



Facultad de Ciencias  
Instituto de Estadística  
Ingeniería en Estadística

# Diseño de un marcador somático artificial para un sistema inteligente de clasificación de riesgo ocupacional en faenas mineras

**Trabajo final presentado por:**  
Javiera Ignacia Arraño Arraño

**Profesor Guía**

Harvey Rosas, Ph.D.

**Profesor Co-Guía**

Daniel Cabrera-Paniagua, Ph.D.

**Proyecto de titulación para optar al:**

grado académico de: *Licenciado en Estadística*

título profesional de: *Ingeniero en Estadística*

minor en: *Minería de datos*

Valparaíso, Chile, 15 de diciembre de 2021

# Resumen

La minería presenta una serie de peligros de ámbito tanto físico derivados de la posible exposición a condiciones que puedan poner en riesgo la salud de los trabajadores, la necesidad por identificarlos y evaluarlos puede generar una disminución de incidentes como de costos que logre impactar el negocio.

Durante los últimos años la utilización de técnicas de *machine learning* ha tomado una relevancia significativa para llevar a cabo diversos estudios en distintas áreas, pero en seguridad industrial o riesgo ocupacional ha sido algo bastante reciente. El presente trabajo de titulación en conjunto de la empresa Previsis, se propone una técnica de *machine learning*, que se incorpora un marcador somático artificial para la clasificación del riesgo ocupacional en el contexto de faenas mineras. Se creó una técnica de reclasificación buscando mejorar la clasificación de la técnica *Bosque aleatorio* de *machine learning*, para así realizar una comparación entre ambos casos, sin reclasificación y con reclasificación. Se obtuvieron resultados positivos al aplicar la reclasificación a cada clase del conjunto de datos, generando una disminución al conjunto de mal clasificados.

# Algunas Palabras

Este trabajo de titulación culmina el término de una etapa que ha marcado mi vida desde el primer momento.

En el año 2016 tuve que emprender un camino que no tenía idea como iba a ser, y como lo iba a enfrentar. Viaje a Valparaíso con una maleta llena de sueños y ganas de salir adelante.

En primer lugar, quiero agradecer a mi madre Marcela Arraño que me ha apoyado incondicionalmente y si hoy estoy terminando esta etapa es gracias a ella, a mis abuelos Silvia y Manuel personas que me han enseñado que a pesar de que el camino se ponga muy difícil nunca hay que rendirse, a mi hermano Sebastián que ha sido mi razón de vivir y de superarme siempre para poder entregarle lo mejor.

Agradecer a mi tía Elsa Alejandra que me ha ayudado hasta el día de hoy, a mi primo Alejandro más que un primo un hermano, agradecer por su gran apoyo, toda una vida batallando junto a él, enseñándome que siempre se puede lograr lo que se lucha por obtener.

Agradecer a mis amigos que fueron partícipes de esta etapa, pero no puedo dejar de mencionar a mi amiga Catalina Figueroa que a pesar de la distancia ha sido un apoyo enorme, agradezco a la vida de poner a una gran persona en mi camino que nunca me ha dejado sola, que ha sido testigo de mis peores batallas, sabe todo lo que luche por llegar a este punto, lo mejor que me pudo dejar la universidad es su amistad.

Gracias a mis profesores guías por confiar en mí, y darme esta oportunidad.

# Confidencialidad

La realización de este trabajo de título está en el marco de un acuerdo de confidencialidad, el cual fue firmado por Javiera Arraño (estudiante), Harvey Rosas (profesor guía) y Daniel Cabrera (profesor co-guía) con la empresa Prevsis para hacer entrega de la información y conjunto de datos. Por lo tanto, los resultados obtenidos en este proyecto de título están resguardados bajo secreto industrial.

# Índice general

<b>Resumen</b>	<b>2</b>
<b>Algunas Palabras</b>	<b>3</b>
<b>Confidencialidad</b>	<b>4</b>
<b>1. Estado del arte</b>	<b>10</b>
<b>2. Marco teórico</b>	<b>13</b>
2.1. Riesgo Ocupacional . . . . .	13
2.2. Hipótesis del marcador somático . . . . .	13
2.3. Machine Learning . . . . .	16
2.4. Bosque Aleatorio . . . . .	16
2.5. Validación cruzada . . . . .	19
2.5.1. Validación cruzada k-Fold . . . . .	19
2.6. Validación cruzada estratificada de K-Fold . . . . .	20
2.7. Métricas de desempeño . . . . .	20
<b>3. Metodología</b>	<b>22</b>
3.1. Conjunto de datos . . . . .	23
3.1.1. Variables de análisis . . . . .	23
3.2. Variable respuesta . . . . .	24
3.3. Técnicas de clasificación . . . . .	24
3.4. Segmentación del conjunto . . . . .	25
3.5. Re-Clasificación . . . . .	25
3.5.1. Umbrales actualizados . . . . .	27
3.6. Marcador somático artificial . . . . .	29
<b>4. Resultados</b>	<b>31</b>
4.1. Clasificación . . . . .	31
4.1.1. Primera clasificación Random Forest Clase “Alerta ” . . . . .	32
4.1.2. Clasificación Random Forest Clase “Alerta” al conjunto de prueba	42
4.1.3. <b>Segunda clasificación Random Forest Clase “Alerta 1”</b> .	45
4.1.4. Clasificación Random Forest Clase “Alerta 1” al conjunto de prueba	53
4.1.5. Tercera clasificación clase “Alerta 2” . . . . .	55
<b>5. Conclusiones</b>	<b>59</b>

# Índice de figuras

2.1. Ilustración de fases de la toma de decisiones . . . . .	15
2.2. Validación cruzada . . . . .	19
2.3. Matriz de confusión . . . . .	21
3.1. Diagrama Umbrales . . . . .	27
3.2. Diagrama reclasificación . . . . .	28
3.3. Incorporación marcador somático . . . . .	30
4.1. División conjunto de datos . . . . .	31
4.2. Matriz de confusión clase alerta sin incluir umbrales . . . . .	33
4.3. Histograma clase 1 mal clasificados clase alerta . . . . .	35
4.4. Histograma clase 0 mal clasificados clase alerta . . . . .	36
4.5. Histograma distancias clase 1 mal clasificados Alerta . . . . .	37
4.6. Histograma distancias clase 0 mal clasificados Alerta . . . . .	37
4.7. Diagrama reclasificación con umbral fijo Alerta . . . . .	38
4.8. Matriz de confusión umbrales fijos Alerta . . . . .	39
4.9. Matriz de confusión umbrales actualizados Alerta . . . . .	40
4.10. Matriz de confusión clase Alerta conjunto de prueba . . . . .	42
4.11. Matriz de confusión de clase Alerta con umbrales actualizados aplicados a conjunto de prueba . . . . .	44
4.12. Matriz de confusión sin umbrales Alerta 1 . . . . .	46
4.13. Histograma clase 1 mal clasificados clase Alerta 1 . . . . .	47
4.14. Histograma clase 0 mal clasificados clase Alerta 1 . . . . .	48
4.15. Histograma distancias Alerta 1 . . . . .	49
4.16. Histograma distancias Alerta 1 . . . . .	49
4.17. Matriz de confusión conjunto de validación con umbrales fijos Alerta 1 . . . . .	50
4.18. Matriz de confusión conjunto de validación con umbrales actualizados Alerta 1 . . . . .	51
4.19. Matriz de confusión conjunto prueba para clase Alerta 1 . . . . .	53
4.20. Matriz de confusión para conjunto de prueba con umbrales actualizados Alerta 1 . . . . .	54
4.21. Matriz de confusión Alerta 2 conjunto validación . . . . .	56
4.22. Matriz de confusión Alerta 2 conjunto prueba . . . . .	57

# Índice de cuadros

3.1. Cantidad de características por nivel . . . . .	24
3.2. Clasificador Random Forest . . . . .	25
3.3. Escenarios de clasificación . . . . .	29
4.1. Clase Alerta . . . . .	32
4.2. División conjunto de datos . . . . .	32
4.3. Métricas clase alerta conjunto de validación . . . . .	34
4.4. Métricas Alerta conjunto validación con reclasificación umbrales fijos Alerta . . . . .	40
4.5. Métricas obtenidas aplicando reclasificación . . . . .	41
4.6. Métricas clase Alerta conjunto de prueba sin reclasificación . . . . .	43
4.7. Métricas clase Alerta de conjunto de prueba con reclasificación . . . . .	45
4.8. Clase Alerta 1 . . . . .	45
4.9. Métricas clase Alerta 1 conjunto validación sin reclasificación . . . . .	47
4.10. Métricas Alerta 1 conjunto validación con reclasificación . . . . .	51
4.11. Métricas Alerta 1 conjunto validación con reclasificación . . . . .	52
4.12. Clase Alerta 2 . . . . .	55
4.13. División conjunto de datos para clase Alerta 2 . . . . .	55
4.14. Métricas Alerta 2 conjunto validación . . . . .	57
4.15. Métricas Alerta 2 conjunto de prueba . . . . .	58

# Introducción

Los accidentes laborales se han vuelto un factor de suma relevancia y consideración, ya que muchas veces se pierden vidas por la falta de control y gestión de los peligros y riesgos presentes en la industria. La organización Internacional del Trabajo (OIT) estima que se llegan a producir anualmente alrededor de 374 millones de accidentes no mortales en el trabajo, y que cerca de 7600 personas mueren por día en el mundo a raíz de estos accidentes laborales. Según los datos presentados por la OIT más de un 40 % de estos accidentes no mortales son en trabajadores entre 18 y 24 años [Sarkar et al. \(2019b\)](#).

El sector de la minería es un sector donde los trabajadores están expuestos constantemente a cualquier tipo de riesgos y accidentes, por lo que incrementar la eficiencia y la productividad de procesos que minimicen estos riesgos sería evidentemente beneficioso, tanto para los trabajadores como para las industrias [Sámano-Ríos et al. \(2019\)](#).

Como incidente se considera cualquier evento relacionado con el trabajo que pueda provocar algún daño o deterioro a la salud. Estos incidentes son causados por múltiples factores, por lo cual se ha buscado constantemente formas de comprender las situaciones y los factores que afectan la ocurrencia y gravedad de ellos. Asimismo, es deseable mejorar la precisión de predicción de la probabilidad de futuros accidentes.

El poder reducir la accidentalidad de trabajadores, es un propósito que muchas veces se busca utilizando algoritmos de inteligencia artificial, ya que así se pueden estudiar las razones porque se produjeron los accidentes.

El presente estudio contribuye a la escasa literatura actual sobre análisis de incidentes ocupacionales en faenas mineras, buscando implementar un marcador somático artificial dentro de una técnica de *machine learning* la cual, en este caso es bosque aleatorio, con el fin de analizar la toma de decisiones de un agente para una posible reclasificación. La hipótesis sugiere que la incorporación de este marcador somático artificial dentro de la técnica de *machine learning* mejora el desempeño de la clasificación.

Como objetivo general de esta investigación se basa en desarrollar una técnica *machine learning* que incorpore un marcador somático artificial para la clasificación del riesgo ocupacional en el dominio de minería, a través de esto se desprende: (i) examinar variables asociadas al riesgo ocupacional en una faena minera, (ii) diseñar un marcador somático artificial para técnicas de *machine learning*, (iii) integrar las variables de riesgo ocupacional con un marcador somático artificial en una técnica *machine learning* en el contexto del riesgo ocupacional, (iv) analizar los resultados experimentales derivados de la evaluación de la técnica de aprendizaje automático en el contexto de riesgo ocupacional.

El trabajo se organiza de la siguiente forma: en el capítulo 1 se presenta el estado de arte, donde exponen diversos estudios relacionado al uso de técnicas *machine learning* en el área del riesgo ocupacional, junto con el uso de un marcador somático artificial. En el capítulo 2 se encuentra el marco teórico, en el cual se abordaran los temas relacionados a riesgo ocupacional, marcador somático, bosque aleatorio, entre otros. En el capítulo 3, está la metodología utilizada en el trabajo. En el capítulo 4 los resultados obtenidos y finalmente en el capítulo 5 las conclusiones y observaciones realizadas.

# Capítulo 1

## Estado del arte

Los accidentes laborales en la industria son un factor de consideración, puesto que se puede llegar a perder muchas vidas debido a una falta de control y gestión de los peligros y riesgos. La utilización de aprendizaje automático de máquinas (ML), es cada vez más frecuente en diferentes áreas, pero en la seguridad industrial viene siendo bastante reciente, y es en esta área donde se utilizan grandes conjuntos de datos, importación relevante basada en actividades de los trabajadores, lugares de ocurrencia de eventos, eventos pasados, capacitaciones, condiciones ambientales, entre otras. Por lo tanto, la utilización de herramientas de ML ayuda a procesar estos datos y clasificarlos simplificando tiempo, y costos [Sarkar and Maiti \(2020\)](#).

[Sámamo-Ríos et al. \(2019\)](#) presentaron una revisión sobre intervenciones de seguridad y salud ocupacional para proteger a los trabajadores de los peligros presentes en el trabajo, donde se pudo evidenciar que se necesita desarrollar y evaluar intervenciones que aborden de forma específica los riesgos que presentan los jóvenes en los trabajos. Por otro lado, [Duarte et al. \(2019\)](#) revisó los accidentes laborales en la industria minera, donde encontraron que los camiones de acarreo, volquetes y cintas transportadoras son los equipos que generan mayor impacto en las tasas de accidentes laborales en las mineras.

Durante los últimos años los sistemas afectivos han tomado mayor protagonismo en diversas áreas de investigación por lo que diversos investigadores han realizado estudios donde se logra evidenciar que estos sistemas generan importantes contribuciones en los avances tecnológicos. Un estudio llevado a cabo por [Verkijika \(2020\)](#) donde se presentó un modelo que busca comprender el papel de las emociones en la aceptación de los sistemas de pago móvil, en el cual se logró mostrar para este caso que el afecto y el arrepentimiento anticipado tienen influencia positiva y de manera significativa en las intenciones conductuales de adoptar pagos móviles, y por otro lado, la influencia de la ansiedad no es significativa. Por otro lado, [Altuwairqi et al. \(2021\)](#), buscaron medir el compromiso de los estudiantes según sus emociones a momento de realizar tareas y desafíos, con el fin de generar mejoras en los procesos de aprendizaje.

Otra contribución sobre la afectividad emocional se refleja en [Kowalczyk et al. \(2019\)](#) estudio que busca un enfoque basado en el seguimiento del estado de las emociones en el monitoreo de conductores, buscando determinar el impacto de la afectividad del conductor en la seguridad de conducción.

Las emociones juegan un papel fundamental al momento de la toma de decisiones de diversas investigaciones. La formulación de la hipótesis del marcador somático, tomó relevancia al momento de explicar el papel de las emociones en la toma de decisiones [Poppa and Bechara \(2018\)](#) . Un ejemplo de uso de marcador somático es el estudio de [Sandor and Gürvit \(2019\)](#) el cual está orientado en la toma de decisiones en adolescentes a través de *Iowa Gambling Task* (IGT), que es un instrumento que se utiliza para evaluar el comportamiento presuntamente defectuoso en la toma de decisiones en grupos de pacientes con ciertas patologías.

Un ejemplo de aplicación de técnicas de *machine learning* es el estudio de [Mosquera et al. \(2021\)](#) donde se propone un sistema de clasificación para la identificación y la prevención de accidentes laborales en bodegas de almacenamiento de fibra en una empresa de pulpa de papel.

Las investigaciones sobre los incidentes ocupacionales deben profundizar en la identificación de las causas, mediante métodos de análisis avanzado. [Davoudi Kakhki et al. \(2019\)](#) buscó clasificar y predecir las causas de los incidentes ocupacionales en las operaciones agro-manufactureras de elevadores de granos en la región del Medio Oeste de los Estados Unidos. Estos entornos de trabajo de fabricación agrícola son peligrosos y los trabajadores están expuestos a sufrir lesiones graves, debido a las actividades que deben realizar, a través de la aplicación lograron confirmar que las redes neuronales artificiales son útiles para estimar los riesgos de seguridad e identificar los factores de riesgo de cada incidente, de esta forma es posible implementar medidas de seguridad para ayudar a prevenir la ocurrencia y severidad de tales incidentes en el ambiente laboral.

Se han aplicado algoritmos optimizados de aprendizaje automático para predecir los resultados de los accidentes laborales, los populares son, máquina de vectores de soporte(SVM) y red neuronal artificial (ANN), así lo ilustra [Sarkar et al. \(2019c\)](#) mediante un estudio que desarrolló un modelo de predicción utilizando técnicas de aprendizaje automático, SVM y ANN para la predicción de resultados de incidentes ocupacionales, donde utilizaron técnicas de optimización, es decir, GA y PSO en los clasificadores.

