



**PROTOTIPO DE SOFTWARE PARA MINERÍA DE IMÁGENES DEL SECTOR SALUD  
EN EL MARCO DE CIUDADES INTELIGENTES**

**EFRAÍN ALBERTO OVIEDO CARRASCAL**

**UNIVERSIDAD PONTIFICIA BOLIVARIANA - ESCUELA DE INGENIERÍAS  
FACULTAD DE INGENIERÍA EN TECNOLOGÍAS DE LA INFORMACIÓN Y LA  
COMUNICACIÓN  
MAESTRÍA EN TECNOLOGÍAS DE LA INFORMACIÓN Y LA COMUNICACIÓN  
MEDELLÍN – SEPTIEMBRE 2015**

**PROTOTIPO DE SOFTWARE PARA MINERÍA DE IMÁGENES DEL SECTOR SALUD  
EN EL MARCO DE CIUDADES INTELIGENTES**

**EFRAÍN ALBERTO OVIEDO CARRASCAL**

Trabajo de grado para optar al título de Maestría en Tecnologías de la Información y la  
Comunicación

Asesora

**ANA ISABEL OVIEDO**

PhD. En Ingeniería Electrónica

**UNIVERSIDAD PONTIFICIA BOLIVARIANA - ESCUELA DE INGENIERÍAS  
FACULTAD DE INGENIERÍA EN TECNOLOGÍAS DE LA INFORMACIÓN Y LA  
COMUNICACIÓN  
MAESTRÍA EN TECNOLOGÍAS DE LA INFORMACIÓN Y LA COMUNICACIÓN  
MEDELLÍN – SEPTIEMBRE 2015**

*DECLARACIÓN ORIGINALIDAD*

*“Declaro que esta tesis (o trabajo de grado) no ha sido presentada para optar a un título, ya sea en igual forma o con variaciones, en esta o cualquier otra universidad”. Art. 82 Régimen Discente de Formación Avanzada, Universidad Pontificia Bolivariana.*

*Efraín A. Oviedo*

FIRMA AUTOR \_\_\_\_\_

## **AGRADECIMIENTOS**

Al grupo de investigación GIDATI que me apoyó durante el desarrollo de la maestría y me dio la oportunidad de hacer parte de este proyecto.

A la doctora Ana Isabel Oviedo, asesora del proyecto, que se encargó de guiarme durante cada una de las etapas realizadas.

A mi familia que me dio el ejemplo para ser lo que soy hoy.

Y principalmente a Dios, quien me da las fuerzas para salir adelante y llena mi vida de bendiciones.

## RESUMEN

Entre las numerosas aplicaciones de la minería de datos se destacan los aportes a los servicios de salud en ciudades inteligentes. Dichas aplicaciones tienen por objetivo mejorar la calidad de vida de los ciudadanos, prevenir enfermedades, facilitar la toma de decisiones y analizar datos provenientes de las instituciones de salud. Sin embargo, un gran porcentaje de los datos requeridos para análisis de las enfermedades se encuentran representados en imágenes o texto. Recientemente algunas herramientas de minería de datos han incluido dentro de sus funcionalidades procesamiento de texto, pero no se ha encontrado evidencia de herramientas que integren el procesamiento de imágenes. Con el objetivo de apoyar el desarrollo de ciudades inteligentes, en este trabajo se presenta un prototipo de software para minería de imágenes provenientes del sector salud. El prototipo desarrollado permite cargar un conjunto de imágenes y realizar dos tipos de análisis de minería: análisis de tipo predictivo y análisis de tipo descriptivo. El prototipo de software desarrollado es evaluado mediante dos casos de estudio, el primero de ellos realiza un análisis predictivo sobre una base de datos de mamografías entre las cuales se encuentran casos normales y casos con anomalías detectadas, el segundo caso de estudio realiza análisis descriptivo sobre un conjunto de imágenes que incluye tres diferentes tipos de exámenes médicos resonancia magnética, radiografía del tórax y mamografías.

**Palabras clave:** Servicio de salud en ciudades inteligentes, minería de imágenes, clasificación supervisada de imágenes, clustering de imágenes.

## ABSTRACT

Among the many applications of data mining are included contributions to citizens life in smart cities. These applications are intended to improve the quality of life of citizens, prevent diseases, facilitate decision making and the analyzing data from health institutions. However, a large percentage of the data required for analysis of diseases is found in images or text. Recently some data mining tools have included in their features the text processing, but no evidence has been found of tools that integrate image processing. With the aim of supporting the development of smart cities, this paper work a software prototype of image mining from the health sector. The developed prototype allows you to load a set of images and perform two types of mining analysis: predictive and descriptive analysis. The developed prototype is evaluated by two case studies, the first one performed a predictive analysis on a mammograms database, and the second case study performed a descriptive analysis on a set of images which includes three different types of medical exams: MRI, chest X-ray and mammography.

Key words: Health Service in smart cities, mining images, supervised image classification, clustering of images.

## GLOSARIO

**AWS** (Amazon Web Services): Servicios Web ofrecidos por Amazon que ofrecen soluciones relacionadas en la computación en la nube.

**EC2** (Elastic Compute Cloud): Servicio de AWS que ofrece capacidad informática en la nube, utilizadas generalmente como servidores.

**IaaS** (Infrastructure as a Service): Categoría de la computación en la nube que ofrece infraestructura de cómputo como un servicio en la nube.

**IDC** (International Data Corporation): Corporación internacional dedicada al análisis e inteligencia de mercados en el sector de la informática y de las comunicaciones.

**IoT** (Internet of Things): Interconexión de los objetos cotidianos al Internet.

**mSalud**: Servicios de salud utilizando dispositivos móviles

**PaaS** (Platform as a Service): Categoría de la computación en la nube que ofrece plataformas de desarrollo en un ambiente virtualizado.

**S3** (Simple Storage Services): Servicio de AWS que ofrece capacidad de almacenamiento de forma segura, duradera y fácil de escalar.

**SaaS** (Software as a Services): Categoría de la computación en la nube que ofrece aplicaciones de software como un servicio en la nube.

**TICs** (Tecnologías de la Información y la comunicación): Conjunto de tecnologías utilizadas para la gestión, transformación y comunicación de la información.

**UML** (Unified Modeling Language): Lenguaje de modelado comúnmente utilizado para el desarrollo de software.

**Curva ROC** (Receiver Operating Characteristic curve): Curva característica operativa del receptor, consiste en una representación gráfica que permite evaluar un modelo de clasificación en minería de datos.

## CONTENIDO

1.	INTRODUCCIÓN.....	1
1.1	Minería de Imágenes del Sector Salud.....	1
1.2	Aportes.....	2
1.3	Divulgación.....	3
1.4	Organización de la Tesis.....	4
	<b>Parte I: Análisis del Problema.....</b>	<b>5</b>
1.	DESCRIPCIÓN DEL PROBLEMA .....	6
2.	JUSTIFICACIÓN .....	8
3.	OBJETIVOS.....	9
	<b>Parte II: Revisión Literaria.....</b>	<b>10</b>
1.	MARCO CONCEPTUAL.....	11
1.1	Ciudades Inteligentes.....	11
1.2	Minería de Datos.....	14
1.3	Procesamiento de Imágenes.....	21
2.	ESTADO DEL ARTE .....	25
2.1	Minería de datos para el desarrollo de ciudades inteligentes.....	25
2.2	Minería de imágenes en el área de la salud .....	25
2.3	Minería de imágenes en los diagnósticos médicos.....	27
2.4	Discusión .....	28
	<b>Parte III: Solución Propuesta.....</b>	<b>30</b>
1.	MINERÍA DE IMÁGENES DEL SECTOR SALUD .....	31
1.1	Requisitos de la Minería de Imágenes del Sector Salud .....	31
1.2	Metodología de Minería Multimedia KDM (Knowledge Discovery in Multimedia) .....	32
2.	DESARROLLO DE SOFTWARE.....	36
2.1	Metodología de desarrollo de software .....	36

2.2	Lenguaje de programación .....	38
2.3	Base de datos .....	38
2.4	Librerías utilizadas.....	38
<b>Parte IV: Documento de Requisitos de Software IEEE830 .....</b>		<b>41</b>
1.	PRESENTACIÓN .....	42
1.1	Propósito.....	42
1.2	Alcance.....	42
1.3	Personal involucrado.....	43
1.4	Resumen .....	44
2.	DESCRIPCIÓN GENERAL.....	45
2.1	Perspectiva del producto .....	45
2.2	Funcionalidad del producto .....	46
2.3	Características de los usuarios .....	46
2.4	Restricciones .....	47
2.5	Suposiciones y dependencias .....	47
2.6	Evolución previsible del sistema .....	47
3.	REQUISITOS ESPECÍFICOS.....	49
3.1	Requisitos comunes de los interfaces.....	49
3.2	Interfaces de usuario .....	49
3.3	Interfaces de hardware.....	51
3.4	Interfaces de software .....	51
3.5	Interfaces de comunicación .....	52
3.6	Requisitos funcionales .....	52
3.7	Requisitos no funcionales .....	58
<b>Parte V: Documento de Análisis y Diseño con UML.....</b>		<b>60</b>
1.	DESCRIPCIÓN DE LA HERRAMIENTA SOFTWARE .....	61
2.	DIAGRAMA DE CASOS DE USO.....	62
3.	DIAGRAMA DE CLASES .....	66
4.	DIAGRAMA DE SECUENCIA .....	68
5.	DIAGRAMA DE ACTIVIDADES .....	70

