Language Engineering* (The MIT Pr). Cambridge, MA. <https://doi.org/10.1017/S1351324912000290>
- B. Lakshmi, K.Parish Venkata Kumar, Dr. K.Nageswara Rao, S. B. (2013). Simulation of Artificial Noses for the Automated Detection and Classification of Organic Compounds. *International Journal of Computer Science and Information Technologies*, 4, 233–237.
- Brownlee, J. (2018). A Gentle Introduction to k-fold Cross-Validation. Recuperado el 8 de mayo de 2019, de <https://machinelearningmastery.com/k-fold-cross-validation/>
- Brownlee, Jason. (2016). *Master Machine Learning Algorithms. Discover How They Work and Implement Them From Scratch. Machine Learning Mastery With Python.*
- Brownlee, Jason. (2017). Long Short-Term Memory Networks With Python. *Machine Learning Mastery With Python.*
- Burkov, A. (2019). The Hundred-Page Machine Learning Book-Andriy Burkov. *Expert Systems*. <https://doi.org/10.1111/j.1468-0394.1988.tb00341.x>
- Chawla, N. V., Bowyer, K. W., Hall, L. O., & Kegelmeyer, W. P. (2002). SMOTE: Synthetic minority over-sampling technique. *Journal of Artificial Intelligence Research*. <https://doi.org/10.1613/jair.953>
- Chollet, F. (2018). *Deep Learning with Python*. Manning Publications Co.
- Dunn, K. G. (2019). *Process Improvement using Data*. Recuperado de [learnche.mcmaster.ca/pid](http://learnche.mcmaster.ca/pid)

- Eriksson, L., Johansson, E., Kettaneh-Wold, N., Trygg, C., Wikström, C., & Wold, S. (2006). *Multi- and Megavariate Data Analysis Part I. Basic Principles and Applications*. Umetrics AB.
- Espinosa, C., & Venegas, C. (2017). *Optimización e integración de una nariz electrónica autónoma embebida en un sistema robótico para la identificación de sustancias explosivas como TNT y pólvora base doble en ambientes controlados*. Universidad de Fuerzas Armadas ESPE.
- Flach, P. (2012). *MACHINE LEARNING The Art and Science of Algorithms that Make Sense of Data*. Cambridge, New York, Melbourne, Madrid, Cape Town, Singapore, São Paulo, Delhi, Mexico City: Cambridge University Press.
- Fort, G., & Lambert-Lacroix, S. (2005). Classification using partial least squares with penalized logistic regression. *Bioinformatics*. <https://doi.org/10.1093/bioinformatics/bti114>
- Fu, J., Li, G., Qin, Y., & Freeman, W. J. (2007). A pattern recognition method for electronic noses based on an olfactory neural network. *Sensors and Actuators, B: Chemical*. <https://doi.org/10.1016/j.snb.2007.02.058>
- Gardner, J. W. (2013). Review of Conventional Electronic Noses and Their Possible Application to the Detection of Explosives. En *Electronic Noses & Sensors for the Detection of Explosives*. [https://doi.org/10.1007/978-1-4020-2800-7\\_1](https://doi.org/10.1007/978-1-4020-2800-7_1)
- Geladi, P., & Kowalski, B. R. (1986). Partial least-squares regression: a tutorial. *Analytica Chimica Acta*. [https://doi.org/10.1016/0003-2670\(86\)80028-9](https://doi.org/10.1016/0003-2670(86)80028-9)
- Geron, A. (2017). *Hands-On Machine Learning With Scikit-Learn & Tensor Flow*. *Hands-on Machine Learning with Scikit-Learn and TensorFlow*. <https://doi.org/10.3389/fninf.2014.00014>
- Guaman, A. V., Lopez, P., & Torres-Tello, J. (2019). Multivariate Discrimination Model for TNT and Gunpowder Using an Electronic Nose Prototype: A Proof of Concept.