Asimismo, está el estudio de [Ayhan and Tokdemir \(2020\)](#), donde desarrollaron un novedoso modelo para predecir los resultados de los incidentes de construcción utilizando el análisis de agrupamiento de clases latentes (LCCA) y las redes neuronales artificiales (ANN) para determinar las acciones preventivas necesarias tomando datos reales recopilados desde varios sitios de construcción, presentando acciones preventivas prácticas para evitar incidentes en este rubro.

Un estudio realizado por [Sarkar et al. \(2019a\)](#) explica que los accidentes por resbalón, tropiezo o caída (STF) son los principales causantes de lesiones, por lo tanto, los investigadores llevaron a cabo un estudio que tenía como objetivo predecir la ocurrencia de los (STF) a través de clasificadores de árbol de decisión.

Basándose en el riesgo ocupacional, está el estudio de [Zhu et al. \(2021\)](#) relacionado a los accidentes de construcción utilizando técnicas de aprendizaje automático, para llevar a cabo un análisis y evaluación del impacto de diversos factores en la predicción de la gravedad de estos accidentes. El estudio utilizó algoritmos desde regresión logística, árboles de decisión, máquina de vectores de soporte, entre otros similares, así se generaron

diversas conclusiones que pueden ser utilizadas para realizar una mejora en la seguridad de la construcción urbana.

Otro estudio relacionado con estas técnicas es [Xu et al. \(2020\)](#) el que tuvo como objetivo mejorar la comprensión de los accidentes de helicópteros y su historial de seguridad a través de técnicas de aprendizaje automático, generando una comparación del rendimiento de distintos clasificadores y utilizar el mejor de estos para el análisis y la prevención de accidentes.

A pesar de la diversidad de estudios relacionados con el riesgo ocupacional y los accidentes laborales a través de técnicas de *machine learning*, no se ha profundizado la utilización de la afectividad emocional o el uso de marcadores somáticos para abarcar esta área basándose en la toma de decisiones.

# Capítulo 2

## Marco teórico

### 2.1. Riesgo Ocupacional

Según [Rabeiy et al. \(2018\)](#) el riesgo ocupacional se define como la posibilidad de ocurrencia de un evento en el ambiente de trabajo, que pueda generar un daño y con consecuencias de distinta severidad.

La minería presenta una variedad de peligros tanto físicos como de exposición para la salud de los trabajadores. Por lo tanto, la identificación y evaluación de los riesgos es cada vez más importante para evitarlos.

### 2.2. Hipótesis del marcador somático

La afectividad es definida como el grado de reacción que tienen las personas ante estímulos internos o externos, por lo tanto, es la capacidad para sentir emociones y sentimientos y las decisiones de las personas están ligadas a los sentimientos [Cosentino et al. \(2016\)](#).

El proceso racional en la toma de decisiones se basa en evaluar los pros y contras de la información y así responder de forma más adaptativa, comparado con la toma de decisiones en un procesamiento afectivo, en los cuales sus decisiones se llevan a cabo a partir de informaciones pasadas asociadas al estímulo presente.

La emoción asalta al cuerpo, altera las constantes como el ritmo cardíaco y la circulación de hormonas en sangre. Ésta es pública por ser detectable y muchas veces visible, precede al sentimiento más privado y racional. Ante un gran peligro, el miedo llega primero en forma de calor, palpitaciones, temblores y después se afirma a la conciencia real del miedo y su causa, el mundo interior se ve afectado y modificado tanto por el entorno como por el propio medio interno como lo son los pensamientos o recuerdos [Márquez et al. \(2012\)](#).

Según [Linguist and Bartol \(2013\)](#) los sentidos tales como la visión, el oído, el olfato, el gusto y el tacto construyen para el cerebro una representación del mundo exterior y las emociones intentan configurar el estado del organismo en la congruencia de percep-

ciones externas y pensamientos. La emoción y los cambios fisiológicos que la acompañan quedan asociados en el cerebro a la situación que los ha provocado y resurgirán cuando se reproduzcan estas circunstancias.

Esta es la explicación que hace el médico y neurólogo, António C. Rosa Damásio sobre el término marcadores somáticos. Antonio Damásio, nació en Lisboa en el año 1944 y fue el autor del libro *El error de Descartes* (*Descartes' error: Emotion, rationality and the human brain*) el cual en su primera edición tenía como subtítulo “La razón de las emociones”.

Según Damásio (1994) la toma de decisiones está guiada por cambios homeostáticos evidentes que el cuerpo genera, éste por sí mismo envía señales que son transformadas en cambios corporales inmediatos, de esta manera anticipando la toma de decisiones, y evidentemente los resultados de dichas preferencias, de manera que disminuye el proceso racional.

Esta hipótesis del marcador somático es una teoría que rompe con el dualismo entre cerebro y cuerpo, se termina con la idea de que el cerebro es una parte diferente al resto del cuerpo, sino una unidad corporal.

Entonces Damásio (1994), describe que el marcador somático consigue tomar atención al resultado negativo al que la acción determinada pueda conducir, funcionando como una alarma que hace referencia a tener cuidado si se llegase a tomar dicha opción, ya que su resultado es peligroso. Por lo tanto, puede llevarse a rechazar de forma inmediata, y hacer que se elija otra alternativa.

Los marcadores somáticos son una demostración de sentimientos generados a partir de emociones secundarias, dichas emociones se conectan mediante un aprendizaje. En ocasiones, estos marcadores, pueden operar de manera encubierta esto quiere decir, sin llegar a la conciencia, pueden utilizar un bucle o ciclo, ya que al momento de tomar decisiones van destacando algunas opciones, tomándolas como peligrosas o favorables, y eliminando apresuradamente de una prueba posterior.

En resumen, explica que las elecciones de acciones cuya consecuencia sea inmediata son negativas, pero generan efectos positivos en resultados futuros, por lo que se entiende que la experiencia es el medio a través del cual los marcadores somáticos se van adquiriendo. Estas experiencias están articuladas por dos aspectos, primero es el aspecto interno, el cual regula las preferencias personales y las respuestas psicosomáticas, lo que significa que son respuestas de carácter innato, están preparadas para que el organismo pueda sobrevivir, evitando el peligro, dolor y buscando situaciones de tranquilidad.

Recientemente Cabrera et al. (2020) llevaron a cabo un estudio que incorpora marcadores somáticos artificiales en la toma de decisiones de agentes autónomos, donde se realiza una explicación sobre las fases en la toma de decisiones de los agentes autónomos.

Presentando una cantidad de 5 pasos que son ilustrados en la siguiente Figura:

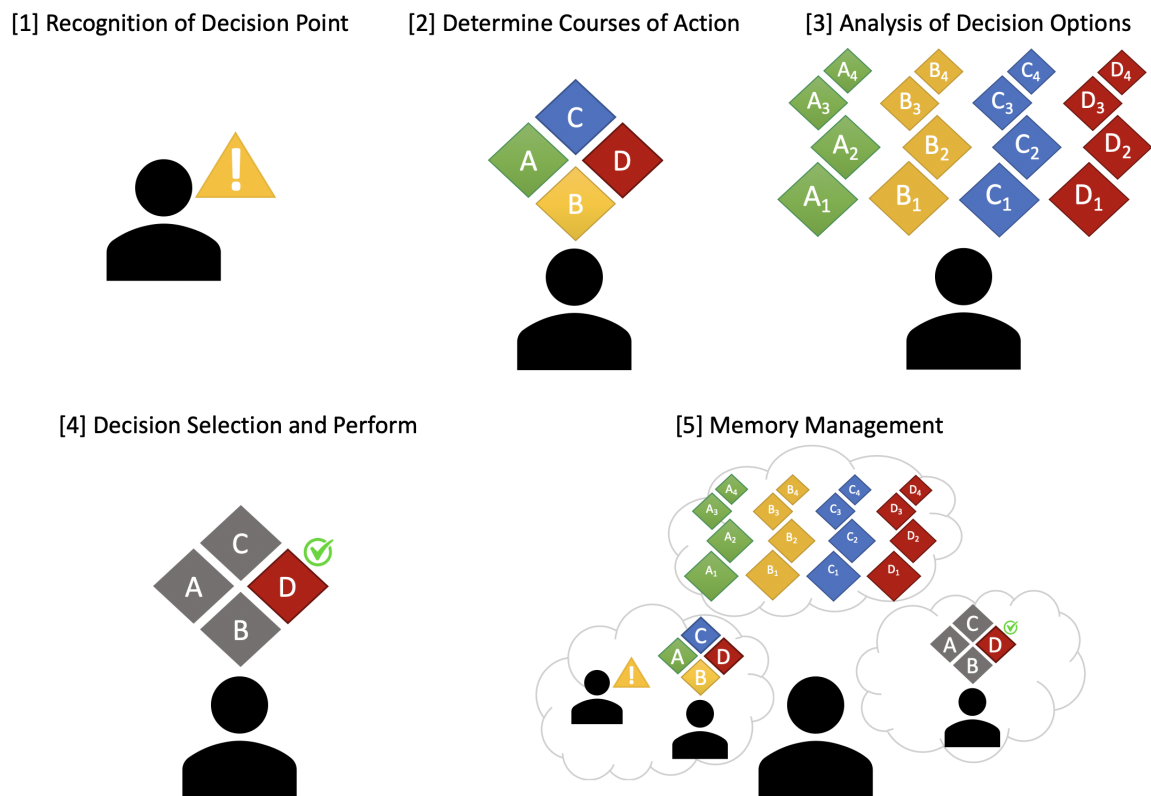


Figura 2.1: Ilustración de fases de la toma de decisiones

Los pasos presentados en la Figura 2.1 son los siguientes:

- Reconocimiento del punto de decisión:** Esta etapa explica que el marcador somático artificial al momento de detectar un estímulo, actúa como mecanismo que capta atención del agente autónomo y busca interpretarlo. Estos estímulos en diversos casos pueden ser asociados a una amenaza u oportunidad, o está relacionado con la presencia de algún sentimiento. Asimismo, puede que no exista relación pasada.
- Determinación de los cursos de acción:** Los marcadores somáticos artificiales actúan con el fin de generar opciones de decisión, presentando alternativas. El agente autónomo tiene la oportunidad de asociar las acciones o situaciones pasadas relacionadas al estímulo presentado en la actualidad, generando una posible solución a un problema específico.
- Análisis de las opciones de decisión:** Un marcador somático artificial participa en la opciones de decisión. La relación existente entre un objetivo a alcanzar y la opción disponible es posible ser representada de forma numérica, fortaleciendo la opción de decisión establecida dentro del conjunto de opciones presentes.
- Selección y ejecución de decisiones:** Es posible que un marcador somático artificial sea partícipe de la selección final de una opción de decisión, mediante uso de recompensas o castigos, esto indica que este proceso de análisis de opciones puede sugerir un camino, y esta fase de selección y ejecución puede ayudar la ejecución de un camino diferente.

- **Gestión de la memoria:** El manejo de la memoria es fundamental debido a que, las actividades somáticas del último proceso de la toma de decisiones pueden registrarse en la memoria de trabajo de agente autónomo, y así ser relacionadas posteriormente. La decisión que se toma y las consecuencias afectivas son guardadas en la memoria a largo plazo, así el agente autónomo tendrá una asociación somática, y ampliará su experiencia actual.

## 2.3. Machine Learning

El *machine learning* es un subcampo de las ciencias de la computación y las matemáticas aplicadas, es considerado una rama de la inteligencia artificial, el cual tiene como objetivo desarrollar técnicas que permitan que las computadoras aprendan, a través de la experiencia para ir mejorando su desempeño [Haykin \(2009\)](#). Se busca encontrar algoritmos que sean capaces de convertir conjuntos de datos en modelos que sean capaces de generalizar comportamientos e inferencias para un amplio conjunto de datos.

Existen **tipos de aprendizajes**, de manera resumida se explican 3: supervisado, no supervisado y por refuerzo.

Según [De Wilde \(1997\)](#):

En el primer caso, el aprendizaje supervisado es considerado una técnica para deducir una función a partir de datos de entrenamiento, estos consisten en pares de objetos normalmente vectores, una componente de este par serían los datos de entrenamiento y por otro lado, los resultados deseados. La salida de la función puede ser un valor numérico o una etiqueta de clase. El objetivo de este tipo de aprendizaje consta en crear una función que sea capaz de predecir el valor que corresponde a cualquier objeto de entrada válida después de haber visto una serie de ejemplos, que son los datos de entrenamiento. Para aquello se tiene que generalizar a partir de los datos presentados a las situaciones no vistas previamente.

En el segundo caso, el aprendizaje no supervisado es cuando no se dispone de *outputs* para el entrenamiento. Sólo se conoce los datos de entrada, pero no existen datos de salida los cuales corresponden a un determinado input. Por lo que, solo se puede describir la estructura de los datos, para intentar encontrar algún tipo de organización que simplifique el análisis. Este tipo de aprendizaje tendría un carácter exploratorio.

Y por último se tiene el aprendizaje por refuerzo el cual tiene como finalidad determinar cuáles acciones se deben escoger en un entorno dado con el objetivo de maximizar alguna noción de “recompensa”.

## 2.4. Bosque Aleatorio

Bosque Aleatorio o en inglés *Random Forest* es una técnica de aprendizaje supervisado que genera múltiples árboles de decisión sobre un conjunto de datos de entrenamiento, los resultados que se pueden obtener se van combinando a fin de obtener un único y robusto modelo, en comparación con resultados de cada árbol por separado. Estos bosques se

forman mediante un algoritmo que introduce una aleatoriedad para reducir la correlación entre los árboles [Breiman \(2001\)](#).

Este clasificador de bosque aleatorio busca construir múltiples árboles de decisión correlacionados para luego promediarlos. El clasificador de bosque aleatorio es un método de aprendizaje de árbol de decisión, basado en conjuntos con el objetivo de mejorar la precisión con una variación reducida y potencialmente evitar sobreajuste [Liu \(2004\)](#).

La idea fundamental de este procedimiento se basa en que para el árbol  $k$ , se genera un vector aleatorio  $\Theta_k$  independiente de los vectores aleatorios pasados  $\Theta_1, \dots, \Theta_{k-1}$  pero con la misma distribución, para lograr hacer crecer un árbol se debe utilizar el conjunto de entrenamiento  $\Theta_k$ , resultando un clasificador  $h(x, \Theta_k)$  donde,  $x$  es un vector de entrada [Breiman \(2001\)](#).

*Un bosque aleatorio es un clasificador que consiste en una colección de árboles estructurados  $h(x, \Theta_k), k = 1, \dots$  donde  $\Theta_k$  son vectores aleatorios independientes e idénticamente distribuidos (i.i.d) y cada árbol emite un voto unitario para la clase más popular en la entrada  $x$  [Breiman \(2001\)](#).*

Para llevar a cabo el algoritmo de bosque aleatorio para una clasificación se utilizan los siguientes pasos según: [Hastie et al. \(2009\)](#) :

- Para  $k = 1$  a  $K$ :
  1. Seleccionar una muestra de inicio  $Z^*$  de tamaño  $N$  de datos de entrenamiento.
  2. Crear un árbol de bosque aleatorio  $T_K$  en relación a los datos de entrada, reiterando los siguientes pasos para cada nodo terminal del árbol, hasta alcanzar el tamaño mínimo de nodo  $n_{min}$ .
    - Seleccionar  $m$  variables al azar de las  $p$  variables.
    - Se realiza la votación eligiendo la clase más popular.
    - Dividir el nodo en dos nodos hijos.
- Salida del conjunto de árboles  $(T_K)_{i=1}^K$

El promedio de  $K$  i.i.d variables aleatorias, cada una con varianza  $\sigma^2$  tiene varianza  $\frac{1}{K}\sigma^2$ , si las variables son simplemente distribuidos de forma idéntica pero no necesariamente independientes con correlación positiva por pares  $p$ , la varianza del promedio sería:

$$p\sigma^2 + \frac{1-p}{K}\sigma^2 \quad (2.1)$$

A medida que aumenta  $k$ , el segundo término desaparece, pero el primero permanece y, por lo tanto, el tamaño de la correlación de pares de árboles en bolsas limita los beneficios de promediar.

La idea de bosques aleatorios es mejorar la reducción de la varianza del ensacado reduciendo la correlación entre los árboles, sin aumentar demasiado la varianza. Esto se puede lograr en el proceso de crecimiento de los árboles mediante la selección aleatoria de las variables de entrada [Hastie et al. \(2009\)](#).

Específicamente, al hacer crecer el árbol en un conjunto de datos de arranque:

Antes de cada división, seleccione  $m \leq p$  de las variables de entrada al azar como candidatas para la división.

Entonces, cuando ya está construido el bosque, se utiliza para realizar la predicción, siendo puesta la media entre las predicciones de cada árbol en el caso de tener un problema de regresión. En el caso de clasificación la predicción será la clase más votada entre todos los árboles del bosque.