<b>Parte VI: Documento de Pruebas .....</b>	<b>72</b>
1. PLAN DE PRUEBAS.....	73
2. PRUEBAS DEL PROTOTIPO DE SOFTWARE .....	75
<b>Parte VII: Casos de Estudio, Conclusiones y Trabajo Futuro.....</b>	<b>85</b>
1. CASO DE ESTUDIO 1: PREDICCIÓN DE ANORMALIDADES EN MAMOGRAFÍAS.....	86
1.1 Descripción del caso de estudio.....	86
1.2 Aplicación de la metodología KDM al caso de estudio .....	86
1.2.1 Selección de las imágenes.....	86
1.2.2 Preprocesamiento de las imágenes .....	87
1.2.3 Indexamiento .....	89
1.2.4 Modelamiento.....	89
1.2.5 Evaluación .....	89
2. CASO DE ESTUDIO 2: CLUSTERING DE TIPOS DE EXÁMENES MÉDICOS.....	92
2.1 Descripción del caso de estudio.....	92
2.2 Aplicación de la metodología KDM al caso de estudio .....	92
2.2.1 Selección de las imágenes.....	92
2.2.2 Preprocesamiento de las imágenes .....	93
2.2.3 Indexamiento .....	94
2.2.4 Modelamiento.....	94
2.2.5 Evaluación .....	94
3. CONCLUSIONES .....	96
4. TRABAJO FUTURO .....	98
<b>BIBLIOGRAFÍA .....</b>	<b>99</b>
<b>ANEXOS .....</b>	<b>106</b>
1. Manual de usuario del prototipo de software.....	106

## LISTA DE ILUSTRACIONES

Ilustración 1: Metodología propuesta KDM.....	33
Ilustración 2: Mapa conceptual: minería de datos para la salud .....	45
Ilustración 3: Diagrama de casos de uso.....	62
Ilustración 4: Funciones del administrador.....	63
Ilustración 5: Editar perfil.....	64
Ilustración 6: Aplicación de la minería de imágenes.....	65
Ilustración 7: Tipos de análisis y técnicas utilizadas .....	65
Ilustración 8: Diagrama de clases.....	66
Ilustración 9: Diagrama de secuencia .....	68
Ilustración 10: Diagrama de actividades.....	70
Ilustración 11: Pantalla de acceso al sistema.....	75
Ilustración 12: Pruebas de acceso al sistema.....	76
Ilustración 13: Pantalla de inicio .....	76
Ilustración 14: Pantalla de administración de usuarios .....	77
Ilustración 15: Crear y editar usuarios .....	78
Ilustración 16: Editar perfil.....	79
Ilustración 17: Pantalla de configuraciones .....	79
Ilustración 18: Crear configuración.....	80
Ilustración 19: Pantalla preparar imágenes .....	81
Ilustración 20: Pantalla extraer características.....	82
Ilustración 21: Pantalla de análisis predictivo- Entrenamiento .....	83
Ilustración 22: Pantalla de análisis descriptivo .....	84
Ilustración 23: Mamografías de la base de datos MIAS.....	87
Ilustración 24: Segmentación de Mamografías .....	88
Ilustración 25: Árbol de decisión caso de estudio 1 .....	91
Ilustración 26: Imágenes caso de estudio 2.....	93
Ilustración 27: Índice de la silueta caso de estudio 2 .....	95

## LISTA DE TABLAS

Tabla 1: Conjunto de entrenamiento utilizado en el análisis predictivo discreto .....	15
Tabla 2: Conjunto de entrenamiento utilizado en el análisis predictivo continuo .....	15
Tabla 3: Conjunto de datos utilizados en el análisis descriptivo .....	16
Tabla 4: Pre procesamiento de Multimedia .....	34
Tabla 5: Indexamiento de multimedia .....	34
Tabla 6: Diferencias entre metodologías ágiles y tradicionales.....	37
Tabla 7: Personal involucrado en el proyecto.....	43
Tabla 8: Características de los usuarios de prototipo de software.....	47
Tabla 9: Requisito funcional 1.....	52
Tabla 10: Requisito funcional 2.....	53
Tabla 11: Requisito funcional 3.....	54
Tabla 12: Requisito funcional 4.....	54
Tabla 13: Requisito funcional 5.....	55
Tabla 14: Requisito funcional 6.....	55
Tabla 15: Requisito funcional 7.....	56
Tabla 16: Requisito funcional 8.....	56
Tabla 17: Requisito funcional 9.....	57
Tabla 18: Matriz de confusión caso de estudio 1 .....	90
Tabla 19. Medidas de evaluación del clasificador caso de estudio 1 .....	90

# 1. INTRODUCCIÓN

## 1.1 Minería de Imágenes del Sector Salud

Las ciudades inteligentes tienen como objetivo mejorar la calidad de vida de los ciudadanos. Para ello se utilizan las tecnologías de la información y la comunicación (TICs) como herramientas para transformar y mejorar los procesos y actividades de la administración (Rodríguez & Gil, 2014). En los últimos años, este concepto ha tenido una gran acogida alrededor del mundo y se ha elevado considerablemente el número de ciudades que se han preocupado por realizar actividades de investigación y desarrollo al respecto. Como muestra de esta tendencia, en (Chen, 2014) se realiza un estudio de los planes de desarrollo de 415 ciudades inteligentes de diferentes países del mundo.

El *International Data Corporation* (IDC)<sup>1</sup> propone algunas áreas en las cuales las ciudades inteligentes deben centrar sus esfuerzos (Achaerandio, Curto, Bigliani, & Gallotti, 2012). Estas áreas son: gobierno, construcción, movilidad, energía, medio ambiente y servicios. El servicio de salud en las ciudades inteligentes se enmarca en las áreas de gobierno / servicio y se preocupa por mejorar la calidad de vida de los ciudadanos, para ello se requiere realizar actividades de prevención, diagnóstico, análisis y toma de decisiones que contribuyan a mejorar la salud de las personas. Para realizar estas actividades el personal de la salud utiliza diferentes recursos como: exámenes de laboratorio, imágenes diagnósticas y exámenes especializados (Ortiz, 2013).

Buscando aportar al servicio de salud en una ciudad inteligente, en la actualidad se desarrollan numerosas aplicaciones que permiten analizar datos provenientes de las instituciones de salud para apoyar la toma de decisiones en la ciudad. Dichas aplicaciones son desarrolladas con técnicas de minería de datos, las cuales permiten descubrir conocimiento en grandes volúmenes de datos.

La minería de datos es un campo de las ciencias de la computación que utiliza las técnicas de la inteligencia artificial y la estadística para descubrir patrones a partir de grandes volúmenes de datos (Mena, 1999). Típicamente un proceso de minería de datos implica varias etapas: selección del conjunto de datos, análisis de los datos, transformación de los datos de entrada, extracción de conocimiento, interpretación y evaluación de los datos. En la actualidad, la minería de datos se

---

<sup>1</sup> Empresa multinacional dedicada a la investigación de mercados en las áreas de tecnología de la información y telecomunicaciones < <https://www.idc.com/>>

utiliza en diversas áreas entre ellas: finanzas, banca, negocios, canasta de mercado y la salud. Este trabajo se enfoca específicamente en la aplicación de la minería de datos a la salud.

Teniendo en cuenta que gran parte de la información disponible para analizar la salud de un paciente se encuentran representada en imágenes, en este trabajo se aborda el problema de incorporar la información contenida en imágenes diagnósticas en procesos de minería de datos, lo cual recibe el nombre de minería de imágenes.

La minería de imágenes permite descubrir patrones característicos a partir de un conjunto de imágenes (Fernández, Miranda, Guerrero, & Piccoli, 2006). La evolución de esta área se debe en parte a los avances tecnológicos que han facilitado la captura y almacenamiento de imágenes digitales en las llamadas bases de datos de imágenes. El proceso de la minería de imágenes, generalmente, involucra varias etapas: recuperación de imágenes, indexamiento de imágenes, modelamiento y evaluación.

## **1.2 Aportes**

En este proyecto se desarrolla una herramienta software para realizar minería de imágenes del sector salud. A continuación se presentan los aportes del trabajo.

- **Levantamiento de requisitos desde el estado del arte**

El levantamiento de requisitos de software se realiza básicamente para determinar las necesidades, objetivos y actividades que debe realizar la herramienta, y comúnmente son especificados por el usuario. Los requisitos de la herramienta desarrollada en este trabajo tiene la particularidad de que dependen del estado del arte revisado, ya que no se encuentran antecedentes de herramientas de minería para imágenes. Este trabajo hace parte de un proyecto realizado en el grupo de investigación GIDATI llamado “Plataforma de Minería de Datos Estructurados y No Estructurados - Caso de Estudio Salud Pública”. El estado del arte permitió definir: 1) tipos de análisis más comunes en la minería de datos del sector salud, 2) formas de caracterización de imágenes del sector salud, 3) técnicas de minería más utilizadas y 4) medidas de evaluación.