[https://doi.org/10.1007/978-3-030-11890-7\\_28](https://doi.org/10.1007/978-3-030-11890-7_28)

- Jacome, S. (2019). *Repotenciación de la nariz electrónica con dopaje automático del proyecto de investigación 2016-PIC-009*. Universidad de las Fuerzas Armadas-ESPE.
- Lee, L. C., Liong, C. Y., & Jemain, A. A. (2018). Partial least squares-discriminant analysis (PLS-DA) for classification of high-dimensional (HD) data: A review of contemporary practice strategies and knowledge gaps. *Analyst*. <https://doi.org/10.1039/c8an00599k>
- López Hernández, I. P. (2016). *Desarrollo de un prototipo electrónico de sensado químico, para la detección de trinitrotolueno (TNT) y pólvora base doble en un ambiente controlado*. Universidad de las Fuerzas Armadas-ESPE.
- López, P., Triviño, R., Calderón, D., Arcentales, A., & Guamán, A. V. (2017). Electronic nose prototype for explosive detection. En *2017 CHILEAN Conference on Electrical, Electronics Engineering, Information and Communication Technologies, CHILECON 2017 - Proceedings*. <https://doi.org/10.1109/CHILECON.2017.8229657>
- Lundh, F. (1999). An Introduction to Tkinter. *Review Literature And Arts Of The Americas*.
- Ministerio del Interior. (2019). Guías y canes, dupla estratégica del Grupo de Intervención y Rescate. Recuperado de <https://www.ministeriointerior.gob.ec/guias-y-canec-dupla-estrategica-del-grupo-de-intervencion-y-rescate/>
- Mujica, L. E., Rodellar, J., Fernández, A., & Güemes, A. (2011). Q-statistic and t2-statistic pca-based measures for damage assessment in structures. *Structural Health Monitoring*. <https://doi.org/10.1177/14759217110388972>
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., ... Duchesnay, É. (2011). Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*.
- Rao, D., & McMahan, B. (2019). *Natural Language Processing with PyTorch*. Sebastopol: O'Reilly

Media.

Raschka, S., & Mirjalili, V. (2017). *Python Machine Learning - Second Edition*. *Materials Research Bulletin*. [https://doi.org/10.1016/0025-5408\(96\)80018-3](https://doi.org/10.1016/0025-5408(96)80018-3)

Rojas, R. (1996). *Neural Networks: A Systematic Introduction*. Berlin, Heidelberg: Springer-Verlag.

Salazar, J. (2018). *Creación de un modelo de calibración multivariante para determinar el límite de detección de un prototipo nariz electrónica para medición de sustancias explosivas*. Universidad de las Fuerzas Armadas-ESPE.

Samarasinghe, S. (2006). *Neural Networks for Applied Sciences and Engineering*. *Neural Networks for Applied Sciences and Engineering*. Auerbach Publications.

Schafer, R. W. (2011). What is a savitzky-golay filter? *IEEE Signal Processing Magazine*. <https://doi.org/10.1109/MSP.2011.941097>

Trogler, W. C. (2013). Luminescent Inorganic Polymer Sensors for Vapour Phase and Aqueous Detection of TNT. En *Electronic Noses & Sensors for the Detection of Explosives*. [https://doi.org/10.1007/978-1-4020-2800-7\\_3](https://doi.org/10.1007/978-1-4020-2800-7_3)

Vallejo, Z., & Zurita, D. (2017). *Diseño e implementación de un sistema mecatrónico portable con dopaje automático para detección de muestras explosivas*. Universidad de las Fuerzas Armadas-ESPE.

Vallejos Brito, D. A. (2019). Evaluación del desempeño de un UAV, equipado con dispositivos de sensado para la identificación de componentes químicos volátiles empleados en el procesamiento ilícito de sustancias.

Walt, D., & Sternfeld, T. (2013). Optical Microsensor Arrays for Explosives Detection. En J. Gardner & J. Yinon (Eds.), *Electronic Noses & Sensors for the Detection of Explosives* (pp. 81–92). Dordrech: Springer. [https://doi.org/https://doi.org/10.1007/1-4020-2319-7\\_6](https://doi.org/https://doi.org/10.1007/1-4020-2319-7_6)

- Wang, Q. Q., Liu, K., Zhao, H., Ge, C. H., & Huang, Z. W. (2012). Detection of explosives with laser-induced breakdown spectroscopy. *Frontiers of Physics*. <https://doi.org/10.1007/s11467-012-0272-x>
- Wold, S., Sjöström, M., & Eriksson, L. (2001). PLS-regression: A basic tool of chemometrics. En *Chemometrics and Intelligent Laboratory Systems*. [https://doi.org/10.1016/S0169-7439\(01\)00155-1](https://doi.org/10.1016/S0169-7439(01)00155-1)
- Zhang, L., & Peng, X. (2016). Time series estimation of gas sensor baseline drift using ARMA and Kalman based models. *Sensor Review*. <https://doi.org/10.1108/SR-05-2015-0073>
- Zuñiga, D. (2018). *Filtro Savitzky y Golay de Segundo Grado en datos de Tomografía de Resistividad Eléctrica 2D*. Universidad Nacional Autónoma de México.