Random Forest depende de dos parámetros fundamentales:

1. Ntree: Número de árboles que forman el bosque.
2. Mtree: Número de variables  $p$  que se seleccionan en cada nodo.

La tasa de error de Random Forest tiene relación con los parámetros mencionados, ya que al reducir  $p$  variables, se reduce la correlación entre los árboles, entonces en cada nodo se tiene menos posibilidades para elegir. Sin embargo, al reducir  $p$  se reduce de igual manera la precisión del árbol.

En la práctica el valor de Mtree dependerá del problema, al disminuir la correlación entre los árboles, disminuirá la varianza, por lo que será más preciso el árbol.

Se recomiendan los siguientes valores,  $\sqrt{p}$  para un problema de clasificación y  $\frac{p}{3}$  para un problema de regresión.

Ntree también tiene efecto en la precisión de la predicción, a mayor número de árboles mejor será la predicción, puesto que el número de datos será el promedio mayor [De Wilde \(1997\)](#).

## 2.5. Validación cruzada

La validación cruzada es un método utilizado para evaluar y comparar algoritmos de aprendizaje, es necesario validar la estabilidad del modelo.

Esta estrategia permite estimar la capacidad predictiva de los modelos cuando se aplican a nuevas observaciones, entonces el conjunto de datos se divide en dos segmentos, uno de ellos utilizado para aprender o entrenar un modelo, y el otro para validar el modelo.

### 2.5.1. Validación cruzada k-Fold

El método *k-Fold Cross-Validation* es un proceso iterativo, el cual trata sobre dividir los datos de manera aleatoria en  $k$  grupos de aproximadamente el mismo tamaño,  $k - 1$  grupos se emplean para entrenar el modelo y uno de los grupos se emplea como validación. Este proceso se reitera  $k$  veces utilizando un grupo distinto como validación en cada iteración [Refaeilzadeh et al. \(2009\)](#).

El proceso genera  $k$  estimaciones del error cuyo promedio se emplea como estimación final, para ver la efectividad total del modelo.

Cada punto de datos llega a estar en un conjunto de validación exactamente una vez y llega a estar en un conjunto de entrenamiento  $k - 1$  veces, lo que disminuye significativamente el sesgo, porque se utiliza la mayoría de los datos para el ajuste [Refaeilzadeh et al. \(2009\)](#).

Además se reduce la varianza, porque la mayoría de los datos también se utilizan en el conjunto de validación.

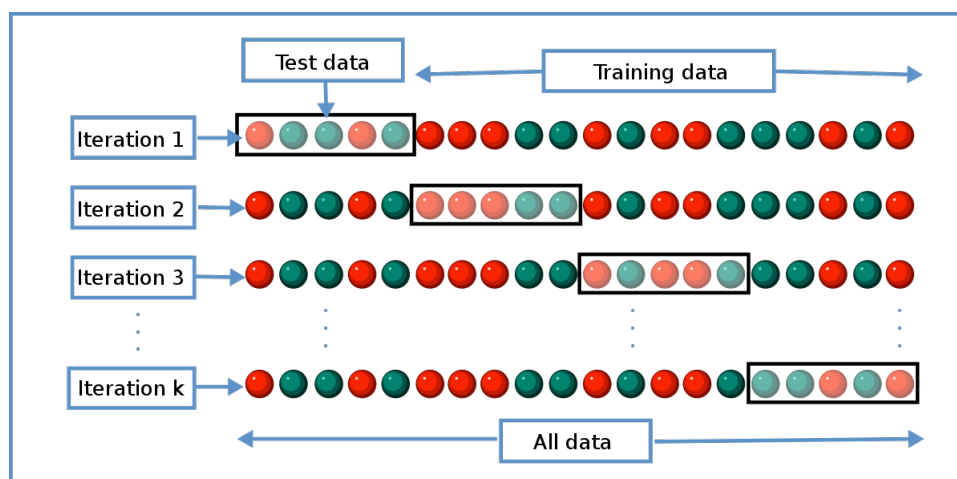


Figura 2.2: Validación cruzada

## 2.6. Validación cruzada estratificada de K-Fold

En los casos donde existe un desequilibrio en la variables de respuesta, se lleva a cabo una ligera variación en la técnica de validación cruzada de K-fold, buscando que cada pliegue contenga aproximadamente el mismo porcentaje de muestra para cada clase que el conjunto de completo, o el valor de respuesta medio es aproximadamente igual en todos los pliegues [Refaeilzadeh et al. \(2009\)](#).

## 2.7. Métricas de desempeño

Las métricas de desempeño pueden brindar resultados satisfactorios para el modelo, debido a que pueden evaluar el desempeño de este [Breiman \(2001\)](#).

### ■ Matriz de confusión

Entrega una matriz con el objetivo de describir el rendimiento completo del modelo.

Cada columna de la matriz representa el número de predicciones de cada clase, mientras que cada fila representa a las instancias en la clase real. Uno de los beneficios de las matrices de confusión es que facilitan ver si el sistema está confundiendo dos clases.

La matriz de confusión esta compuesta por:

- **Verdaderos positivos (VP):** Son aquellos casos que el modelo los predijo como 1 y la salida real también es 1.
- **Verdaderos negativos (VN):** Son aquellos casos que el modelos predijo como 0 y la salida real fue 0.
- **Falsos positivos (FP):** Son aquellos casos que el modelo predijo como 1 y la salida real fue 0.
- **Falsos negativos (FN):** Son aquellos casos que el modelo predijo como 0 y la salida real fue 1.

		Actual Values	
		Positive (1)	Negative (0)
Predicted Values	Positive (1)	TP	FP
	Negative (0)	FN	TN

Figura 2.3: Matriz de confusión

- **Accuracy:** Mide la frecuencia con la que el clasificador hace la predicción correcta. Es la relación entre la cantidad de predicciones correctas y la cantidad total de predicciones.

$$Accuracy = \frac{VP+VN}{VP+VN+FP+FN}$$

- **Precision:** Entrega la proporsión de eventos identificados correctamente como positivos con respecto a todos los verdaderos positivos y falsos positivos.

$$Precision = \frac{VP}{VP+FP}$$

- **Recall:** Entrega la proporsión de eventos que son identificados correctamente como positivos respecto al total de positivos verdaderos.

$$Recall = \frac{VP}{VP+FN}$$

# Capítulo 3

## Metodología

Dentro de este capítulo se presentarán los procesos metodológicos para el desarrollo del proyecto, las respectivas consideraciones y todas las herramientas que se utilizarán. En primer lugar, se presenta la herramienta con la cual se llevará a cabo el análisis, se utiliza el *Software Python 3*, a través de los siguientes paquetes y módulos:

- **Pandas:** Es una librería que proporciona herramientas de análisis y manipulación de datos de alto rendimiento utilizando sus potentes estructuras de datos. Permite leer y escribir con facilidad ficheros en formato CSV, Excel y bases de datos SQL. Además de acceder a los datos mediante índices o nombres para filas y columnas.
- **Numpy:** Es la librería principal para la información científica, proporcionando grandes estructuras de datos, implementación de matrices, y matrices multidimensionales.
- **Math:** Esta librería ofrece funciones matemáticas para uso en el campo de los números reales, tales como funciones numéricas, de potencia y logarítmicas, trigonométricas y de conversión de ángulos y constantes.
- **Scikit-Learn:** Es una de las librerías que cuenta con algoritmos de clasificación, regresión, clustering y reducción de dimensionalidad, utilizada para programar y estructurar los sistemas de análisis de datos y modelado estadístico.
- **Keras:** librería diseñada específicamente para hacer experimentos con redes neuronales, permitiendo crear prototipos de manera fácil y rápida. Proporcionando modelos complejos de aprendizaje profundo.
- **Matplotlib.pyplot:** Es una librería de Python especializada en la creación de gráficos en dos dimensiones.
- **RandomForestClassifier:** clasificador de bosque aleatorio para ajustar varios clasificadores de árboles de decisión en varias submuestras del conjunto de entrenamiento.
- **Classification report:** crea un informe para las métricas de la clasificación, con precision, recall, f1-score.
- **Confusion matrix:** calcula la matriz de confusión para evaluar la precisión de una clasificación.

- **Accuracy score:** métrica que calcula la precisión del subconjunto con las etiquetas predichas.

## 3.1. Conjunto de datos

Se realiza un análisis exploratorio del conjunto de datos, estos datos provienen de una faena minera y están relacionados a registros de incidentes, tanto información previa, como posterior.

Se procedió a realizar la exploración de los datos mencionados, con el fin de conocer sus características. El conjunto de datos en primera instancia contiene un total de 32 variables, donde una de ellas fue tomada como variable respuesta, la cual trata sobre el Nivel de Criticidad del incidente.

En la exploración se pudo apreciar los tipos de variables, la cantidad de valores faltantes, la presencia de variables no estructuradas relacionadas a la descripción del incidente. Además, ver las relaciones entre las variables, para ver el comportamiento de estas.

Como preprocesamiento de los datos, se determinaron las variables a utilizar las cuales solo debían tener relación con lo sucedido antes y en el incidente, por lo tanto no se ocuparon las variables relacionadas a lo posterior a esto, ya que no eran relevantes en este estudio.

### 3.1.1. Variables de análisis

- **Variable Dependiente:**
  - **Nivel de Criticidad antes de la investigación:** Es un indicador proporcional del riesgo que permite establecer una jerarquía o prioridades en los incidentes, es una variable categórica (Alto, Medio y Bajo.)
- **Variables Independientes:**
  - **Fecha Incidente:** Fecha en la cual se registra el incidente, variable temporal día, mes y año.
  - **Compañía:** Contiene el nombre de las compañías donde ocurrió el incidente, variable cualitativa nominal.
  - **Nivel Organizacional Registra Incidente:** Nivel organizacional asociado al incidente.
  - **Clasificación Incidente:** El incidente puede estar clasificado en tres secciones, seguridad y salud, daño material y medio ambiente, variable cualitativa nominal.
  - **Tipo Causa:** El tipo de causa del incidente está dividido en 4 secciones, acción individual, condición de la tarea, defensa ausente y factor organizacional, variable cualitativa nominal.

- **Lugar de Ocurrencia:** Lugar donde se registró el incidente, variable cualitativa nominal.
- **Lesión:** Tipo de lesión provocada por el incidente, variable cualitativa nominal.

Este conjunto de datos contiene un total de 2643 registros, y para este análisis se utilizarán las variables presentadas anteriormente y además se incorporaron dos variables que fueron creadas a partir de la fecha del incidente las cuales son:

- **Último festivo:** Esta variable contiene la cantidad de días que han transcurrido desde el último festivo, por lo que es cuantitativa.
- **Próximo festivo:** Esta variable contiene la cantidad de días que faltan para el próximo festivo, por lo que de igual forma es cuantitativa.

Cabe mencionar, que debido a que el trabajo está bajo confidencialidad no es posible dar más información sobre las variables a utilizar.

Así se procedió a codificar las etiquetas de las columnas de cadenas, para así proceder con el análisis.

## 3.2. Variable respuesta

A través de la exploración de los datos, se fue evidenciando un desbalanceo en los datos, donde la variable respuesta nivel de criticidad tiene una mayor concentración de registros en el nivel medio, el cual fue codificado como 1, mientras que existe una baja cantidad de registros en niveles bajo y alto los cuales fueron codificados como 0 y 2 respectivamente.

## 3.3. Técnicas de clasificación

Se aplicó un modelo de clasificación Random Forest sobre los datos. En el caso de la variable respuesta esta cuenta con 3 niveles bajo, medio y alto.

En el cuadro 3.1 se entregan las cantidades de características por nivel y su respectiva etiqueta en el conjunto de datos, se puede apreciar el desbalanceo que se mencionó anteriormente.

Nivel	Cantidad de características	Etiqueta
Bajo	323	0
Medio	2236	1
Alto	84	2

Cuadro 3.1: Cantidad de características por nivel

Se aplicará una clasificación binaria por lo que la variable se dividirá en tres clases, lo que dará un total de 3 clasificaciones que son explicadas en las siguientes tablas:

Clasificador	Clases	Nombre clase	Etiqueta	Cantidad de características
Random Forest	Bajo vs Alto y Medio	Alerta	Bajo = 1 Alto y Medio = 0	323 2320
Random Forest	Medio vs Alto y Bajo	Alerta 1	Medio = 1 Alto y Bajo = 0	2236 407
Random Forest	Alto vs Medio y Bajo	Alerta 2	Alto = 0 Bajo y Medio = 1	84 2643

Cuadro 3.2: Clasificador Random Forest

### 3.4. Segmentación del conjunto

En el análisis se dividió el conjunto en 3 partes, **Entrenamiento**, **Validación** y **Prueba**, con conjunto de datos de entrenamiento se entrena un modelo, el aprendizaje ocurre en este conjunto de datos, por otro lado, el conjunto de datos de validación se utiliza para ajustar los modelos ya entrenados, es donde a veces se elige el modelo final para posteriormente probarlo utilizando los datos de prueba. Finalmente, el conjunto de prueba es para probar la predicción de nuestro modelo en este subconjunto.

### 3.5. Re-Clasificación

La reclasificación es una técnica que busca aprender de cuando el modelo se equivoca, en ella se obtendrá un grado de equivocación y de acuerdo con este se tratará de reclasificar.

Para llevar a cabo esta reclasificación de los datos se debe realizar un análisis previo en el conjunto de validación, es ahí donde se realizarán las simulaciones para obtener los umbrales necesarios, estos umbrales son los valores que se calcularán de manera heurística y serán partícipe de las condiciones necesarias para generar la reclasificación.

Entonces se entrena y se obtienen las clasificaciones binarias 0,1. A partir de estas clasificaciones se extraen los vectores mal clasificados, los cuales serían los 0 mal clasificados (el modelo predijo 1, pero realmente eran 0), y los 1 mal clasificados (el modelo predijo 0, pero realmente eran 1).

Para ambos casos, tanto 0 y 1 mal clasificados se tendrán sus probabilidades de clasificación  $P(output = 0)$  y  $P(output = 1)$ , estas probabilidades son obtenidas directamente a través del método *predic-proba* el cual viene incorporado en el paquete Sklearn del software Python, este método devuelve las probabilidades de clase para punto de datos. En la predicción se tendrá una matriz que contiene la etiqueta predicha para cada punto de datos, en el *predict-proba* se calculan las probabilidades y devuelve una matriz de listas que contienen aquellas probabilidades de clase para los puntos de datos de entrada.

Por lo tanto, a partir de las probabilidades se tiene:

- En el caso de los 0 mal clasificados la  $P(output = 1) > P(output = 0)$  se obtuvo que la probabilidad de 1 es mayor.
- En el caso de los 1 mal clasificados la  $P(output = 1) < P(output = 0)$  se obtuvo que la probabilidad de 0 es mayor.

A partir de los vectores mal clasificados y sus respectivas probabilidades de clasificación se obtendrán dos umbrales un umbral de probabilidad de clasificación y otro umbral de distancia al centroide los cuales serán fijos. El centroide es el vector promedio de las variables independientes, en este caso se toman los vectores mal clasificados y se calcula el promedio.

Ya con los umbrales de probabilidad de clasificación y los centroides de cada caso, se extraerán los “sospechosos” que serían considerados los vectores con probabilidad de clasificación por debajo del umbral obtenido, estos son los que podrían ser candidatos a una reclasificación.

Con los vectores sospechosos para cada caso se lleva a cabo el cálculo de la distancia entre los vectores al centroide y con las distancias ya calculadas se fijará un umbral de distancias. Cuando el vector esté muy cerca del centroide quiere decir que está más cerca del error.

A partir de lo explicado anteriormente se desprenden las siguientes notaciones con su definición:

- $P(output = 0)$  es la probabilidad de clasificación cuando la predicción es 0.
- $P(output = 1)$  es la probabilidad de clasificación cuando la predicción es 1.
- $U_1$  umbral de probabilidad de clasificación para 1.
- $U_0$  umbral de probabilidad de clasificación para 0.
- $C_1$  centroide de vectores mal clasificados para 1.
- $C_0$  centroide de vectores mal clasificados para 0.
- $\delta_1$  umbral de distancia al centroide para los 1 mal clasificados.
- $\delta_0$  umbral de distancia al centroide para los 0 mal clasificados.

En el caso de las probabilidades de clasificación, se obtienen al momento de realizar dicha clasificación se obtienen 2 probabilidades una para cero y otra para uno, esto quiere decir, que la que tenga mayor probabilidad será el valor entregado por el modelo en la clasificación.