- **Metodología de minería multimedia**

Aunque las metodologías de minería de datos en su mayoría incluyen fases de preparación de los datos, dicha preparación sólo incluye análisis estadísticos y transformaciones. Para analizar datos multimedia es necesario incluir etapas de pre procesamiento e indexamiento, donde se pueda representar la información multimedia en vectores de características que puedan ser procesados por las técnicas de minería de datos. En este trabajo se realiza una modificación a la metodología

planteada en KDD (Knowledge Discovery in Databases), la cual fue llamada KDM (Knowledge Discovery in Multimedia), esta metodología incluye etapas de: selección, preprocesamiento, indexamiento, transformación, modelamiento y evaluación.

- **Desarrollo de un prototipo de herramienta software**

Dentro de las investigaciones revisadas sobre minería de datos aplicada a la salud, se ha notado cierta preferencia por la herramienta WEKA para aplicar las técnicas de la minería de datos. Aunque WEKA recientemente incluye métodos para el procesamiento de textos en español, no presenta soporte para minería de imágenes. Otras herramientas como RapidMiner, SPSS Modeler y SAS Analytics también incluyen módulos para minería de textos, pero no para minería de imágenes. En este trabajo se desarrolla un prototipo de una herramienta software para minería de imágenes que incluye las siguientes etapas: preparación de datos, selección de características, aplicación de la minería y evaluación.

- **Aplicación hacia otras áreas**

A pesar de que el prototipo de software ha sido diseñado y desarrollado con el objetivo de aplicar minería sobre imágenes médicas, también es posible utilizarlo para realizar minería sobre imágenes de cualquier otra área. Las imágenes utilizadas para ello deben estar en alguno de los siguientes formatos: JPG, PNG o BMP.

### 1.3 Divulgación

Durante la ejecución del trabajo, se desarrollaron los siguientes documentos:

- **Artículo en Revista:**

MINERÍA DE DATOS: APORTES Y TENDENCIAS EN EL SERVICIO DE SALUD DE CIUDADES INTELIGENTES. Autores: Efraín Alberto Oviedo, Ana Isabel Oviedo y Gloria Liliana Vélez. Revista Politécnica enero-junio de 2015, vol 11, ISSN 1900-2351(Impreso) - ISSN 2256-5353 (En línea), pág 111 – 120. Disponible web:

<[http://politecnicojic.edu.co/index.php?option=com\\_content&view=article&id=362&Itemid=339](http://politecnicojic.edu.co/index.php?option=com_content&view=article&id=362&Itemid=339)>

- **Artículo en Revista:**  
HACIA LA CONSTRUCCIÓN DE UNA METODOLOGÍA Y UNA HERRAMIENTA DE MINERÍA MULTIMEDIA. Artículo en desarrollo.
- **Registro de Software:**  
PLATAFORMA DE MINERÍA MULTIMEDIA. En formulación.

#### **1.4 Organización de la Tesis**

La organización del documento es la siguiente. En la parte I se presenta el análisis del problema, abordando la descripción, justificación y los objetivos planteados para el desarrollo del proyecto. La parte II consiste en una revisión literaria, que incluye el marco conceptual y el estado del arte. En la parte III se plantea la solución propuesta. En la parte IV se describen los requisitos del proyecto utilizando el estándar de requisitos de software IEEE830. En la parte V se presenta el documento de análisis y diseño basado en el lenguaje de modelado unificado UML. En la parte VI se presenta el documento de pruebas. Finalmente en la parte VII se presentan dos casos de estudio, conclusiones y trabajo futuro.

# Parte I: Análisis del Problema

## 1. DESCRIPCIÓN DEL PROBLEMA

De acuerdo al banco mundial (Banco Mundial, 2015), se estima que para el 2030 el 60% de la población mundial vivirá en áreas urbanas, con un aumento diario de 180.000 personas a la población urbana. Este aumento trae consigo un gran reto para las ciudades que deben prepararse para abordar esta situación, en temas trascendentales como transporte, alimentación, vivienda, salud, entre otros.

Para afrontar este tipo de situaciones, las ciudades están experimentando un cambio hacia lo que conocemos como ciudades inteligentes, donde la ciudad empieza a ser vista como un medio innovador que facilita el desarrollo y progreso del país (Vivas, Britos, García, & Cambarieri, 2013). Uno de los principales objetivos que se buscan con ciudades inteligentes consiste en utilizar las tecnologías de la información y la comunicación (TICs) como motor de desarrollo para el bienestar de sus habitantes (Rodríguez & Gil, 2014). Siendo los servicios de salud un área de gran interés.

Dentro de los servicios de salud, existe la necesidad de diagnosticar el estado de los pacientes, en esta labor cada vez es más común utilizar imágenes diagnósticas que permiten al personal de la salud, visualizar e identificar posibles anomalías en el cuerpo del paciente. Es el caso de radiografías, ecografías, mamografías, escenografías, resonancias magnéticas, entre otros exámenes, que requieren el análisis de una imagen para diagnosticar el estado de salud del paciente. La minería de imágenes puede ser de gran utilidad para apoyar el análisis de las imágenes diagnósticas.

La minería de imágenes permite descubrir patrones característicos a partir de un gran conjunto de imágenes, con estos patrones es posible apoyar los servicios del área de salud, donde se cuenta con grandes cantidades de imágenes en las cuales hay información implícita de gran utilidad. Esta información se extrae mediante técnicas de procesamiento de imágenes.

Para la implementación de sistemas de minería de imágenes, se requiere de conocimientos en áreas como la minería de datos y el procesamiento digital de imágenes, y en general, el personal de la salud, que es quien tiene acceso a las imágenes diagnósticas, no cuenta con estos conocimientos, lo que dificulta la implementación de dichos sistemas ya que no se ha encontrado evidencia de plataformas de minería con soporte a procesamiento de imágenes que faciliten esta labor.

De esta manera, en el proyecto se integran las áreas de minería de datos (Kumari & Sunila, 2011), procesamiento de imágenes (Fernández, Miranda, Guerrero, & Piccoli, 2006) y servicios de salud en ciudades inteligentes, para desarrollar un prototipo de software de minería de imágenes que permita analizar imágenes diagnósticas provenientes del sector salud. Se espera que dicho prototipo apoye la toma de decisiones relacionadas con salud, mediante tareas de minería de

imágenes.

La aplicación de la minería de imágenes al sector salud debe superar un desafío: la ausencia de herramientas de minería con soporte a procesamiento de imágenes. Las plataformas tradicionales de minería de datos sólo soportan datos estructurados y algunas de ellas están empezando a soportar texto. Una de las herramientas más destacadas es WEKA, la cual ha sido ampliamente utilizada en aplicaciones de la salud, sin embargo no cuenta con soporte para minería de imágenes, de hecho en la revisión literaria no se encontró ninguna plataforma de minería con soporte a imágenes.

La construcción de una herramienta de minería de imágenes de la salud debe superar los siguientes problemas:

- Ausencia de metodologías que guíen el proceso de minería de imágenes para definir los requisitos del software.
- No se tienen definidas técnicas y métodos de minería específicos para el procesamiento de imágenes, lo que dificulta la definición de requisitos del software.
- El procesamiento digital de imágenes presenta cierto grado de dificultad debido entre otras cosas, a la gran cantidad de información disponible en una imagen (Peláez, 2014).
- Las imágenes de la salud comúnmente contienen información sobre el paciente y su diagnóstico. Dicha información actualmente es protegida por normas nacionales impidiendo su análisis. Esto implica que dicha información debe ser eliminada de las imágenes antes de realizar procesos de minería.

Como aporte a la solución, se presenta en este trabajo el desarrollo de una herramienta software para minería de imágenes del sector salud.

## 2. JUSTIFICACIÓN

Para resolver las necesidades del sector salud en las ciudades inteligentes es común que se utilicen técnicas como la minería de datos con el fin de generar nuevo conocimiento relacionado con prevención, diagnóstico, tratamiento de enfermedades y procedimientos de recuperación (Timarán & Yépez, 2012). En los últimos años se ha utilizado minería de datos para abordar problemas de la salud en áreas como: tratamiento de cáncer (Timarán & Yépez, 2012), enfermedades cardiovasculares (Solarte & Castro, 2012), enfermedades tumorales primarias (Shaikh, 2014), hipertensión arterial (Pérez, 2012), efectos del dengue (Amin, Takib, Raza, & Javed, 2014), diabetes (Hernández, 2014), esclerosis múltiple (Filipuzzi, Rodrigo, Graffigna, Isoardi, & Noceti, 2012), entre otras.

Teniendo en cuenta la importancia de las imágenes en el sector salud, disponer un software que facilite la implementación de la minería de imágenes puede ayudar a mejorar la efectividad de los diagnósticos tempranos y a evitar remisiones innecesarias de pacientes. Estas dos tareas son muy importantes ya que hay enfermedades como el caso del Alzheimer en las cuales los métodos actuales de diagnóstico temprano no superan el 70% de efectividad (Romero, 2011). Adicionalmente, hoy en día es común que se realicen remisión de pacientes a entidades de salud especializadas para realizar diagnósticos (Ortiz, 2013), estos traslados podrían evitarse si se dispone de herramientas de software para realizar diagnósticos tempranos a partir de imágenes.

La minería de imágenes es mucho más que una simple extensión de la minería de datos al campo de las imágenes. En la minería de imágenes se realiza la extracción de patrones de un conjunto de imágenes, no se trata simplemente de extraer características específicas de una única imagen, tampoco se trata de extraer las características más relevantes de las imágenes para luego aplicar técnicas de la minería de datos. Un sistema de minería de imágenes debe realizar el procesamiento de imágenes, extraer características, realizar el indexado, aplicar métodos de minería de datos y evaluar los resultados obtenidos.

De acuerdo con lo anterior y con el fin de contribuir al mejoramiento del servicio de salud en las ciudades inteligentes, en este trabajo se plantea el desarrollo de un prototipo de software que permita realizar tareas de minería de imágenes provenientes del sector salud, específicamente las tareas de clasificación y agrupación de imágenes. La evaluación de la herramienta se realizará mediante un caso de estudio, sin embargo el objetivo de este proyecto no es hacer un estudio de minería para terceros, este es un proyecto de desarrollo de software.

### 3. OBJETIVOS

#### **Objetivo General**

Desarrollar un prototipo de software para minería de imágenes provenientes del sector salud como apoyo a los servicios ofrecidos en Ciudades Inteligentes.

#### **Objetivos Específicos**

- Diseñar una herramienta software para minería de imágenes del sector salud en el marco de ciudades inteligentes.
- Desarrollar un prototipo de la herramienta software.
- Validar el prototipo de software con un caso de estudio.

# Parte II: Revisión Literaria

## 1. MARCO CONCEPTUAL

La minería de imágenes es un área que ha despertado mucho interés en los últimos años, usualmente se define como la búsqueda de información no trivial en un importante conjunto de imágenes (Thamilselvan & Sathiaseelanb, 2015). Entre sus usos más frecuentes se encuentran la clasificación de imágenes supervisada y no supervisada (Clustering). Para realizar estas actividades se extraen características de las imágenes relacionadas con su color, forma y textura (Chidambaranathan, 2015). Algunos autores visualizan la minería de imágenes como la suma de dos áreas: la minería de datos y el procesamiento de imágenes (Dua, 2014).

En esta sección se presenta una revisión literaria en varias líneas de trabajo: ciudades inteligentes como contexto general, minería de datos como área general y procesamiento de imágenes como área de interés.

### 1.1 Ciudades Inteligentes

En los últimos años el concepto de ciudades inteligentes ha tenido una gran acogida alrededor del mundo, a tal punto que, se ha elevado considerablemente el número de ciudades que se han preocupado por realizar actividades de investigación y desarrollo al respecto. Como muestra de ello en (Chen, 2014) se realiza un estudio de los planes de desarrollo de 415 ciudades inteligentes de diferentes países del mundo.

En general, el objetivo de las ciudades inteligentes es mejorar la calidad de vida de los ciudadanos, para ello se utilizan las tecnologías de la información y la comunicación (TICs) como herramientas para transformar y mejorar los procesos y actividades de la administración (Rodríguez & Gil, 2014).

El *International Data Corporation* (IDC) propone algunas áreas en las cuales las ciudades inteligentes deben centrar sus esfuerzos (Achaerandio, Curto, Bigliani, & Gallotti, 2012). Estas áreas son: gobierno, construcción, movilidad, energía, medio ambiente y servicios. Los servicios de salud en las ciudades inteligentes se enmarcan en las áreas de gobierno y servicio, buscando prevenir enfermedades y mejorar la salud de los ciudadanos. En la actualidad se desarrollan numerosas aplicaciones que permiten analizar datos provenientes de las instituciones de salud para apoyar la toma de decisiones en la ciudad. Algunas de estas aplicaciones son desarrolladas con técnicas de minería de datos, las cuales permiten descubrir conocimiento en grandes volúmenes de datos.

### **1.1.1 Tendencias en ciudades inteligentes**

En (Namiot & Schneps-Schneppe, 2013) se presenta una mirada al software de las ciudades inteligentes desde el punto de vista del desarrollo. En este estudio se resalta una definición muy interesante sobre ciudades inteligentes, planteada por Forrester, una empresa líder en investigación tecnológica y de mercado, según esta definición, en una ciudad inteligente se requiere del uso combinado de factores como sistemas de software, infraestructura de servidores, infraestructura de red y dispositivos del cliente con el fin de conectar los componentes críticos de infraestructura de la ciudad y servicios.

Los componentes involucrados en esta definición permiten visualizar algunas tendencias en cuanto a tecnologías facilitadoras que se vienen utilizando en el desarrollo de las ciudades inteligentes, a continuación se mencionan dichas tendencias.

- **Internet de las cosas**

El Internet de las cosas, IoT (Internet of Things), es un nuevo paradigma donde todos los electrodomésticos, sensores y actuadores se interconectan entre sí a través de Internet, de modo tal que se pueda tener acceso a ellos en cualquier momento y en cualquier lugar (Gómez, Huete, Hoyos, Perez, & Grigori, 2013).

Este concepto es retomado en (Kamel & Al-Shorbaji, 2014) donde se mencionan algunas aplicaciones del Internet de las cosas en las ciudades inteligentes. Una de estas aplicaciones se refiere a los servicios de salud tele-asistidos, donde se utilizan sensores inalámbricos para monitorear la salud de las personas, incluso se disponen de prendas de vestir inteligentes que contienen dichos sensores y envían la información censada a servidores remotos a través de Internet.

Otro ejemplo de la influencia del Internet de las cosas en las ciudades inteligentes es presentado en (Zanella, Bui, Castellani, Vangelista, & Zorzi, 2014) donde se resumen los resultados de la implementación de un proyecto de ciudad Inteligente en la ciudad de Italiana Padova. Es tal la relación entre ambos conceptos, que incluso en (Jin, Gubbi, Marusic, & Palaniswami, 2014) se presenta un Framework para desarrollar ciudades inteligentes a través del concepto de Internet de las cosas.

- **Computación en la nube**

Un complemento muy interesante para el concepto de Internet de las cosas es presentado en (Petrolo, Loscri, & Mitton, 2014), se trata de la computación en la nube, un nuevo paradigma que permite ofrecer servicios de computación a través de Internet. Con la computación en la nube los usuarios pueden acceder a los servicios en cualquier momento a través de una conexión a Internet, una característica muy interesante de estos servicios es que se pueden ofrecer bajo demanda, esto quiere decir que se paga solamente por el

consumo efectuado y los recursos destinados para ofrecer el servicio se comportan de forma flexible y adaptativa de acuerdo a la demanda. Estos servicios se pueden dividir en tres capas:

- **Infraestructura como servicio (IaaS):** La infraestructura como servicio corresponde a la capa inferior y consiste en ofrecer componentes de hardware, tales como: capacidad de cómputo y espacio de almacenamiento, a través de servicios en la nube. Hoy en día existen diversas empresas dedicadas a suministrar este tipo de servicios (Serrano, Gallardo, & Hernantes, 2015), ejemplo de ello son los productos de Amazon Web Services (AWS) que cuenta con la solución Amazon Elastic Compute Cloud (Amazon EC2) para ofrecer capacidades de cómputo y con el servicio Amazon Simple Storage Service (Amazon S3) para satisfacer necesidades de almacenamiento en la nube. En diversas ocasiones, AWS ha sido considerado como líder en el cuadrante mágico de Gartner (Leong, Toombs, Gill, Petri, & Haynes, 2015) en cuanto a la infraestructura como servicio
- **Plataforma como servicios (PaaS):** Esta es la capa del medio y básicamente se enfoca en realizar una abstracción de un ambiente de desarrollo previamente configurado para trabajar con determinar arquitectura o tecnología y ofrecerlo como un servicio en la nube. Gracias a este tipo de servicios los desarrolladores pueden realizar sus actividades desde cualquier computador con acceso a internet sin preocuparse por instalar el software necesario para configurar el ambiente de desarrollo que requieren. Una de las ofertas de plataforma como servicio que ha tenido un buen recibimiento en los últimos años es Google App Engine (Beyer, Dresler, & Wendler, 2014), se trata de un servicio en la nube que ofrece un entorno de desarrollo donde se puede realizar y desplegar aplicaciones en lenguajes como Python y Java.
- **Software como servicio (SaaS):** Hace referencia a la capa de más alto nivel y consiste en ofrecer una aplicación como un servicio, al cual se accede a través de Internet. Bajo este modelo, el cliente ya no tiene la necesidad de instalar la aplicación en cada uno de sus equipos, simplemente accede a ella a través de cualquier dispositivo que tenga conexión a Internet. El modelo de software como servicio ha tenido un gran crecimiento y una gran aceptación a nivel mundial, debido entre otros factores a un cambio de modelo negocio de compra de licencia, utilizada tradicionalmente, a pago por uso, cambio que ha permitido una salida rápida a producción (Perozo & Boscán, 2014).