A continuación se presenta un diagrama explicando el paso a paso para obtener los umbrales que se ocuparán en la reclasificación:

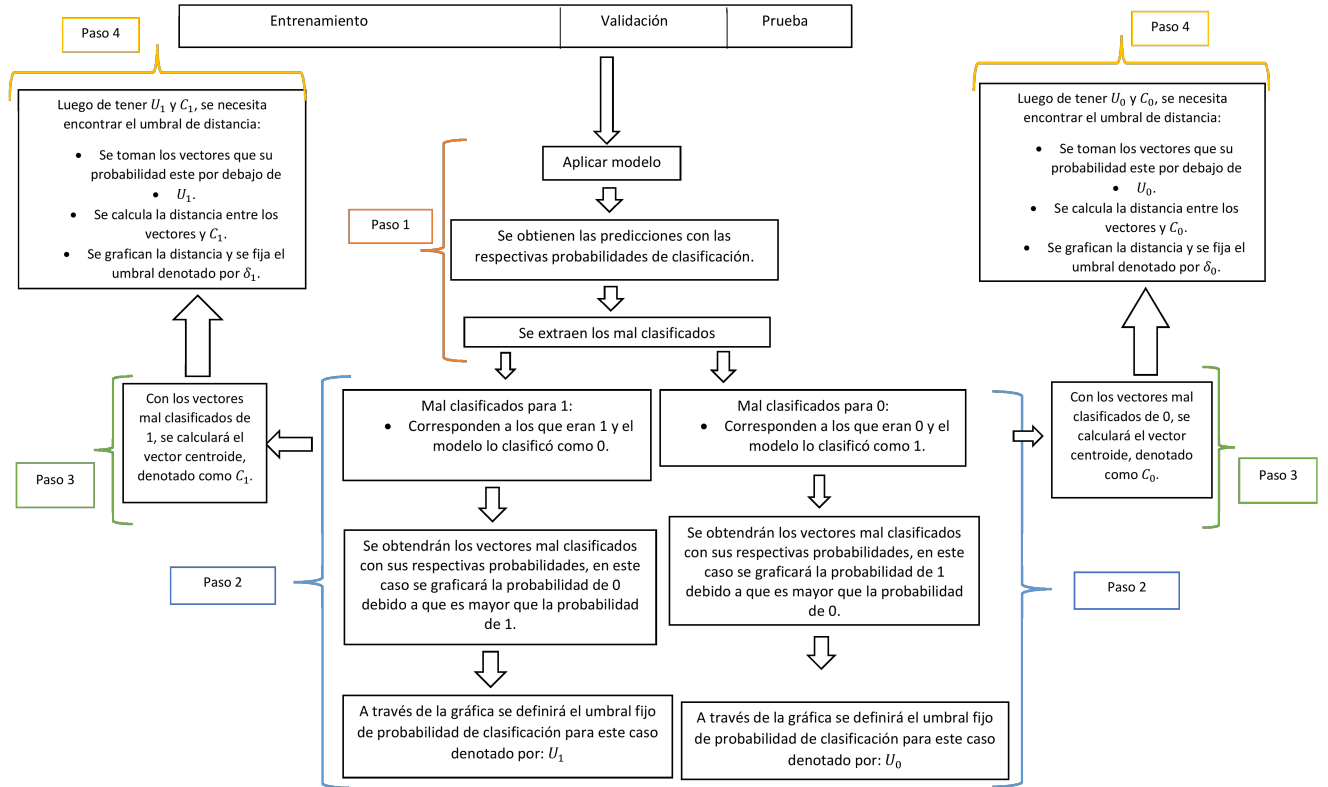


Figura 3.1: Diagrama Umbrales

Luego de obtener los umbrales fijos, se da paso a la reclasificación, se presenta un diagrama de flujo relacionado a los pasos que se deben realizar:

### 3.5.1. Umbrales actualizados

Se hace referencia a umbrales actualizados porque los umbrales que eran fijos, ahora se actualizarán mediante la siguiente ecuación para cada caso:

- Caso de umbrales de probabilidad  $(U_1, U_0)$ :

$$U_{i*} = \left(1 - \frac{(-1)^\beta}{k}\right)(U_i \cdot w) \cdot b \quad (3.1)$$

donde,

- $U_{i*}$  viene siendo la notación para definir la actualización del umbral.
- $\beta = [\hat{y} - y]^2$  sería la diferencia entre el valor de  $y$  clasificado con su respectivo valor real al cuadrado, ya que se aplicó el modelo al conjunto de validación, por lo tanto, se podrá ver los elementos predichos y los reales  $(y, \hat{y})$ .
- $k$  es el número de vectores en el conjunto de validación.
- $U_i$  es el umbral de probabilidad de clasificación fijado en la etapa anterior.
- $w = U(P(\text{output} = i), 1)$  es un número aleatorio entre la probabilidad de clasificación del caso que corresponda y 1.

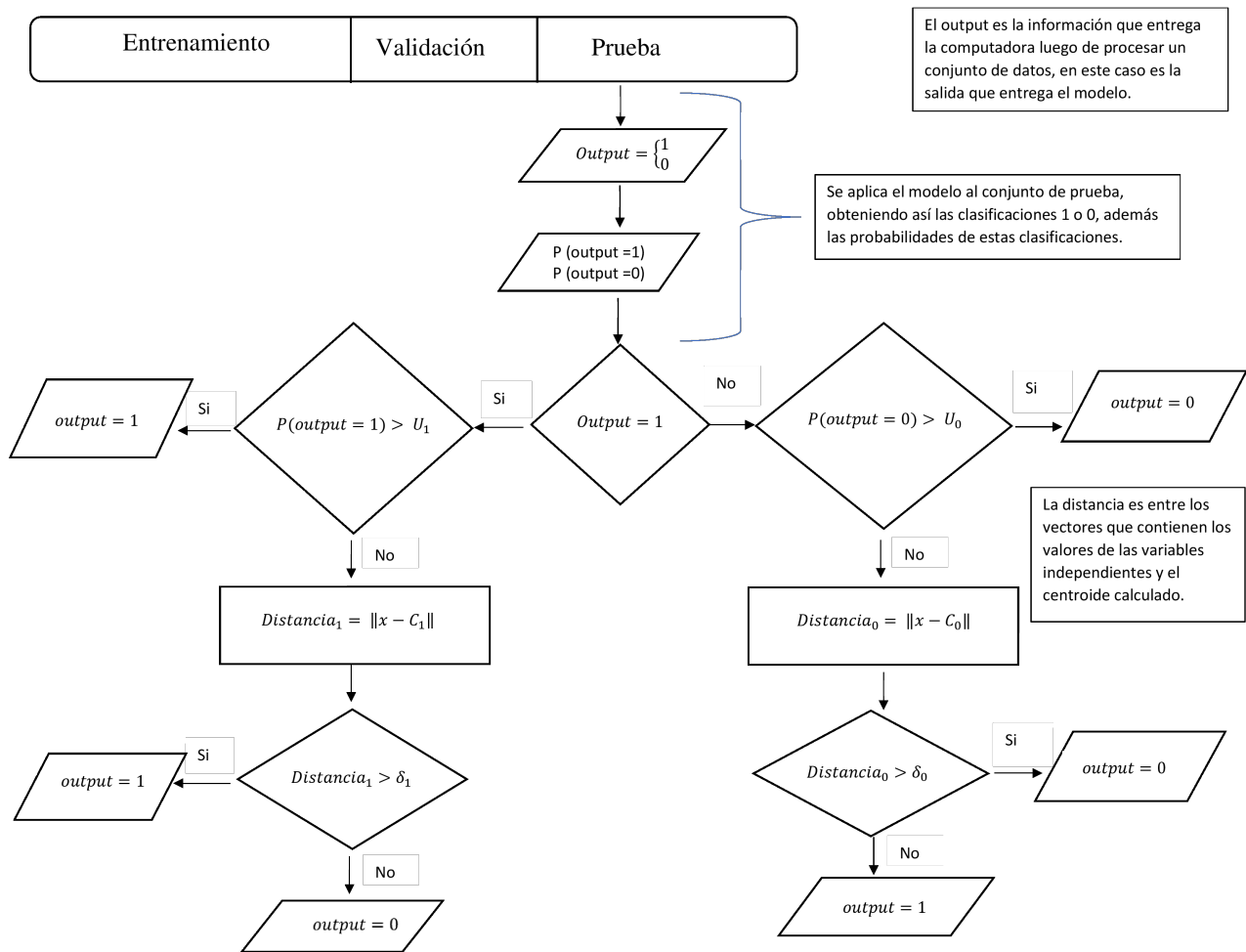


Figura 3.2: Diagrama reclasificación

- $b$  es una constante entre  $U_i$  y 1.

En el caso de los umbrales de distancia a centroide  $(\delta_1, \delta_0)$ :

$$\delta_i^* = \left(1 - \frac{(-1)^\beta}{k}\right) (\delta_i \cdot w) \cdot b \quad (3.2)$$

donde,

- $\delta_i^*$  notación para definir el umbral de distancia actualizado.
- $\beta = [\hat{y} - y]^2$  sería la diferencia entre el valor de  $y$  clasificado con su respectivo valor real al cuadrado, ya que se aplicó el modelo al conjunto de validación, por lo tanto, se podrá ver los elementos predichos y los reales  $(y, \hat{y})$ .
- $\delta_i$  es el umbral de distancia al centroide fijado en la etapa anterior.
- $w = U(P(x = i), 1)$  es un número aleatorio entre la probabilidad de clasificación del caso que corresponda y 1.

- $b$  es una constante entre  $U_i$  y 1.

Entonces se realiza el cambio de umbrales pasando de  $U_i$  como estaba en el caso presentado en (3.2) a  $U_i^*$  que viene siendo el umbral actualizado.

Este proceso es llevado a cabo con el conjunto de validación, los resultados finales de estos umbrales son lo que se utilizarán en el proceso de reclasificación en el conjunto de prueba.

A partir de esto, para cada clasificador se tendrán 3 escenarios:

Clasificador	Escenario
Random Forest	Sin umbral de corte establecido
Random Forest	Umbral fijo establecido
Random Forest	Umbral actualizado

Cuadro 3.3: Escenarios de clasificación

### 3.6. Marcador somático artificial

La incorporación del marcador somático artificial (MSA) se refleja en que se tiene un estímulo de entrada el cual en este caso es la intensidad de la mala clasificación en el modelo, que tan lejos están del umbral de clasificación.

Este estímulo se debe evaluar en una función de activación y en base a esta función hay una intensidad de la reacción la cual tiene una consecuencia en la entidad artificial. El error es un aspecto dentro de la evaluación de la reacción somática artificial.

La función de activación es aquella que dice que por debajo del umbral se aplica la reclasificación o no, está en función del estímulo, y ante el mismo estímulo la entidad artificial no reacciona de la misma manera, la consecuencia no es la misma, es por ello que esta función presenta una parte aleatoria.

La consecuencia se ve reflejada al final del proceso de reclasificación, ya que se incluye la selección y ejecución de decisiones. En la parte de la desigualdad presente al final del diagrama 3.2 la cual si no se cumple se realiza el cambio de decisión.

De manera que cuando se cometa mucho error el mecanismo de reclasificación indique que se debe realizar el cambio. Por lo tanto, la reacción somática artificial puede ser más intensa si el error es más grande.

El algoritmo presente en la figura 3.2 se utiliza para la toma de decisiones, en primer lugar, se identifica el evento obtenido en la clasificación, a partir del resultado se obtendrán dos casos cuando la predicción sea 1 o 0. En cada caso se analiza la intensidad del error cometido, por un lado, a través de la probabilidad de clasificación, y luego con su distancia hacia el vector promedio de los eventos mal clasificados.

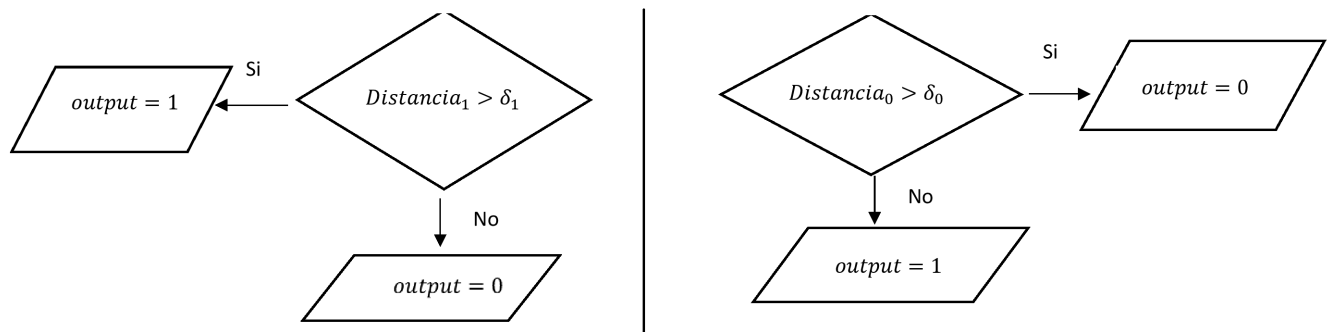


Figura 3.3: Incorporación marcador somático

La figura 3.3 refleja la parte final del proceso de reclasificación para ambos casos, a través de la desigualdad presente entre la distancia calculada y el umbral de distancia hacia el centroide se refleja el error cometido en cada caso, entonces cuando no se cumple dicha desigualdad se debe realizar el cambio en la clasificación obtenida.

El error se puede ver reflejado a través de las probabilidades de clasificación:

Por ejemplo, para el caso de 0 mal clasificados (cuando el valor real era 0 pero se predijo 1) al tener:

$$(P(output = 0)) = 0,1 \text{ y } P(output = 1) = 0,9,$$

se está considerando la probabilidad de  $P(x = 1)$ , entonces hay grado de equivocación muy alto.

Distinto es:

$$P(output = 0) = 0,49 \text{ y } P(output = 1) = 0,51,$$

en donde se presenta un grado de equivocación más leve.

Por lo tanto, cuando se presenta un grado de equivocación muy alto se aplicará el mecanismo de reclasificación para cambiar esa mala clasificación y que pueda aprender de ese error.

# Capítulo 4

## Resultados

En este capítulo se presentan los resultados obtenidos mediante la metodología planteada en el capítulo anterior.

### 4.1. Clasificación

Se utilizan las variables relacionadas directamente con el incidente antes de su investigación, posteriormente el conjunto se dividió en distintos porcentajes:

- **Entrenamiento:** 50 %.
- **Validación:** 25 %.
- **Prueba:** 25 %.

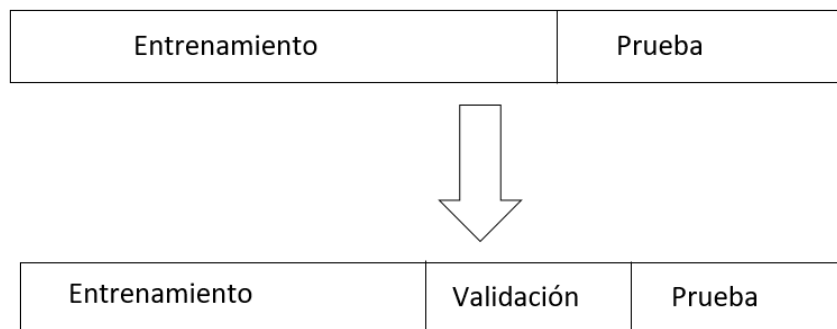


Figura 4.1: División conjunto de datos

Como aparece en 4.1 del total de datos se realizó una división en dos partes en primera instancia, entrenamiento y prueba para luego, a partir de los datos de entrenamiento se divide en validación, entonces de un total de 2643 registros de incidentes la división quedo de la siguiente manera:

- **Entrenamiento:** 1486 registros.
- **Validación:** 496 registros.
- **Prueba:** 661 registros.

#### 4.1.1. Primera clasificación Random Forest Clase “Alerta ”

- La clase alerta consta de las siguientes cantidades de características:

Clase	Niveles	Características
0	Medio y Alto	2320
1	Bajo	323

Cuadro 4.1: Clase Alerta

- Se dividió el conjunto de datos en tres partes quedando con las siguientes cantidades:

División	Características
Entrenamiento	1486
Validación	496
Prueba	661

Cuadro 4.2: División conjunto de datos

En primer lugar, para la búsqueda de los umbrales  $(U_1, U_0, \delta_1, \delta_0)$  tanto de probabilidad como de distancia se ocupó el conjunto de validación, llevado a cabo los pasos descritos en la metodología.