- **Aplicaciones móviles**

Otra tendencia que se viene presentando en las ciudades inteligentes, consiste en ofrecer servicios a través de aplicaciones móviles. Como el presentado en (Latorre, 2014) donde se

diseña una aplicación para tal fin sobre el sistema operativo Android, argumentando que en el 2013 el 81% de los teléfonos inteligentes del mundo utilizaban este sistema operativo. En dicho estudio se resalta la importancia que tiene hoy en día el desarrollo de aplicaciones para las ciudades inteligentes y la forma como los dispositivos móviles pueden apoyar esta labor, específicamente presentan dos tipos de servicios para ciudades inteligentes a través de dispositivos móviles, el primero de ellos consiste en un buscador que permite ubicar en tiempo real lugares libres para parquear, y el segundo es un notificador instantáneo de emergencias.

Los dispositivos móviles presentan cierto tipo de ventajas para el desarrollo de aplicaciones dentro el contexto de ciudades inteligentes, por este motivo, áreas como la salud por dispositivos móviles (mSalud) han empleado esta estrategia para ofrecer servicios de salud a través de estos dispositivos. Un ejemplo de esto es presentado en (Cabrera, y otros, 2014) donde se desarrolló un sistema de mensajería móvil que tiene como objetivo apoyar el control de la diabetes en México. Este sistema se encarga de recordar a los pacientes el horario para tomar los medicamentos y la asistencia a las citas médicas de control, además se envía información que ayuda a promover estilos de vida saludables, entregando tips para el cuidado de la salud.

## **1.2 Minería de Datos**

La minería de datos se puede definir como el proceso a través del cual se descubre conocimiento no trivial en forma de patrones, asociaciones, cambios, anomalías y estructuras significantes de grandes cantidades de datos almacenados en bases de datos, bodegas de datos u otros repositorios de información. Para realizar este proceso se suele utilizar técnicas de la inteligencia artificial y la estadística (Mena, 1999).

### **1.2.1 Tipos de Análisis**

En la minería de datos se pueden desarrollar dos tipos de análisis, a saber, predictivo y descriptivo. Dichos análisis permiten desarrollar diferentes tareas como la clasificación (Kumari & Sunila, 2011), la predicción (Shaikh, 2014), la segmentación (Lovrić, Milanović, & Stamenković, 2014) y la asociación (Amin, Takib, Raza, & Javed, 2014).

- **Análisis predictivo**

Algunas de las aplicaciones comúnmente desarrolladas con análisis predictivo son: predecir riesgos, predecir activación de nuevos clientes, predicción de ventas, entre otras (Riquelme, Ruiz, & Gilbert, 2006). Este tipo de análisis se caracteriza porque requiere un conjunto de entrenamiento, el cual está formado por un histórico de datos.

En el análisis predictivo se pueden desarrollar tareas de predicción discreta y predicción continua. La predicción discreta también recibe el nombre de clasificación, donde el conjunto de entrenamiento está conformado por los diferentes atributos y clases predefinidas para cada registro, como se muestra en la Tabla 1.

<b>Id</b>	<b>Atributo 1</b>	<b>Atributo 2</b>	<b>...</b>	<b>Atributo n</b>	<b>Clase</b>
1	10	Alto		35	Cliente Oro
2	35	Bajo		54	Cliente Plata
3	43	Medio		28	Cliente Bronce
4	26	Bajo		65	Cliente Bronce
5	87	alto		32	Cliente Oro

**Tabla 1: Conjunto de entrenamiento utilizado en el análisis predictivo discreto**

En la predicción continua, en el conjunto de entrenamiento los registros están conformados por atributos y una variable de predicción continua, como se muestra en la Tabla 2.

<b>Id</b>	<b>Atributo 1</b>	<b>Atributo 2</b>	<b>...</b>	<b>Atributo n</b>	<b>Predicción</b>
1	10	Alto		35	54.678
2	35	Bajo		54	100.500
3	43	Medio		28	27.000
4	26	Bajo		65	9.800
5	87	alto		32	23.600

**Tabla 2: Conjunto de entrenamiento utilizado en el análisis predictivo continuo**

- **Análisis descriptivo**

Algunas de las aplicaciones más comunes del análisis descriptivo son: análisis del perfil de los clientes, segmentación de tipos de clientes o productos, detección de anomalías, detección de reglas que condicionen la venta de productos, entre otras (Riquelme, Ruiz, & Gilbert, 2006). En este tipo de análisis se pueden desarrollar tareas de agrupación o clustering, y de asociación. El conjunto de datos requerido está conformado por los diferentes atributos que se desean analizar para encontrar similitudes o asociaciones entre los datos. Un ejemplo de un conjunto de datos para análisis descriptivo se presenta en la Tabla 3.

Id	Atributo 1	Atributo 2	...	Atributo n
1	10	Alto		35
2	35	Bajo		54
3	43	Medio		28
4	26	Bajo		65
5	87	Alto		32

Tabla 3: Conjunto de datos utilizados en el análisis descriptivo

### 1.2.2 Técnicas de Minería de Datos

Existen diversas técnicas de minería de datos (Han & Kamber, 2006), la elección de una de ellas depende básicamente de dos condiciones: el tipo de atributos y el objetivo de la minería (Riquelme, Ruiz, & Gilbert, 2006). De forma general, las técnicas se pueden agrupar en técnicas supervisadas y no supervisadas. Aunque existen gran cantidad de técnicas, a continuación se presenta una breve descripción de las técnicas más utilizadas en las aplicaciones de minería de datos.

- **Técnicas supervisadas**

Las técnicas supervisadas son aplicadas en el análisis predictivo. Algunas técnicas supervisadas son: Redes Neuronales, Árboles de Decisión, Máquinas de Soporte Vectorial, Métodos de Regresión, Método Bayesiano y Métodos basados en Ejemplos.

**Las Redes Neuronales** imitan el funcionamiento del cerebro humano para realizar tareas de aprendizaje. Tienen una arquitectura organizada en capas de neuronas, las cuales tienen pesos asignados a sus interconexiones. El aprendizaje de la red consiste en ajustar los pesos mediante una regla que indica cómo modificar los pesos en función de los datos de entrenamiento (Wiener, Jan, & Weigend, 1995).

**Los Árboles de Decisión** representan reglas en una estructura de árbol, en la cual los nodos internos son configurados con los atributos, las ramas representan los valores del atributo y las hojas del árbol identifican las clases. La clasificación se realiza descendiendo en el árbol hasta alcanzar una hoja, la cual indica la clase a la cual pertenece cada registro de la base de datos (Apté & Weiss, 1997). También existen árboles de predicción, que permiten analizar la salida en variables continuas.

**Las Máquinas de Soporte Vectorial** mapean los datos de entrada a un espacio de características de más alta dimensión, donde se puede construir un hiperplano que separe los datos que pertenecen a la clase, de los que no pertenecen a ella. El mapeo de los datos se realiza por medio de una función kernel (por ejemplo: lineal, polinomial, función de base radial, sigmoial, entre otras). Los datos más próximos al hiperplano de separación

son conocidos como muestras críticas o vectores soporte del modelo (Joachims, Hofmann, Yue, & Yu, 2009).

**La Regresión** también es utilizada como técnica supervisada en la minería de datos. La regresión lineal permite predecir la salida continua de una variable dependiente. Por su parte, la regresión logística es utilizada para predecir la clase a la que pertenece cada registro de la base de datos según variables predictoras independientes entre sí (Yang & Liu, 1999).

**Los métodos Bayesianos** se basan en el teorema de Bayes para pasar de la probabilidad a priori de un suceso P (suceso) a la probabilidad a posteriori P (suceso/observaciones). El aprendizaje en el clasificador bayesiano consiste en estimar las diferentes probabilidades en términos de sus frecuencias sobre los registros de la base de datos. La probabilidad de que un registro pertenezca a una clase está dada por el teorema de probabilidad condicional de Bayes (Wettig, Grünwald, Roos, Myllymäki, & Tirri, 2002).

**Los métodos basados en ejemplos** también son llamados clasificadores perezosos ya que no realizan ninguna labor en la etapa de entrenamiento, sólo almacenan los datos. El algoritmo de los K – vecinos más cercanos, KNN (K-Nearest Neighbor), es el más utilizado. Cuando se tienen nuevos datos para clasificar, el algoritmo busca los k registros más cercanos según funciones de distancia. Finalmente, el algoritmo asigna la clase a la que pertenece la mayoría de los registros vecinos (Yang & Liu, 1999).

- **Técnicas no supervisadas**

Las técnicas no supervisadas son aplicadas en el análisis descriptivo. Algunas técnicas no supervisadas son: Método Particional, Método Jerárquico, Método Probabilístico, Redes Neuronales y Reglas de Asociación.

**Los método particionales** dividen el conjunto de datos en un número predefinido de clusters (Xu & Wunsch, 2005) (Filippone, Camastra, Masulli, & Rovetta, 2008). K-means es el algoritmo más popular por su simplicidad y eficacia. El objetivo del algoritmo es encontrar k centroides, uno por cada cluster, de tal manera que los centroides sean lo más alejados posibles según funciones de distancia y los datos son asociados al centroide más cercano (Steinley, 2006).