Se entrena el modelo sin establecer ningún umbral y se lleva a cabo la predicción con el conjunto de validación obteniendo la siguiente matriz de confusión:

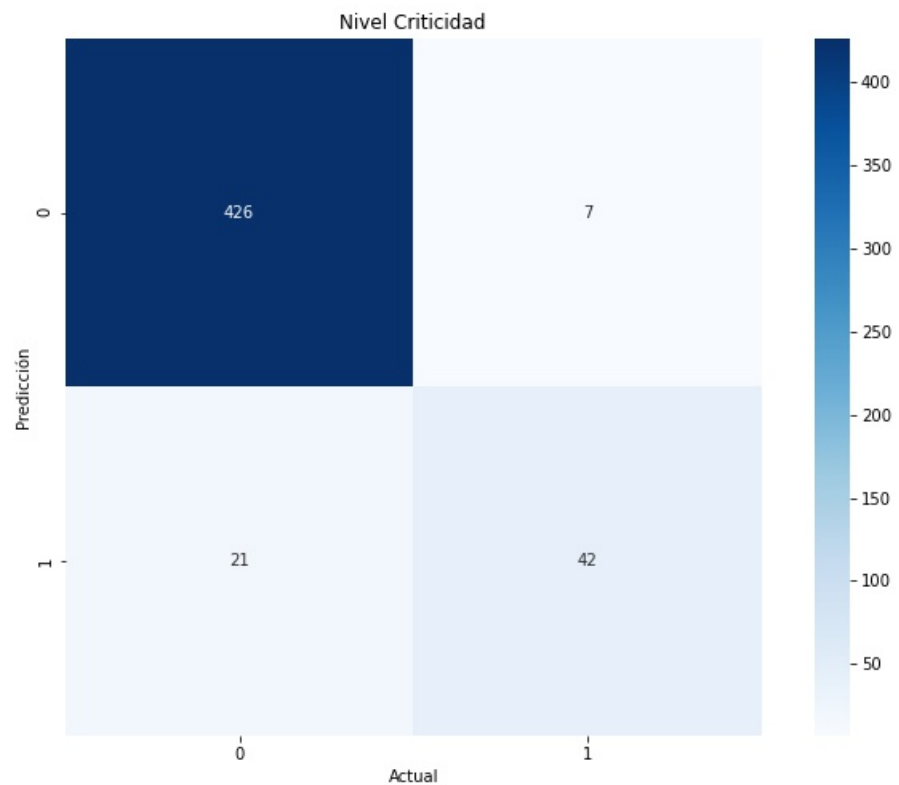


Figura 4.2: Matriz de confusión clase alerta sin incluir umbrales

Obteniendo los siguientes resultados:

- **Verdaderos positivos:** Son **426**, los cuales el modelo los predijo como 0 y la salida real igual es 0.
- **Verdaderos negativos:** Son **42**, los cuales el modelo los predijo como 1 y la salida real igual es 1.
- **Falsos positivos:** Son **7**, los cuales el modelo los predijo como 0 y la salida real igual es 1.
- **Falsos negativos:** Son **21**, los cuales el modelo los predijo como 1 y la salida real igual es 0.

Además, se incluyen las métricas obtenidas:

Clase	Precisión	Recall	f1-score
0	0.95	0.98	0.97
1	0.86	0.67	0.75

Cuadro 4.3: Métricas clase alerta conjunto de validación

Con *Accuracy* de 0,94.

- Como el objetivo es generar una reclasificación, a través del conjunto de validación se realizará la búsqueda de los umbrales, para aplicar esta técnica al conjunto de prueba.

Entonces el primer paso es tomar los vectores mal clasificados que en este caso serían los falsos positivos y los falsos negativos presentes en 4.2.

El total de eventos mal clasificados es de **29**, se puede decir que hay una buena predicción, el objetivo se basa en la reclasificación buscando disminuir el error presente en esta predicción.

A partir de los mal clasificados se tienen:

- Para 1 mal clasificados (eran 1 pero el modelo entregó 0) : **21**.
- Para 0 mal clasificados (eran 0 pero el modelo entregó 1) : **7**.

Se extraen por separado los vectores con sus respectivas probabilidades de clasificación, para obtener el umbral de probabilidad de clasificación.

En el caso de los 1 mal clasificados la probabilidad de cero será mayor, por lo que se genera un histograma con estas probabilidades, y se extrae el umbral fijo:

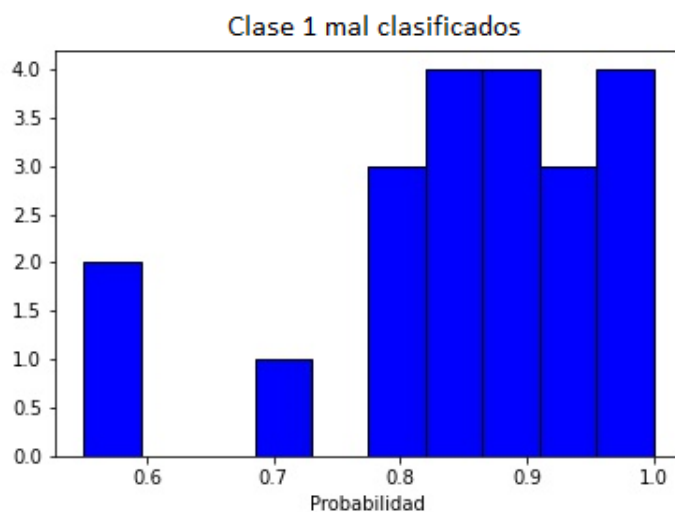


Figura 4.3: Histograma clase 1 mal clasificados clase alerta

Las probabilidades de los 1 mal clasificados van entre 0,59 y 1, donde hay una gran concentración cerca de 1. El promedio está en 0,86.

Luego de una revisión al histograma y el análisis de las probabilidades se determina heurísticamente el umbral  $U_1 = 0,86$  utilizando el criterio que por debajo de esa probabilidad se tomarán los vectores que se consideran sospechosos y pueden ser reclasificados.

Por otro lado, están los 0 mal clasificados que son aquellos que eran 0 pero el modelo entregó como 1. En este caso, la probabilidad de 1 será mayor por lo que se realiza un histograma con dicha probabilidad para obtener el umbral fijo  $U_0$ .

Se presenta el histograma:

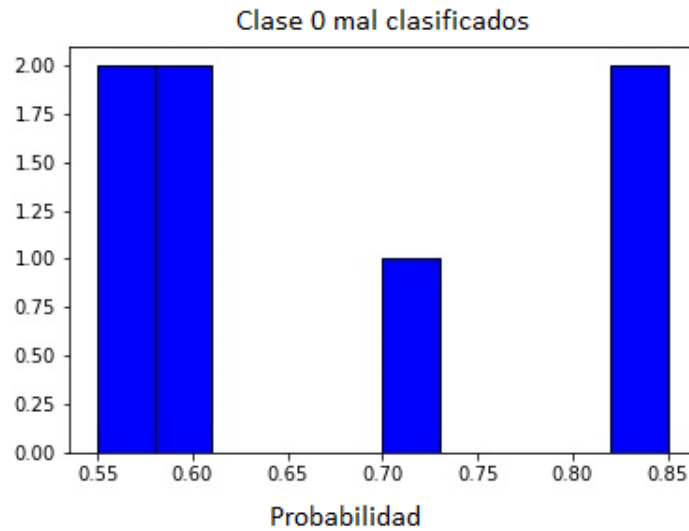


Figura 4.4: Histograma clase 0 mal clasificados clase alerta

La probabilidad de 0 mal clasificados está entre 0,55 y 0,85 con un promedio de 0,67.

Utilizando el mismo criterio anterior, el umbral quedaría  $U_0 = 0,67$ .

- Ya se han cumplido dos pasos del diagrama 3.1 los cuales son la obtención de vectores mal clasificados y sus umbrales fijos. A continuación toca realizar el paso 3 que es el cálculo de los centroides ( $C_1, C_0$ ) para cada caso.

Al obtener este vector centroide para ambos casos, se puede comenzar el paso 4 que es la búsqueda de los umbrales de distancia al centroide. Primero se deben determinar los vectores que se considerarán sospechosos, como se mencionó anteriormente todos aquellos que estén por debajo del umbral de probabilidad de clasificación. A todos estos vectores se le calcula la distancia euclidiana con el centroide, para así realizar la gráfica de distancia para ambos casos y extraer un umbral de distancia fijo ( $\delta_1, \delta_0$ ).

Entonces se presentan los histogramas para ambos casos:

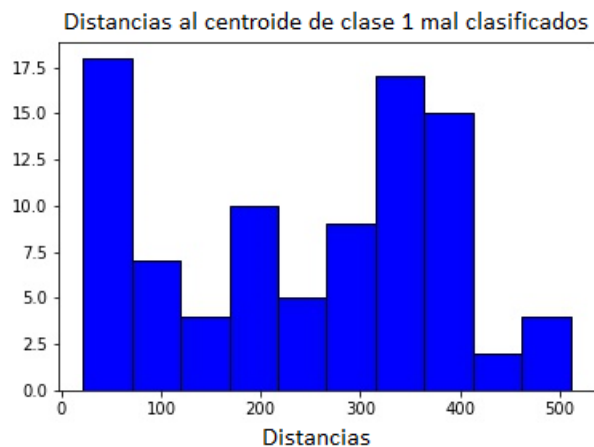


Figura 4.5: Histograma distancias clase 1 mal clasificados Alerta

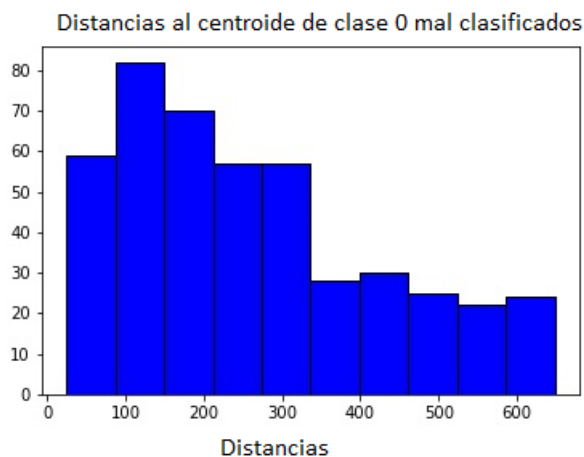


Figura 4.6: Histograma distancias clase 0 mal clasificados Alerta

En el caso, de los 0 mal clasificados las distancias van entre 28, 244 y 506, 922 con un promedio de 242, 21, así el umbral fijo de distancia quedaría en  $\delta_1 = 242$ .

En el caso, de los 1 mal clasificados las distancias van entre 24, 53 y 648, 159 con un promedio de 264, 455, así el umbral fijo de distancia quedaría en  $\delta_1 = 264$ .

De esta manera se tienen finalmente los 4 umbrales necesarios para aplicar la reclasificación de **umbrales fijos**.

Ya con la matriz de confusión sin umbrales 4.2, ahora se obtiene la matriz de confusión aplicando el diagrama de flujo presentado en 3.2 en el software *Python*, con el objetivo de verificar la efectividad de la reclasificación.

El diagrama de reclasificación con los umbrales fijos quedaría de la siguiente forma:

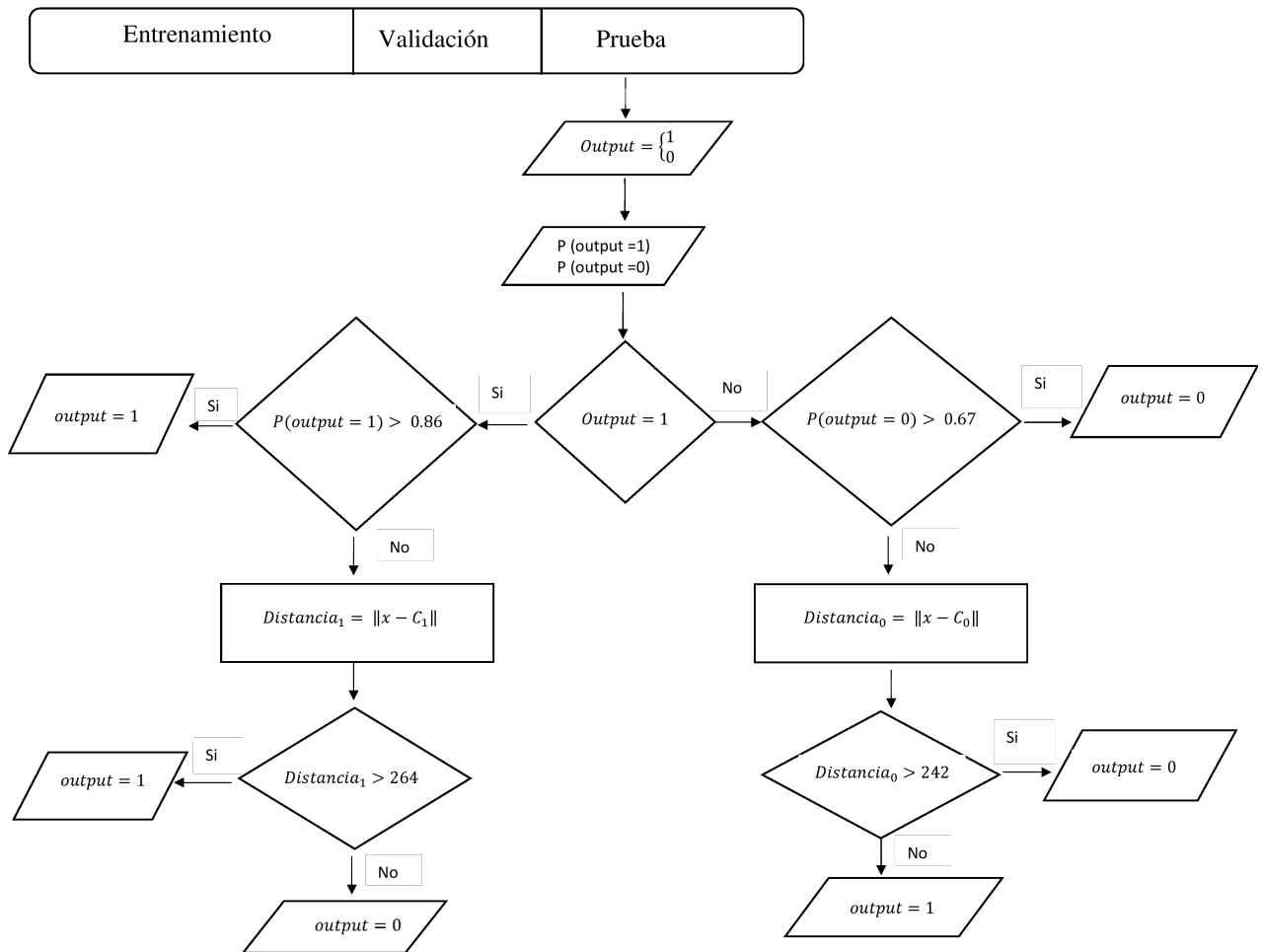


Figura 4.7: Diagrama reclasificación con umbral fijo Alerta

Se obtiene la siguiente matriz de confusión:

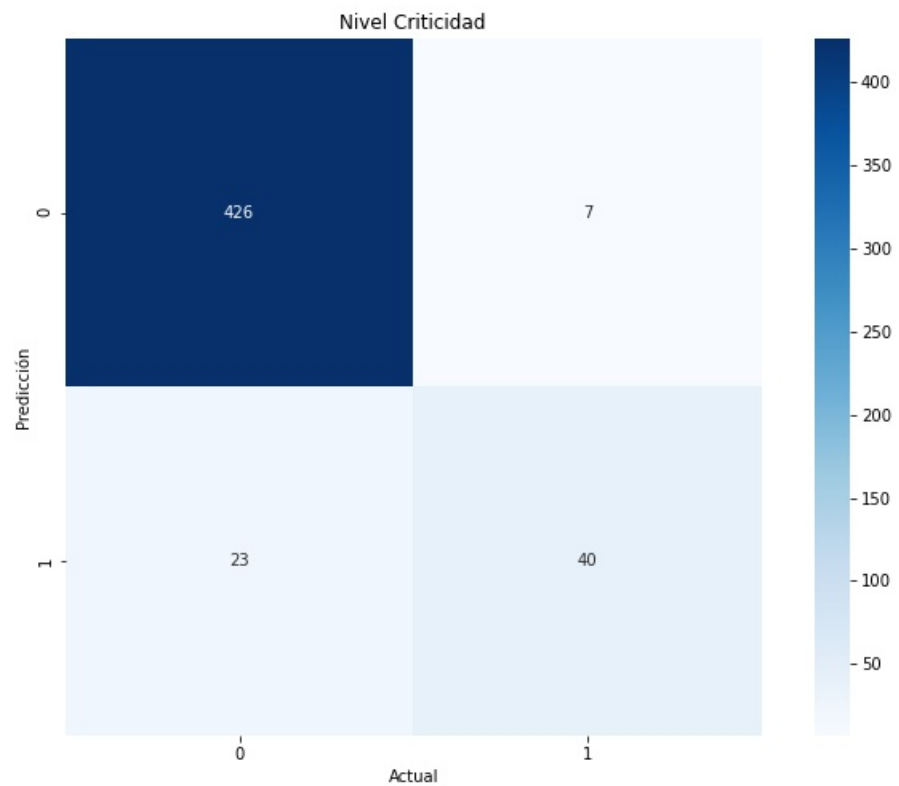


Figura 4.8: Matriz de confusión umbrales fijos Alerta

Obteniendo los siguientes resultados:

- **Verdaderos positivos:** Son **426**, los cuales el modelo los predijo como 0 y la salida real igual es 0.
- **Verdaderos negativos:** Son **40**, los cuales el modelo los predijo como 1 y la salida real igual es 1.
- **Falsos positivos:** Son **7**, los cuales el modelo los predijo como 0 y la salida real igual es 1.
- **Falsos negativos:** Son **23**, los cuales el modelo los predijo como 1 y la salida real igual es 0.

Clase	Precisión	Recall	f1-score
0	0.95	0.98	0.97
1	0.85	0.63	0.73

Cuadro 4.4: Métricas Alerta conjunto validación con reclasificación umbrales fijos Alerta

En el cuadro 4.4 se encuentran las métricas de esta aplicación, con *Accuracy* de 0,94.

Se ve que los resultados no son favorables al aplicar la técnica con los umbrales fijos, debido a que en vez de disminuir los eventos mal clasificados, aumentan.

Se llega al último paso antes de aplicar la reclasificación que es encontrar el umbral actualizado, ya con los umbrales fijos y centroides de cada caso, se aplica la función 3.2 y 3.1 en los cuales se obtendrá  $(U_1^*, U_0^*, \delta_1^*, \delta_0^*)$  que serían los umbrales finales, en los que se aplicaría la reclasificación en el conjunto de prueba.

Para ver la validación del modelo se aplica la reclasificación al conjunto de validación en primera instancia, se obtiene la siguiente matriz de confusión:

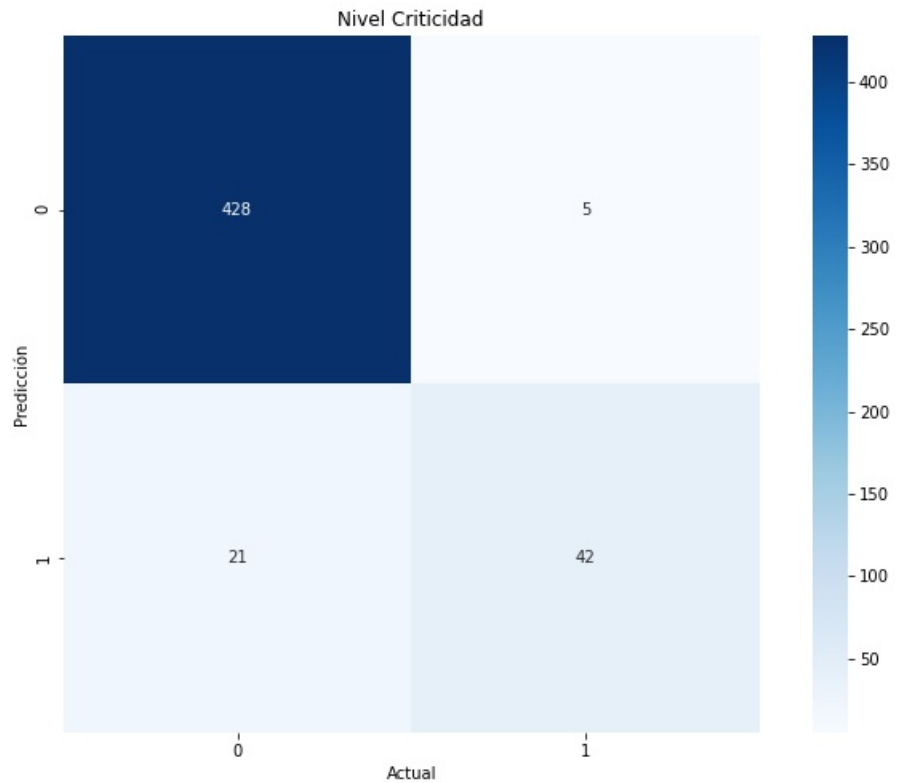


Figura 4.9: Matriz de confusión umbrales actualizados Alerta

Obteniendo los siguientes resultados:

- **Verdaderos positivos:** Son **428**, los cuales el modelo los predijo como 0 y la salida real igual es 0.
- **Verdaderos negativos:** Son **42**, los cuales el modelo los predijo como 1 y la salida real igual es 1.
- **Falsos positivos:** Son **5**, los cuales el modelo los predijo como 0 y la salida real igual es 1.
- **Falsos negativos:** Son **21**, los cuales el modelo los predijo como 1 y la salida real igual es 0.

Clase	Precisión	Recall	f1-score
0	0.95	0.99	0.97
1	0.89	0.67	0.76

Cuadro 4.5: Métricas obtenidas aplicando reclasificación

El cuadro 4.5 presenta las métricas de la clasificación con reclasificación al conjunto de validación. Además, con *Accuracy* de 0,95.

Se presenta una mejora en los verdaderos positivos, debido a que desde un principio el modelo aplicado tenía una muy buena predicción, a pesar de que la baja en los mal clasificados fue muy pequeña si hubo una disminución que era lo que se buscaba, en estos casos ocurre debido a lo que se explicaba anteriormente, por la intensidad de los errores cometidos.

### 4.1.2. Clasificación Random Forest Clase “Alerta” al conjunto de prueba

Luego de calcular los umbrales en el conjunto de validación se procede a aplicar los pasos de reclasificación al conjunto de prueba, todo el proceso realizado con el conjunto de validación tenía como objetivo analizar el error presentado, por lo tanto con los resultados obtenidos en ese proceso se aplica la reclasificación al conjunto de prueba.

Entonces, antes de aplicar la reclasificación se aplica el modelo y se obtiene la siguiente matriz de confusión:

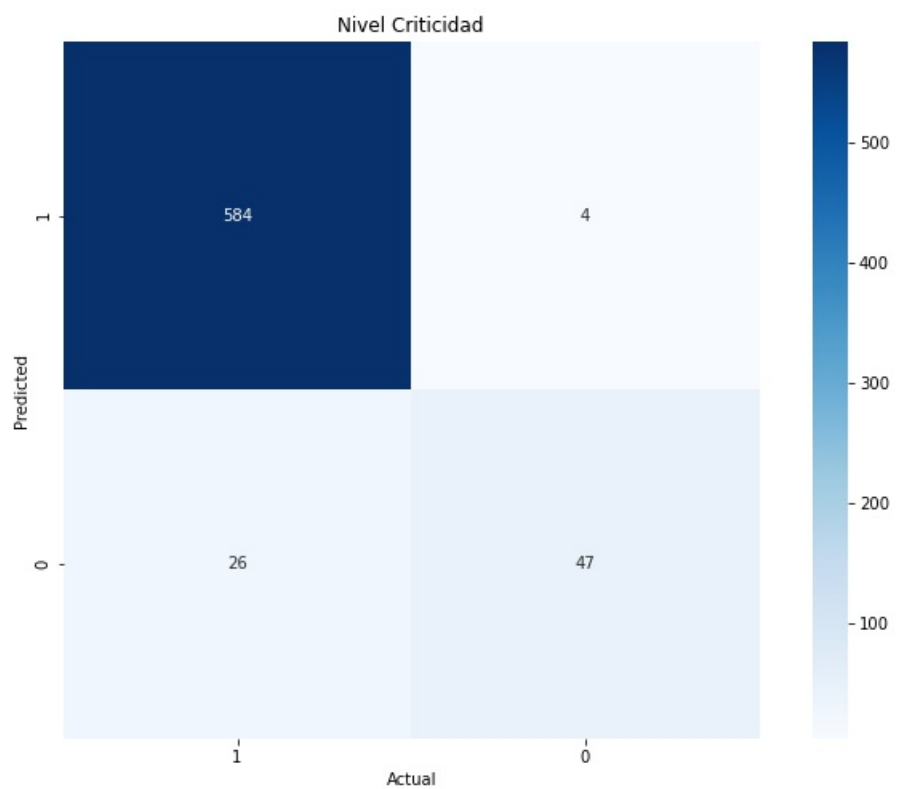


Figura 4.10: Matriz de confusión clase Alerta conjunto de prueba

Obteniendo los siguientes resultados:

- **Verdaderos positivos:** Son **584**, los cuales el modelo los predijo como 1 y la salida real igual es 1.
- **Verdaderos negativos:** Son **47**, los cuales el modelo los predijo como 0 y la salida real igual es 0.
- **Falsos positivos:** Son **4**, los cuales el modelo los predijo como 1 y la salida real igual es 0.
- **Falsos negativos:** Son **26**, los cuales el modelo los predijo como 0 y la salida real igual es 1.

Clase	Precisión	Recall	f1-score
0	0.96	0.99	0.98
1	0.92	0.66	0.77

Cuadro 4.6: Métricas clase Alerta conjunto de prueba sin reclasificación

El cuadro 4.6 presenta las métricas de la clasificación sin reclasificación al conjunto de prueba. Además, con *Accuracy* de 0,96.

Luego se realiza la reclasificación y se obtiene la siguiente matriz de confusión:

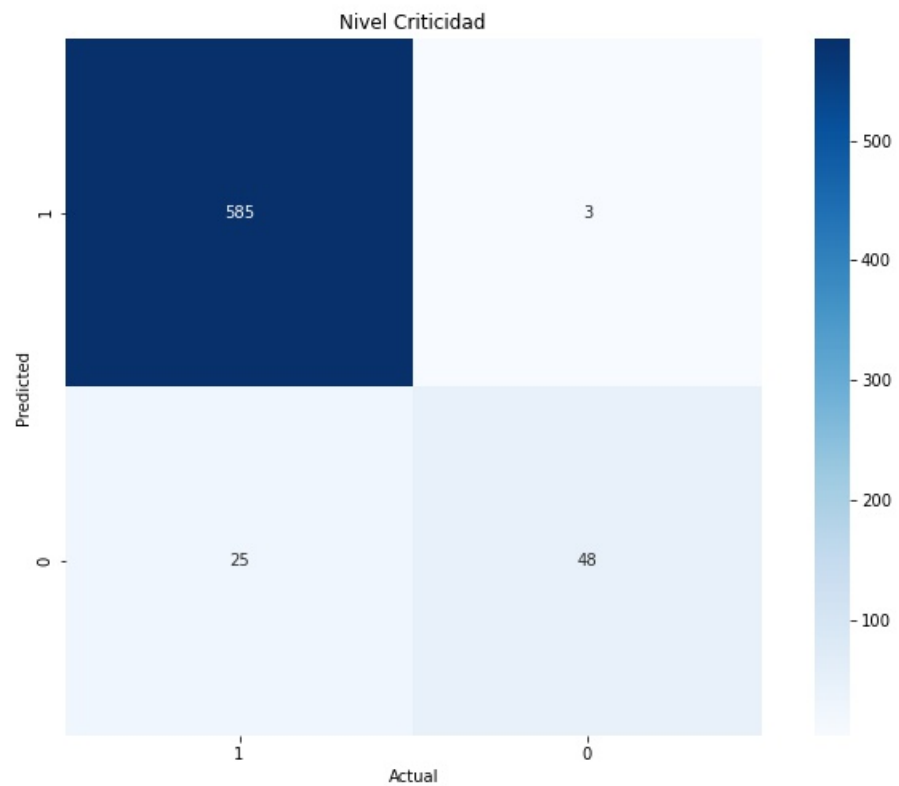


Figura 4.11: Matriz de confusión de clase Alerta con umbrales actualizados aplicados a conjunto de prueba

Obteniendo los siguientes resultados:

- **Verdaderos positivos:** Son **585**, los cuales el modelo los predijo como 1 y la salida real igual es 1.
- **Verdaderos negativos:** Son **48**, los cuales el modelo los predijo como 0 y la salida real igual es 0.
- **Falsos positivos:** Son **3**, los cuales el modelo los predijo como 1 y la salida real igual es 0.
- **Falsos negativos:** Son **25**, los cuales el modelo los predijo como 0 y la salida real igual es 1.

Clase	Precisión	Recall	f1-score
0	0.96	0.99	0.98
1	0.94	0.66	0.77

Cuadro 4.7: Métricas clase Alerta de conjunto de prueba con reclasificación

El cuadro 4.7 presenta las métricas de la clasificación con reclasificación al conjunto de prueba. Además, con *Accuracy* de 0,96.

De la misma manera que en el conjunto de validación existe una mejora en la clasificación bajando 1 en los verdaderos negativos y 1 en los falsos positivos, por ende, el proceso de reclasificación se ha realizado correctamente.

### 4.1.3. Segunda clasificación Random Forest Clase “Alerta 1”

- La clase alerta 1 consta de las siguientes cantidades de características:

Clase	Niveles	Características
0	Medio	2236
1	Bajo y Alto	407

Cuadro 4.8: Clase Alerta 1

De igual manera como se realizó en la clase Alerta, se procedió a la búsqueda de los umbrales necesarios para llevar a cabo la reclasificación.

Se aplica el modelo sin ningún umbral establecido al conjunto de validación para realizar el análisis de los umbrales:

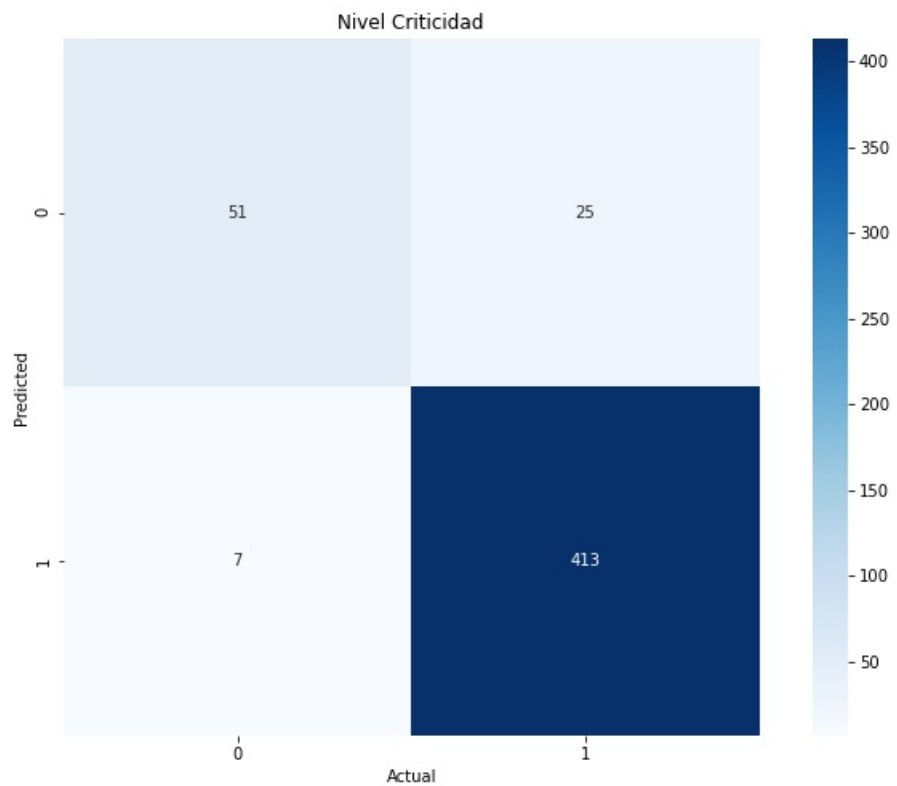


Figura 4.12: Matriz de confusión sin umbrales Alerta 1

Obteniendo los siguientes resultados:

- **Verdaderos positivos:** Son **413**, los cuales el modelo los predijo como 1 y la salida real igual es 1.
- **Verdaderos negativos:** Son **51**, los cuales el modelo los predijo como 0 y la salida real igual es 0.
- **Falsos positivos:** Son **7**, los cuales el modelo los predijo como 1 y la salida real igual es 0.
- **Falsos negativos:** Son **25**, los cuales el modelo los predijo como 0 y la salida real igual es 1.

Clase	Precisión	Recall	f1-score
0	0.87	0.68	0.76
1	0.94	0.98	0.96

Cuadro 4.9: Métricas clase Alerta 1 conjunto validación sin reclasificación

El cuadro 4.9 presenta las métricas de la clasificación sin reclasificación al conjunto de validación. Además, con *Accuracy* de 0,94.