**Los métodos jerárquicos** permiten encontrar estructuras de clustering de forma recursiva, utilizando dendogramas o árboles binarios. En el dendograma la raíz representa la población, los nodos intermedios simbolizan la proximidad entre los datos y las hojas representan los datos de la población (Xu & Wunsch, 2005) (Filippone, Camastra, Masulli, & Rovetta, 2008).

**Los métodos probabilísticos** asumen que los datos son generados de acuerdo a distribuciones de probabilidad (Xu & Wunsch, 2005) (François, Ancelet, & Guillot, 2006). Expectation – Maximization es el algoritmo probabilístico más comúnmente usado, el cual asigna una distribución de probabilidad a cada cluster y ajusta los parámetros con los datos.

**Las redes neuronales** también son utilizadas en la búsqueda de clusters. El algoritmo más utilizado es SOMs (Self Organizing Maps), el cual tiene una arquitectura de neuronas hexagonal o rectangular. Las neuronas están conectadas entre sí con una relación de vecindad y se usa la regla de aprendizaje de Kohonen para buscar la neurona más cercana a cada uno de los datos (Meireles, Almeida, & Godoy, 2003).

Las Reglas de Asociación se utilizan para descubrir relaciones frecuentes entre los datos (Slimani & Amor). Apriori es el algoritmo más ampliamente utilizado para detectar asociaciones, el cual se basa en el conocimiento previo de los datos en cada iteración.

### 1.2.3 Metodologías de Minería de Datos

Existen diversas metodologías que proporcionan una serie de pasos a seguir con el fin de realizar una implementación adecuada de la minería de datos. Según sondeos publicados en KDnuggets<sup>2</sup>, las metodologías más utilizadas son: CRISP-DM, SEMMA, KDD y Catalyst.

- **CRISP-DM** (Cross Industry Standard Process for Data Mining) es una metodología de libre distribución que fue concebida desde un enfoque práctico de acuerdo la experiencia de sus creadores: un consorcio de empresas europeas, incluyendo SPSS de IBM. Actualmente CRISP-DM es la guía de referencia más utilizada en el desarrollo de proyectos de minería de datos (Moine, 2013), (Torres, 2013), (Corrales, y otros, 2014). Está constituida por seis fases: entendimiento del negocio, entendimiento de los datos, preparación de los datos, modelado, evaluación y despliegue.
- **SEMMA** (Sample, Explore, Modify, Model and Assess) es la propuesta de SAS Analytics Solutions para desarrollar proyectos de minería de datos. La metodología establece cinco fases: muestreo, exploración, modificación, modelado y evaluación (Azevedo & Rojão, 2008). Se caracteriza por incluir una fase de muestreo estadístico que no se considera en otras metodologías.
- **KDD** (Knowledge Discovery in Database) se conoce como el descubrimiento de conocimiento en bases de datos como un proceso no trivial donde se identifican patrones válidos, novedosos, potencialmente útiles y en última instancia entendibles en los datos (Usama, Piatetsky-Shapiro, Padhraic, & Uthurusamy, 1996). Algunos autores consideran a

---

<sup>2</sup> KDnuggets: Data Mining Community Top Resource <<http://www.kdnuggets.com/>>

la minería de datos como una etapa en el de KDD (Riquelme, Ruiz, & Gilbert, 2006). Sin embargo, según las encuestas de KDnuggets, se está utilizando KDD como metodología para hacer minería de datos.

- **Catalyst** también es conocida como la metodología P3TQ (Product, Place, Price, Time, Quantity). Las relaciones ente estas variables buscan mantener el producto correcto, en el lugar adecuado, en el momento adecuado, en la cantidad correcta y con el precio correcto. Esta metodología plantea la formulación de dos modelos: el modelo de negocios y el modelo de minería de datos (Pyle, 2003).

#### 1.2.4 Plataformas y Herramientas para Minería de Datos

Las plataformas de minería de datos son herramientas que facilitan la aplicación de las técnicas de la minería de datos. Algunas plataformas son: WEKA, RapidMiner, R, SPSS Modeler y SAS Enterprise Miner.

- **WEKA** (*Waikato Environment for Knowledge Analysis*) es una herramienta utilizada con mucha frecuencia en proyectos relacionados con la minería de datos. Esta herramienta ha sido diseñada por un grupo de desarrolladores de la universidad de Waikato en Nueva Zelanda, y se distribuye bajo licencia GNU, es decir que es posible modificar el código fuente para adicionar nuevas funcionalidades (Hall, y otros, 2009). La herramienta WEKA permite realizar tareas de clasificación, regresión, clustering, asociación y visualización. Una de las características más atractivas es su capacidad de extensibilidad, es decir, que añadir nuevas funcionalidades es una tarea sencilla (Azoumana, 2013).
- **RapidMiner** es una herramienta de minería de datos desarrollada en el año 2001 por el departamento de inteligencia artificial de la Universidad de Dortmund. Entre sus principales ventajas se destaca que es multiplataforma, de código abierto y con licencia GPL. RapidMiner permite analizar y extraer datos a través de unos operadores, utilizando para ello un entorno gráfico, como resultado se obtienen patrones y visualizaciones que pueden facilitar la interpretación y generación de conocimiento (Hofmann & Ralf, 2013).
- **R** es un software desarrollado para realizar análisis de datos y presentar como resultado cálculos estadísticos y gráficas que permiten extraer información valiosa de los datos. Fue desarrollada por los científicos Robert Gentleman y Ross Ihaka del departamento de estadística de la Universidad de Auckland de Nueva Zelanda. Esta herramienta permite realizar las siguientes tareas: análisis de modelos lineales y no lineales, pruebas estadísticas clásicas, análisis de series temporales, clasificación y agrupación (Langohr, 2010).

- **SPSS** es un paquete estadístico que contiene una serie de herramientas que permiten realizar análisis de datos. Una de estas herramientas está diseñada para realizar tareas de la minería de datos, se trata de SPSS Modeler (Devi, Rao, Setty, & Rao, 2013), esta herramienta permite desarrollar modelos predictivos orientados a mejorar la toma de decisiones.
- **SAS Analytics** es un paquete comercializado por SAS Institute que permite el modelado predictivo y descriptivo en minería de datos. Esta herramienta se complementa con módulos de visualización, investigación de operaciones, estadística y procesos de calidad.

Se han realizado diversas comparaciones de las herramientas para hacer minería de datos. En (Pehlivanli, 2011) se comparan tres herramientas: RapidMiner, WEKA y R. Dentro de la comparación se presentan algunas estadísticas de las cuales se puede concluir que las herramientas WEKA y RapidMiner son las más descargadas desde Internet y la herramienta WEKA es una de las más populares entre los profesionales.

Una comparación más reciente es realizada en (Tapia, Ruiz, & Chirinos, 2014). Las características comparadas son: área de trabajo, capacidad para integrarse con otro software, tipo de licencia y capacidad para manejar extensa cantidad de registros. Como resultado de la comparación se resalta la herramienta Rapid Miner por permitir un área de trabajo gráfica, tener capacidad para integrarse con otro software, ser de licencia libre y poder manejar extensas cantidades de registros.

### 1.2.5 Validación de modelos de minería de datos

La validación en el contexto de la minería de datos es el proceso en el cual se evalúa el rendimiento de los modelos de minería, utilizando para ello datos reales. Un modelo de base de datos, puede ser evaluado utilizando diferentes criterios.

Para la evaluación de un modelo de clasificación, comúnmente se utiliza la bondad del clasificador, que hace referencia a la capacidad de predicción del modelo. En esta tarea se pueden utilizar los siguientes métodos (Sebastiani, 2005):

- **Matriz de confusión:** Este método permite visualizar en una tabla, la distribución de errores cometidos por un clasificador para cada una de las clases establecidas. Típicamente en las columnas de la matriz se ubican las clases reales y en las filas las clases predichas, de modo tal que el número de datos clasificados correctamente en su respectiva clase es calculado como suma de la diagonal principal de la matriz. Los elementos que se encuentren por fuera de la diagonal principal, corresponden a clasificaciones incorrectas.

- **PRECISION:** Este método mide el porcentaje de documentos que se han clasificado en una clase y realmente pertenecen a ella
- **RECALL:** Este método, relacionado con la cobertura, mide el porcentaje de los documentos que pertenecen a una clase y son correctamente clasificados dentro de ella.
- **Área ROC:** En este método se utilizan curvas para mostrar la habilidad del clasificador para separar las predicciones verdaderas de las falsas. Estas curvas miden la relación entre la tasa de verdaderos positivos y la tasa de falsos positivos.

En la evaluación de los modelos de Clustering, se emplean dos conceptos fundamentales, a saber, la cohesión, entendida como una medición de la cercanía de los elementos de un cluster entre sí, y la separación, entendida como la distancia entre dos cluster diferentes. A partir de estos dos conceptos se establecen criterios para validar los modelos, en la medida en que se minimice la cohesión y se maximice la separación. Algunos de los métodos utilizados para esta medición son (Petrovic, 2006): índice de Dunn, índice de la silueta y el índice de Davies and Bouldin. En los dos primeros un valor elevado del índice implica un buen proceso de clustering, mientras que en el segundo un índice bajo indica un proceso de clustering exitoso.