Siguiendo con el objetivo planteado de analizar los vectores mal clasificados se toman los falsos positivos y los falsos negativos para comenzar con el análisis, se tiene un total de 32 datos mal clasificados:

Entonces a partir de las probabilidades se extraen los umbrales de probabilidad de clasificación:

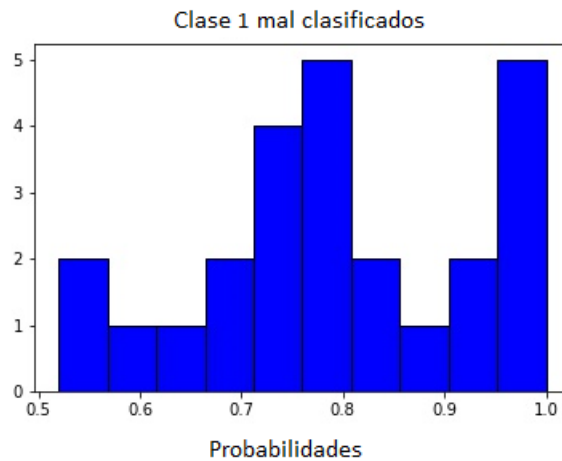


Figura 4.13: Histograma clase 1 mal clasificados clase Alerta 1

En la clase alerta presentada en primer lugar las probabilidades de clasificación para los 1 mal clasificados estaban concentradas muy cerca de 1 en cambio en este caso no se podría identificar una tendencia. Las probabilidades van desde 0,57 hasta 1, con un promedio de 0,82. La gráfica presentada son las probabilidades de 0.

Se define como umbral fijo  $U_1 = 0,82$ . De la misma manera que en la clase Alerta se utiliza el criterio que los vectores por debajo de este umbral sean considerados como sospechosos para una posible reclasificación.

Por otro lado, están los 0 mal clasificados que serían aquellos que el modelo clasificó como 1 pero su salida real era 0, en este caso la probabilidad de 1 sería mayor, se presenta el histogramas de las probabilidades de 1 para este caso:

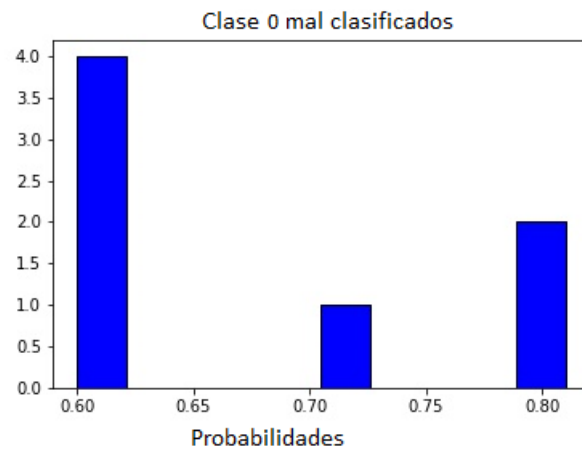


Figura 4.14: Histograma clase 0 mal clasificados clase Alerta 1

Las probabilidades de clasificación para los 0 mal clasificados van de 0,60 hasta 0,88 con un promedio de 0,82.

- El umbral fijo será  $U_0 = 0,82$

Con los umbrales fijos extraídos a partir de los vectores mal clasificados, se procede al cálculo de centroides para cada caso ( $C_1, C_0$ ), los cuales como se mencionó en el caso anterior son el vector promedio de los vectores mal clasificados.

Ya con los vectores calculados se toman los vectores que se encuentran por debajo del umbral fijo en cada caso, los cuales son llamados como sospechosos, y a partir de esos vectores se calcula la distancia euclidiana con el centroide establecido para así graficarlas y utilizarlos como herramientas para ilustrar y definir un umbral fijo de distancia al centroide ( $\delta_1, \delta_0$ ).

- Histogramas de distancias:

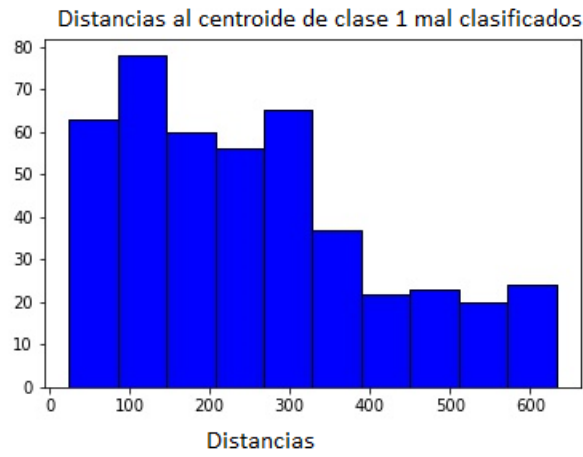


Figura 4.15: Histograma distancias Alerta 1

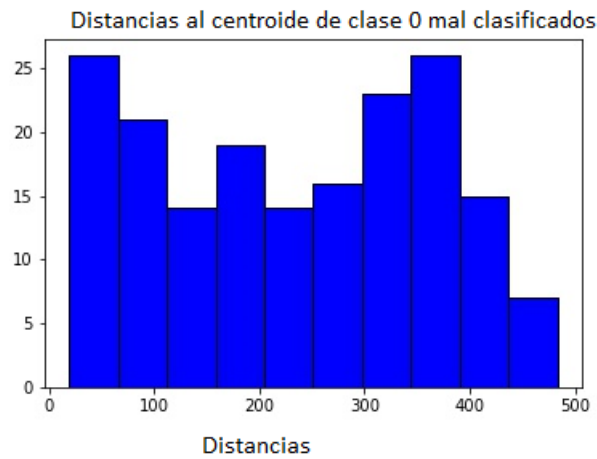


Figura 4.16: Histograma distancias Alerta 1

Entonces, se toma como umbral fijo  $\delta_1 = 230$  y  $\delta_0 = 450$ . De esta manera se tienen finalmente los 4 umbrales necesarios para aplicar la reclasificación de umbrales fijos.

Se aplica la técnica de reclasificación con los umbrales fijos como se hizo con el diagrama mostrado en la ilustración 3.2 a partir de esto se obtiene la respectiva matriz de confusión para umbrales fijos:

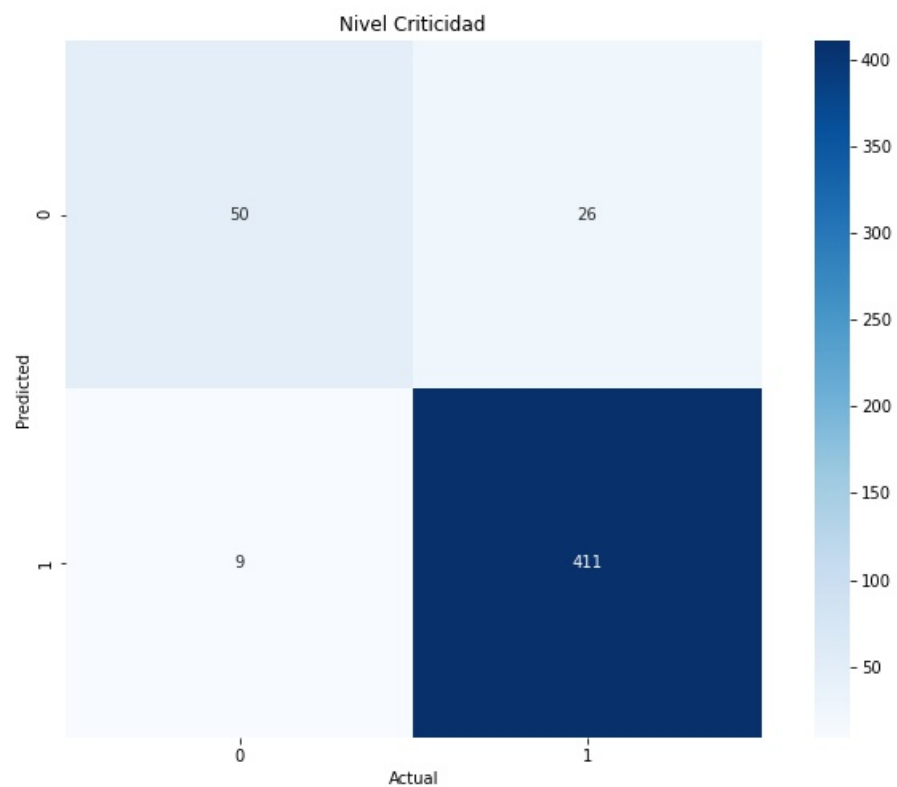


Figura 4.17: Matriz de confusión conjunto de validación con umbrales fijos Alerta 1

- **Verdaderos positivos:** Son **411**, los cuales el modelo los predijo como 1 y la salida real es 1.
- **Verdaderos negativos:** Son **50**, los cuales el modelo los predijo como 0 y la salida real es 0.
- **Falsos positivos:** Son **9**, los cuales el modelo los predijo como 1 y la salida real es 0.
- **Falsos negativos:** Son **26**, los cuales el modelo los predijo como 0 y la salida real es 1.

De igual manera como ocurrió en la clase Alerta cuando se utiliza umbrales fijos no hay una mejora en los eventos mal clasificados.

Clase	Precisión	Recall	f1-score
0	0.85	0.66	0.74
1	0.94	0.98	0.96

Cuadro 4.10: Métricas Alerta 1 conjunto validación con reclasificación

El cuadro 4.10 presenta las métricas de la clasificación con reclasificación al conjunto de validación. Además, con *Accuracy* de 0,93.

Finalmente, se llega al paso final que es encontrar los umbrales actualizados a partir de la función (3.1) y (3.2) para así aplicar al conjunto de prueba.

Se aplica la reclasificación al conjunto de validación de la clase Alerta 1, obteniendo la siguiente matriz de confusión:

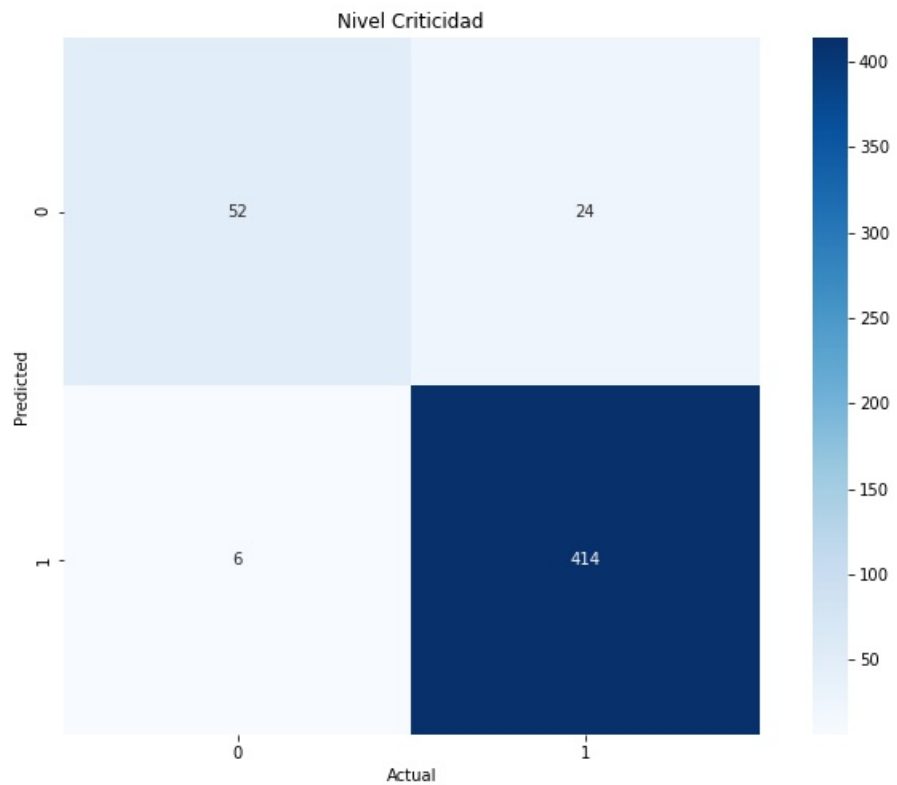


Figura 4.18: Matriz de confusión conjunto de validación con umbrales actualizados Alerta 1

Obteniendo los siguientes resultados:

- **Verdaderos positivos:** Son **414**, los cuales el modelo los predijo como 1 y la salida real es 1.
- **Verdaderos negativos:** Son **52**, los cuales el modelo los predijo como 0 y la salida real es 0.
- **Falsos positivos:** Son **6**, los cuales el modelo los predijo como 1 y la salida real es 0.
- **Falsos negativos:** Son **24**, los cuales el modelo los predijo como 0 y la salida real es 1.

Entonces, hay una mejora en los mal clasificados para este conjunto disminuyendo en 2 datos.

Clase	Precisión	Recall	f1-score
0	0.87	0.67	0.75
1	0.96	0.99	0.96

Cuadro 4.11: Métricas Alerta 1 conjunto validación con reclasificación

El cuadro [4.11](#) presenta las métricas de la clasificación con reclasificación al conjunto de validación. Además, con *Accuracy* de 0,95.

#### 4.1.4. Clasificación Random Forest Clase “Alerta 1” al conjunto de prueba

Se procede a realizar la reclasificación en el conjunto de prueba a partir de lo obtenido en el conjunto de validación.

- Se aplica el modelo al conjunto y se obtiene la siguiente matriz:

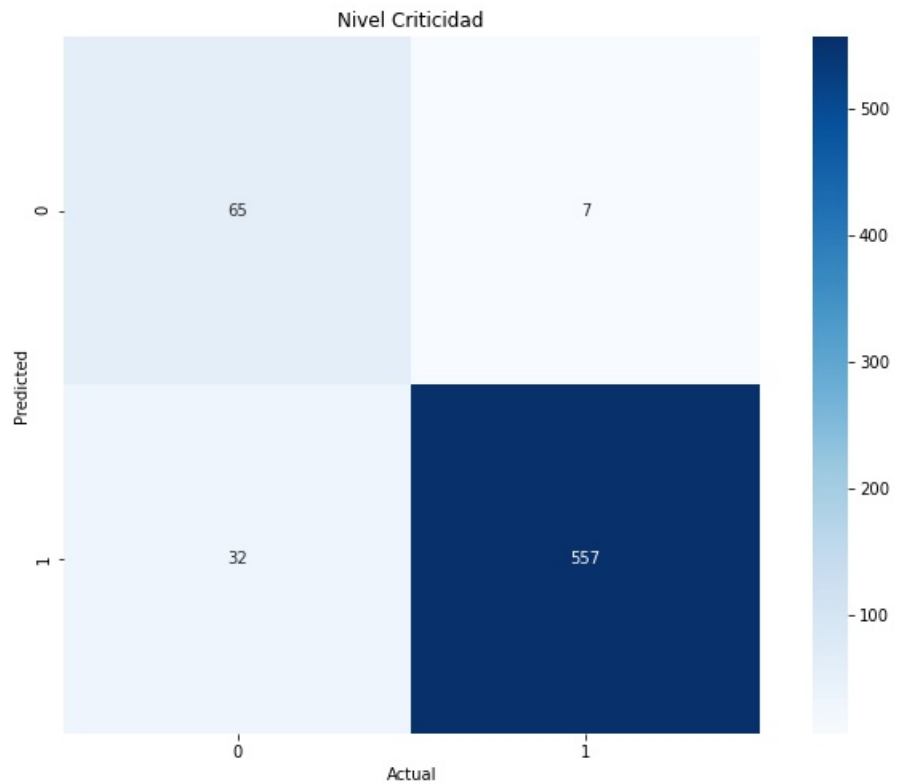


Figura 4.19: Matriz de confusión conjunto prueba para clase Alerta 1

Obteniendo los siguientes resultados:

- **Verdaderos positivos:** Son **557**, los cuales el modelo los predijo como 1 y la salida real es 1.
- **Verdaderos negativos:** Son **65**, los cuales el modelo los predijo como 0 y la salida real es 0.
- **Falsos positivos:** Son **7**, los cuales el modelo los predijo como 1 y la salida real es 0.
- **Falsos negativos:** Son **32**, los cuales el modelo los predijo como 0 y la salida real es 1.

Finalmente, se aplica la reclasificación con los umbrales actualizados y se obtiene la siguiente matriz de confusión:

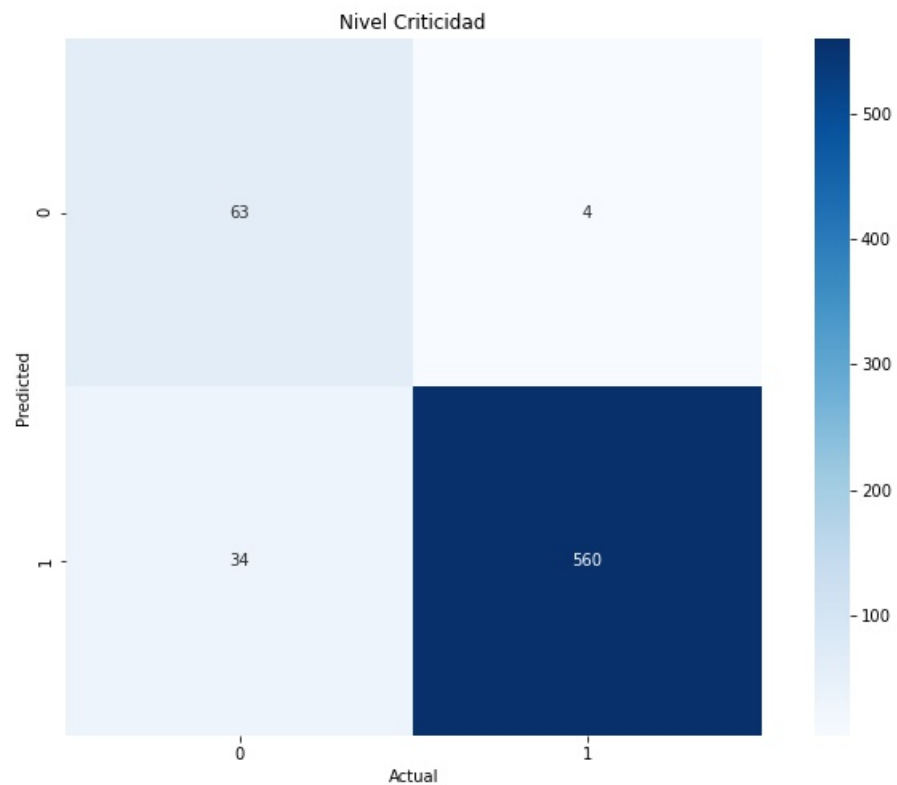


Figura 4.20: Matriz de confusión para conjunto de prueba con umbrales actualizados Alerta 1

Obteniendo los siguientes resultados:

- **Verdaderos positivos:** Son **560**, los cuales el modelo los predijo como 1 y la salida real es 1.
- **Verdaderos negativos:** Son **63**, los cuales el modelo los predijo como 0 y la salida real es 0.
- **Falsos positivos:** Son **4**, los cuales el modelo los predijo como 1 y la salida real es 0.
- **Falsos negativos:** Son **34**, los cuales el modelo los predijo como 0 y la salida real es 1.

Entonces de un total de 39 datos mal clasificados al aplicar el modelo sin agregar ningún umbral, al aplicar la reclasificación baja a un total de 38.

Clase	Niveles	Características
0	Medio y bajo	2559
1	Alto	84

Cuadro 4.12: Clase Alerta 2

#### 4.1.5. Tercera clasificación clase “Alerta 2”

- La clase alerta consta de las siguientes cantidades de características:
- Se dividió el conjunto de datos en tres partes quedando con las siguientes cantidades:

División	Características
Entrenamiento	1486
Validación	496
Prueba	661

Cuadro 4.13: División conjunto de datos para clase Alerta 2

- Al aplicar el modelo Random Forest al conjunto de validación se obtuvo la siguiente matriz de confusión:

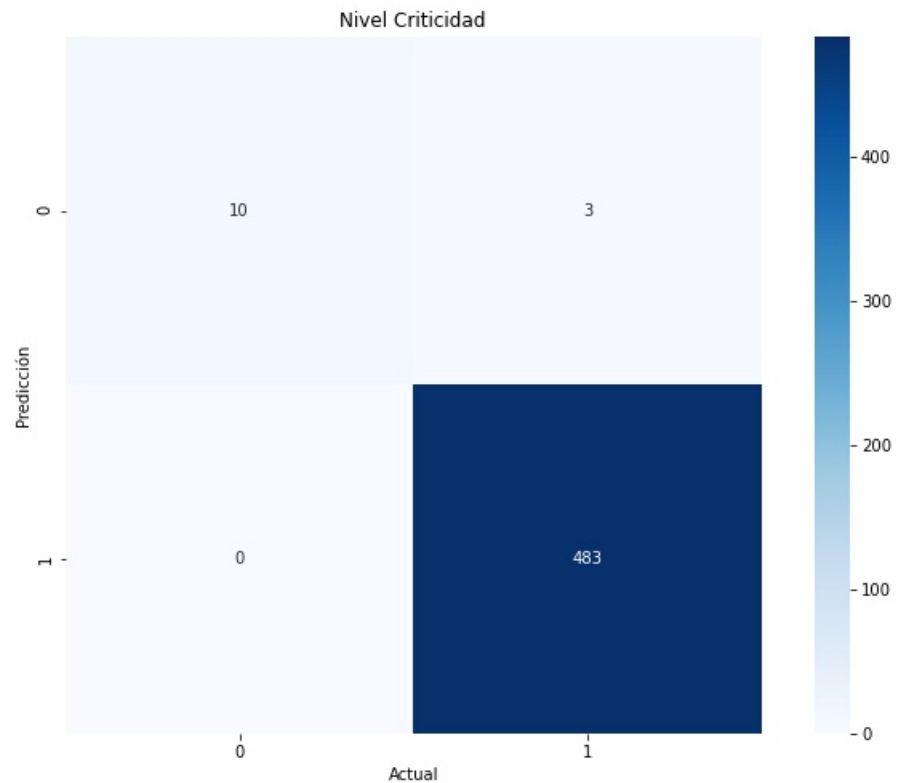


Figura 4.21: Matriz de confusión Alerta 2 conjunto validación

- **Verdaderos positivos:** Son **483**, los cuales el modelo los predijo como 1 y la salida real es 1.
- **Verdaderos negativos:** Son **30**, los cuales el modelo los predijo como 0 y la salida real es 0.
- **Falsos positivos:** Son **0**, los cuales el modelo los predijo como 1 y la salida real es 0.
- **Falsos negativos:** Son **3**, los cuales el modelo los predijo como 0 y la salida real es 1.

Además, se presenta el cuadro con las respectivas métricas de clasificación:

Clase	Precisión	Recall	f1-score
0	1.00	0.77	0.87
1	0.99	1.00	1.00

Cuadro 4.14: Métricas Alerta 2 conjunto validación

Como se puede apreciar en la matriz, solo se tienen 3 datos mal clasificados por lo que no se podrá llevar a cabo un análisis de ellos, se procede a aplicar el modelo al conjunto de prueba para comparar resultados:

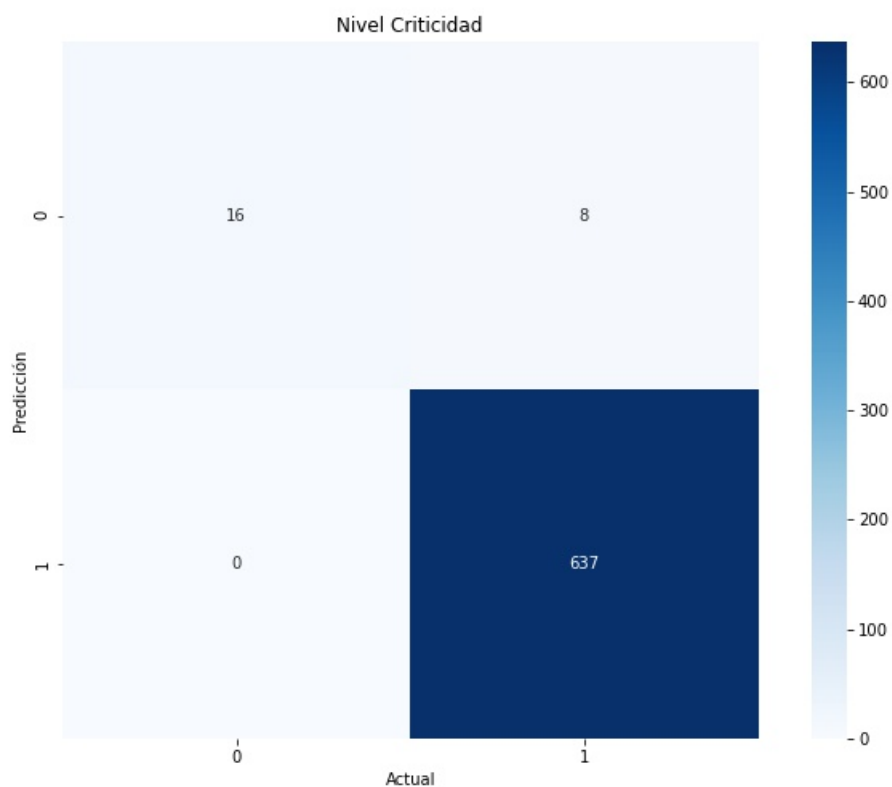


Figura 4.22: Matriz de confusión Alerta 2 conjunto prueba

- **Verdaderos positivos:** Son **637**, los cuales el modelo los predijo como 1 y la salida real es 1.
- **Verdaderos negativos:** Son **16**, los cuales el modelo los predijo como 0 y la salida real es 0.
- **Falsos positivos:** Son **0**, los cuales el modelo los predijo como 1 y la salida real es 0.
- **Falsos negativos:** Son **6**, los cuales el modelo los predijo como 0 y la salida real es 1.

Además, se presenta el cuadro con las respectivas métricas de clasificación:

Clase	Precisión	Recall	f1-score
0	1.00	0.67	0.80
1	0.99	1.00	0.99

Cuadro 4.15: Métricas Alerta 2 conjunto de prueba

De igual manera que en el conjunto de validación es muy baja la cantidad de datos mal clasificados por el modelo, por lo tanto, en esta clase no se llevará a cabo el análisis de mal clasificados ni la aplicación de una reclasificación.

De acuerdo con los resultados obtenidos, se desarrollaron los 3 escenarios planteados en la metodología, el primer escenario relacionado a la aplicación del modelo Random Forest sin ningún tipo de integración de umbral, generó resultados favorables de clasificación en cada una de las clases trabajadas.

Al incorporar la técnica de reclasificación con umbrales fijos que sería el segundo escenario propuesto, arrojó resultados negativos entorno a la clasificación debido a que no se apreció una mejora de los eventos mal clasificados.

Finalmente, al aplicar la técnica de reclasificación, pero con umbrales actualizados, en donde se agregó una componente aleatoria al momento de calcularlo, entregó resultados favorables en la clasificación disminuyendo de forma leve, pero cumpliendo el objetivo planteado.

Por lo tanto, la incorporación del concepto de marcador somático artificial ayuda a mejorar la clasificación en cada una de las clases analizadas.

# Capítulo 5

## Conclusiones

En este trabajo de titulación se presentó un mecanismo de clasificación para un conjunto de datos relacionados al área minera, donde la variable de interés fue el nivel de criticidad. Así mismo se propuso una técnica de machine learning que incorporó el concepto marcador somático artificial.

A partir de lo anterior, se puede concluir que la hipótesis planteada en el comienzo de la investigación - la incorporación del concepto de marcador somático artificial dentro de la técnica de *machine learning* mejora el desempeño de la clasificación - es favorable. Los resultados obtenidos de la clasificación con el uso del MSA indican una mejora en comparación a no incorporar el concepto de MSA.

Respecto a los objetivos propuestos: (i) examinar variables asociadas al riesgo ocupacional en una faena minera, (ii) diseñar un marcador somático artificial para técnicas de *machine learning*, (iii) integrar las variables de riesgo ocupacional con un marcador somático artificial en una técnica *machine learning* en el contexto del riesgo ocupacional, (iv) analizar los resultados experimentales derivados de la evaluación de la técnica de aprendizaje automático en el contexto de riesgo ocupacional, se llevaron a cabo en totalidad. Se implementó el MSA en la técnica de clasificación de bosque aleatorio, de manera que a través de la toma de decisiones se generara un cambio en la predicción con el fin de mejorarla.

Por último, en lo referente a la continuación de este trabajo de titulación, surgieron algunas consideraciones de líneas futuras a poder tratar, las cuales están relacionadas con el presente trabajo, donde se pueden abordar con mayor profundidad como: (i) buscar un método más profundo para el cálculo de umbrales actualizados con el fin de una mayor disminución del error de clasificación, (ii) incorporación de más variables relacionadas directamente a los incidentes.

# Referencias

- Altuwairqi, K., Jarraya, S. K., Allinjawi, A., and Hammami, M. (2021). A new emotion-based affective model to detect student's engagement. *Journal of King Saud University - Computer and Information Sciences*, 33(1):99–109.
- Ayhan, B. U. and Tokdemir, O. B. (2020). Accident analysis for construction safety using latent class clustering and artificial neural networks. *Journal of Construction Engineering and Management*, 146(3):04019114.
- Breiman, L. (2001). Random forests. *Machine Learning*, pages 5–32.
- Cabrera, D., Cubillos, C., Urra, E., and Mellado, R. (2020). Framework for incorporating artificial somatic markers in the decision-making of autonomous agents. *Applied Sciences*, 10(20).
- Cosentino, A., Azzollini, S., Depaula, P., and Castillo, S. (2016). Toma de decisión según racionalidad/afectividad, entrenamiento y saturación cultural en situaciones multiculturales: un estudio experimental con soldados para la paz. *Centro Interamericano de Investigaciones Psicológicas y Ciencias Afines*, 33(2).
- Damásio, A. (1994). *Dacartés'Error*. Grosset/Putnam Book. G. P. Putnam's Sons, Nueva York.
- Davoudi Kakhki, F., Freeman, S. A., and Mosher, G. A. (2019). Use of neural networks to identify safety prevention priorities in agro-manufacturing operations within commercial grain elevators. *Applied Sciences*, 9(21).
- De Wilde, P. (1997). *Neural network models*. 2nd ed.
- Duarte, J., Baptista, J., and Marques, A. (2019). Occupational accidents in the mining industry—a short review. *Studies in Systems, Decision and Control*, pages 61–69.
- Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The elements of statistical learning: Data mining, inference, and prediction*, second edition (springer series in statistics).
- Haykin, S. S. (2009). *Neural networks and learning machines*. Pearson Education, Upper Saddle River, NJ, third edition.
- Kowalczyk, Z., Czubenko, M., and Merta, T. (2019). Emotion monitoring system for drivers. *IFAC-PapersOnLine*, 52(8):200–205. 10th IFAC Symposium on Intelligent Autonomous Vehicles IAV 2019.
- Linguist, S. and Bartol, J. (2013). Two myths about somatic markers. *British Journal for the Philosophy of Science*, 64(3):455–484.

- Liu, P. (2004). *Fuzzy Neural Network Theory and Application*. World Scientific.
- Mosquera, R., Parra, L., Ledesma, A. J., and Bonilla (2021). Predicción de la accidentalidad laboral en la industria de pulpa y papel usando algoritmos de clasificación. *Información tecnológica*, 32:133 – 142.
- Márquez, M. d. R., Salguero, P., Paíno, S., and Alameda, J. R. (2012). La hipótesis del marcador somático y su nivel de incidencia en el proceso de toma de decisiones. *R.E.M.A. Revista electrónica de metodología aplicada*, 18(1):17–36.
- Poppa, T. and Bechara, A. (2018). The somatic marker hypothesis: revisiting the role of the ‘body-loop’ in decision-making. *Current Opinion in Behavioral Sciences*, 19:61–66. Emotion-cognition interactions.
- Rabeiy, R. E., ElTahlawi, M. R., and Boghdady, G. Y. (2018). Occupational health hazards in the sukari gold mine, egypt. *Journal of African Earth Sciences*, 146:209–216. Precambrian Geology of Egypt: Stratigraphy, Geodynamics, and Mineral Resources.
- Refaeilzadeh, P., Tang, L., and Liu, H. (2009). Cross-validation. pages 532–538.
- Sandor, S. and Gürvit, H. (2019). Development of somatic markers guiding decision-making along adolescence. *International Journal of Psychophysiology*, 137:82–91.
- Sarkar, S. and Maiti, J. (2020). Machine learning in occupational accident analysis: A review using science mapping approach with citation network analysis. *Safety Science*, 131:104900.
- Sarkar, S., Raj, R., Vinay, S., Maiti, J., and Pratihar, D. K. (2019a). An optimization-based decision tree approach for predicting slip-trip-fall accidents at work. *Safety Science*, 118:57–69.
- Sarkar, S., Vinay, S., Raj, R., Maiti, J., and Mitra, P. (2019b). Application of optimized machine learning techniques for prediction of occupational accidents. *Computers and Operations Research*, 106:210–224.
- Sarkar, S., Vinay, S., Raj, R., Maiti, J., and Mitra, P. (2019c). Application of optimized machine learning techniques for prediction of occupational accidents. *Computers and Operations Research*, 106:210–224.
- Sámamo-Ríos, M. L., Ijaz, S., Ruotsalainen, J., Breslin, F. C., Gummesson, K., and Verbeek, J. (2019). Occupational safety and health interventions to protect young workers from hazardous work – a scoping review. *Safety Science*, 113:389–403.
- Verkijika, S. F. (2020). An affective response model for understanding the acceptance of mobile payment systems. *Electronic Commerce Research and Applications*, 39:100905.
- Xu, Z., Saleh, J. H., and Subagia, R. (2020). Machine learning for helicopter accident analysis using supervised classification: Inference, prediction, and implications. *Reliability Engineering and System Safety*, 204:107210.
- Zhu, R., Hu, X., Hou, J., and Li, X. (2021). Application of machine learning techniques for predicting the consequences of construction accidents in china. *Process Safety and Environmental Protection*, 145:293–302.