### **1.3 Procesamiento de Imágenes**

Un sistema de minería de imágenes debe incluir la participación y articulación de varias técnicas (Fernández, Miranda, Guerrero, & Piccoli, 2006), como: recuperación de las imágenes, extracción de características, aplicación de minería de datos y reconocimiento de patrones.

#### **1.3.1 Recuperación de las imágenes**

La recuperación de información de imágenes (Image Information Retrieval) es un área de mucho interés que ha tomado fuerza con la evolución del internet, solo hasta antes del año 2000 se encontraba en internet cerca de 180 millones de imágenes que ocupan alrededor de 3 TB (Goodrum, 2000).

Por lo general estos sistemas se basan en tres etapas (Shaban & Khalaf, 2013), (Martínez, 2013): representaciones textuales de las características de la imagen, rasgos visuales de las imágenes y una combinación de las dos anteriores. En la primera etapa se utiliza exclusivamente el código lingüístico (Text-Based Image Retrieval), en la segunda etapa se utiliza la recuperación de imagen basada en contenido (Content-Based Image Retrieval) y en la tercera etapa se utiliza la recuperación visual basada en la semántica (Semantics-Based Video Indexing and Retrieval).

El objetivo de estas técnicas es extraer de una base de datos de imágenes, todas las imágenes similares a una imagen base que se desea consultar. Estos sistemas deben contar con tres módulos fundamentales (Martínez, 2013), a saber, módulo de descripción, módulo de consultas y

módulo de búsquedas. El módulo de descripción tiene como objetivo representar de forma numérica las propiedades o cualidades de las imágenes, estas propiedades pueden ser intrínsecas, como lo son el color, forma, textura, o propiedades extrínsecas, que se dividen en propiedades de nivel medio y propiedades de nivel alto. Las propiedades de nivel medio se refieren a la detección automática de elementos como límites, contornos y objetos, y las propiedades de nivel alto se refieren a la información que se puede extraer de forma subjetiva de una imagen. Por otra parte, el módulo de consulta tiene como objetivo brindarle al usuario una alternativa para consultar imágenes, en este módulo el usuario dispone de varios medios para decirle al sistema el tipo de imágenes que quiere consultar y el sistema le presentan una serie de imágenes relacionadas con la consulta introducida. Finalmente, el módulo de búsqueda utiliza algoritmos de similaridad para extraer las imágenes más relevantes relacionadas con cada consulta.

En (Singh, Thoke, Verma, & Chandraka, 2011) se presenta un sistema de extracción de información de imágenes basado en color y forma. En este estudio se presenta un interés especial por la recuperación de imágenes utilizando consultas incompletas y distorsionadas, como resultado se obtuvo un porcentaje de éxito de 79.87% utilizando una base de datos de 1875 imágenes. Para mejorar el resultado se propone adicionar características de textura.

### 1.3.2 Extracción de características

Una imagen contiene una gran cantidad de información que no puede ser procesada directamente por un sistema de clasificación. Para superar este problema se utilizan técnicas de extracción de características que tienen por objetivo generar nuevas variables que representen la mayor cantidad de información de una imagen (Romero, 2011).

Por lo general, para la extracción se utilizan características de textura, color y forma (Maldonado, 2008).

- **Textura:** No se ha encontrado evidencia de una definición formal para las características de textura en una imagen. Sin embargo, cuando se menciona este concepto se hace énfasis en la repetición de un patrón especial en la imagen. Este tipo de características ha tomado un papel muy importante en el análisis de imágenes médicas (Maldonado, 2008) El análisis de textura se basa en cuatro tipos de métodos, a saber, estadísticos, estructurales, basados en modelos y basados en transformadas.

**Métodos estadísticos:** Es una forma de comparar la textura de dos imágenes a partir de sus píxeles simples. Una forma de utilizar estos métodos es mediante el análisis de las imágenes en escala de grises, a partir de dicho análisis se pueden extraer datos estadísticos como media, mediana y varianza.

**Métodos estructurales:** Estos métodos se basan en patrones primitivos y la forma como estos se agrupan entre sí para formar la textura.

**Métodos basados en modelos:** El objetivo de estos métodos es construir un modelo de la imagen con el fin de describir y sintetizar una textura. Algunos de estos métodos son: las cadenas de Markov y los fractales.

**Métodos basados en transformadas:** Estos métodos consisten en aplicar a las imágenes una serie de transformadas que permitan, entre otras cosas, identificar características presentes en las texturas.

- **Color:** La representación de una imagen se puede realizar en diferentes espacios o modelos de color (García, Terrones, Henarejos, & Martínez, 2014). Estos son representaciones numéricas mediante las cuales es posible describir cualquier color. A continuación se mencionan algunos de los modelos de color más utilizados:

**Modelo de Color RGB.** Este modelo está basado en un sistema de coordenadas cartesianas, donde la representación de cada color se realiza mediante la mezcla por adición de los tres colores de luz primarios, a saber, rojo, verde y azul (Red – Green - Blue). En otras palabras, el modelo RGB se presentan tres sub imágenes, una por cada color de luz primario, la suma de estas dan como resultado una imagen a color.

**Modelo de color HSI.** En este modelo se representa una imagen de acuerdo a tres características fundamentales: el tono, la saturación y la intensidad (Hue, Saturation, Intensity). La función del tono consiste en diferenciar un color de otro, la saturación representa el nivel de luz blanca que se encuentra en el color o la intensidad del tono dentro del color, y la intensidad es una medida de brillo de la luz.

**Modelo de color YIQ.** Este modelo fue creado en el momento de introducción de la televisión a color, este modelo permite convertir una señal de blanco y negro a color y viceversa. La componente Y entrega la imagen en blanco y negro, mientras que las componentes I,Q contienen la información del color.

**Modelo de color YCbCr.** En este modelo se representa una imagen a través de la luminancia y la prominencia. La componente Y hace referencia a la luminancia, y las componentes Cb y Cr contienen información la prominencia, donde Cb representa la diferencia entre la componente azul y un valor de referencia y Cr indica la diferencia entre la componente roja y un valor de referencia.

- **Forma:** La forma de un objeto presente en una imagen se identifica fácilmente en aquellos casos en los que el objeto en cuestión puede aislarse fácilmente en la imagen, sin

embargo, este no es el caso general de imágenes médicas, por lo que la clasificación por forma presenta cierto grado de dificultad. Las características de forma se pueden extraer mediante tres tipos de técnicas, a saber, basadas en el contorno del objeto, basadas en mapas de bordes y basadas en regiones.

**Técnicas basadas en el contorno:** Este tipo de técnicas utilizan funciones de una dimensión para calcular el contorno de la imagen y utilizarlo como una característica diferenciadora.

**Técnicas basadas en mapas de bordes:** Estas técnicas se utilizan frecuentemente en aquellas imágenes en las que es difícil obtener una identificación precisa de los objetos.

**Técnicas basadas en regiones:** En estas técnicas se deja de lado el contorno de la imagen, para ocuparse de la región específica donde se encuentra el objeto a analizar dentro de la imagen.

## 2. ESTADO DEL ARTE

A continuación se presenta el estado del arte en aplicaciones de la minería de datos en ciudades inteligentes, aplicaciones de minería de imágenes en el sector salud y aplicaciones de la minería de datos para analizar diagnósticos médicos.

### 2.1 Minería de datos para el desarrollo de ciudades inteligentes

Es común que para el desarrollo de las ciudades inteligentes se utilicen técnicas como la minería de datos, con el fin de obtener información significativa que permita mejorar un proceso en cuestión. Esta idea es sostenida en (UIT, 2014) donde se presenta un informe que analiza el papel de las Tecnologías de la Información y la comunicación (TICs), en el desarrollo de ciudades inteligentes sostenibles (Smart Sustainable Cities, SSC). En el informe se resaltan tres dimensiones globales de una SSC, a saber, el medio ambiente y la sostenibilidad, los niveles de servicio de la ciudad y la calidad de vida de los ciudadanos.

En este estudio al analizar el papel de las TICs en el desarrollo de una SSC, se tienen en cuenta aspectos relacionados con la predicción de datos, de hecho se afirma que una ciudad inteligente es una ciudad predictiva, donde las decisiones se toman basadas en la información. Y es justamente este campo en el que la minería de datos entra a jugar un papel importante, ya que puede apoyar en la predicción de eventos futuros que permitan aprovechar los datos históricos para identificar patrones de comportamiento y aplicarlos en diversas áreas como los negocios, prevención de desastres, predecir rutas de tráfico a evitar, entre otras.

Un concepto similar es presentado en (Villena, Villanueva, & Serrano, 2013) donde se realiza un modelado predictivo de la contaminación para una ciudad sostenible. En este estudio se utilizan una gran cantidad de datos que han sido recolectados mediante una serie de sensores distribuidos a lo largo de una ciudad, para realizar un modelado predictivo basado en minería de datos. Para la evaluación del modelo se proponen dos tipos de escenarios, el primero de ellos corresponde a la contaminación atmosférica y el segundo a la contaminación acústica. En este estudio se utiliza la metodología CRISP-DM para aplicar la minería de datos y la plataforma WEKA para aplicar la minería de datos.

### 2.2 Minería de imágenes en el área de la salud

Un gran porcentaje de los datos requeridos para el análisis de las enfermedades, se encuentran en imágenes o texto. Sin embargo, la minería de datos se aplica convencionalmente en datos estructurados, es decir información organizada en bases de datos. Los análisis de minería de datos

NO estructurados son un requerimiento nuevo y exigente que permitirá procesar información multimedia, dando lugar a nuevas áreas de interés: minería de texto y minería de imágenes. Estas áreas pretenden desarrollar análisis predictivos y descriptivos a información multimedia de la salud.

La minería de imágenes (Choubassi, Nefian, Kozintsev, Bouguet, & Wu, 2007) (Hsu, Mong, & Ji, 2002) se interesa por extraer patrones característicos a partir de un gran número de imágenes. Una de sus aplicaciones consiste en brindar soluciones a los problemas del sector salud.

El proceso de segmentación es vital en el tratamiento de imágenes médicas. En (Palomino & Garcia, 2011) se hace una revisión del algoritmo de segmentación Watershed, ampliamente utilizado en segmentación de imágenes médicas. Esta técnica permite detectar regiones en una imagen, se ha utilizado con éxito para el análisis de manchas.

Un estudio similar se presenta en (Lorca, Arzola, & Pereira, 2010) donde se analiza la segmentación de imágenes médicas digitales mediante técnicas de Clustering. En este estudio se utilizan las técnicas K-means y Fuzzy K-means para realizar la segmentación de los órganos del cuerpo humano, para ser utilizados en modelos anatómicos 3D. Se concluye que las técnicas de Clustering presentan ventajas frente a otras técnicas de segmentación de imágenes.

Para utilizar las técnicas de la minería de imágenes en el sector salud, es indispensable contar con bases de datos de imágenes, donde se reúnan una gran cantidad de casos de pacientes que padezcan la misma enfermedad, de este modo es posible extraer la información necesaria para realizar las tareas de la minería. En (Santamaría, 2014) se presenta una metodología para la creación de una base de datos de mamografías que tiene como objetivo apoyar la detección de micro calcificaciones. En la actualidad la mamografía es el principal método para la detección de cáncer de mamá, en imágenes de este tipo se pueden identificar micro calcificaciones representadas por pequeños puntos de alta intensidad consideradas como una manifestación temprana del cáncer de mama. En este estudio luego de construir la base de datos, se realiza la clasificación de las imágenes utilizando como técnica las máquinas de soporte vectorial, para esto fue necesario realizar etapas previas de segmentación de la zona mamaria, filtrado y extracción de características. Como resultado final del estudio se presenta una base de datos con 552 imágenes de mamografías y un método de clasificación basado en máquinas de soporte vectorial.

Un aspecto muy importante para realizar el proceso de minería de imágenes, consiste en identificar la técnica que proporcione el mayor porcentaje de exactitud. En (Thamilselvan & Sathiaselvan, 2015) se realiza un ejercicio comparativo de algoritmos de minería de datos aplicados a clasificación de imágenes. Entre los algoritmos a comparar se encuentran las máquinas de soporte vectorial (SVM), redes neuronales (RNA), K-means, entre otros. Como resultado de la evaluación, se obtiene que el algoritmo que presentó mayor precisión son las máquinas de soporte vectorial.

La minería de imágenes también ha ayudado en el proceso de identificación biométrica. En (Valencia, Cruz, Caicedo, & Chamorro, 2014) se presenta un mecanismo de identificación biométrica basado en extracción de características del iris. En este estudio se establecen cuatro etapas fundamentales para realizar un sistema de reconocimiento biométrico usando el iris, a saber, captura, procesamiento, extracción de características y reconocimiento de patrones. Se hace énfasis en la importancia de la limpieza de las imágenes para obtener un buen resultado independientemente de las técnicas utilizadas.

Otra rama de interés para la minería de imágenes consiste en la identificación de características o comportamientos específicos de una persona. En (Emam, 2014) se presenta un sistema que permite detección inteligente de ojos somnolientos. Este sistema fue pensado como una herramienta para ayudar a los conductores de vehículos a mantener su atención en la vía y no quedarse dormidos al volante. Para ello se capturaron cerca de dos mil imágenes de los ojos y se probaron varias técnicas de la minería de datos con el fin de encontrar un sistema que pueda funcionar en tiempo real y dar solución al problema planteado.

### **2.3 Minería de imágenes en los diagnósticos médicos**

La minería de imágenes ha presentado una alternativa muy importante para los diagnósticos médicos, se trata de la posibilidad de realizar diagnósticos asistidos por computadora, gracias a la aplicación de las técnicas de minería de imágenes. Se ha encontrado evidencia de este tipo de aplicaciones orientadas a enfermedades como alzheimer, cáncer, apendicitis, problemas digestivos, entre otras y utilizando imágenes obtenidas mediante técnicas como rayos X, ecografías, resonancia magnética, mamografías entre otras.

En (Ion & Udristoiu, 2015) se presenta un estudio que utiliza reglas semánticas para realizar minería sobre imágenes de endoscopias digestivas, obtenidas de un repositorio gratuito en Internet. Inicialmente se utilizó una base de datos de 200 imágenes previamente clasificadas donde se encontraban cinco diferentes tipos de diagnósticos de exámenes de endoscopias digestivas. Con estas imágenes se establecieron diez y seis reglas semánticas para reconocer estos diagnósticos basadas en características de color y textura. Posteriormente se utilizó una segunda base de datos que contiene un total de 450 imágenes de las cuales 60 corresponden a un diagnóstico específico de interés para el proyecto. Al efectuar el diagnóstico se obtuvo una efectividad del 96.8%.

En (Suganthira, Thamilselvan, Sathiaselvan, & Lakshmi Prabha, 2015) se presenta un estudio de tres diferentes algoritmos de clasificación: árboles de decisión, algoritmos genéticos y k-means, utilizados para el análisis de imágenes biomédicas. En este estudio se enfatiza en la importancia de los rayos X para el diagnóstico de pacientes, técnica que se viene utilizando hace más de cien años y hoy en día sigue siendo importante para la detección de anomalías en el cuerpo humano, analizando como dichas anomalías bloquean el haz de rayos X. El estudio se centra

particularmente en el cáncer de pulmón, considerada la principal causa de muerte por cáncer en Estados Unidos, enfermedad a la que anualmente se le invierten casi diez mil millones de dólares en tratamientos. Para el estudio se utilizaron imágenes pertenecientes a más de 30 mujeres con edades entre los 25 y 70 años. Como resultado se obtuvo una precisión de 86% utilizando árboles de decisión, 91% utilizando algoritmos genéticos y 93% utilizando k-means.

En (Peláez, 2014) se presenta un estudio que analiza los métodos para la clasificación automática de imágenes de resonancia magnética del cerebro. En este estudio se enfatiza en la importancia del procesamiento digital de imágenes para la eliminación de ruido intrínseco generado por las herramientas de captura de imágenes y se plantea dentro de los objetivos identificar cuál técnica de clasificación presenta mayor eficiencia. Como resultado, se obtuvo un mayor éxito utilizando las máquinas de soporte vectorial como técnica de clasificación.

En (Balu & Devi, 2012) se presenta un sistema de detección automática de apendicitis utilizando técnicas de minería de imágenes aplicadas a imágenes ecográficas. La apendicitis aguda es una enfermedad que requiere de atención inmediata en donde se diagnostica si se debe realizar una cirugía abdominal, para este diagnóstico se utilizan imágenes obtenidas de ecografías, que brindan algunas facilidades para el paciente en tanto que no son invasivas, se realizan rápidamente y tienen un bajo costo. Sin embargo, hoy en día el análisis de estas ecografías es una actividad que toma un tiempo considerable, razón por la cual un sistema de detección automático para detección de apendicitis aguda es de gran interés en tanto disminuya los tiempos del diagnóstico de forma significativa. En este trabajo se emplea la técnica de la distancia euclidiana para realizar el diagnóstico de apendicitis aguda utilizando ecografías de 44 pacientes con edades entre los 16 y los 52 años. La evaluación de la clasificación es realizada mediante la matriz de confusión, donde se destaca un 86% de sensibilidad, entendida como el porcentaje de pacientes que tienen la enfermedad y se diagnostica adecuadamente, y un 81% de especificidad, referida al porcentaje de pacientes que no tienen la enfermedad y el resultado del diagnóstico es negativo.

En (Romero, 2011) se presenta un estudio orientado al diagnóstico temprano de la enfermedad de Alzheimer utilizando imágenes tomográficas cerebrales. Se realiza inicialmente una tarea de segmentación que busca limitar la región de interés y posteriormente se realiza una etapa de extracción de características y calificación supervisada de imágenes, utilizando técnicas como las máquinas de soporte vectorial y las redes neuronales. El diagnóstico temprano de esta enfermedad es una tarea de gran interés ya que la enfermedad afecta aproximadamente a 30 millones de personas en todo el mundo y las técnicas de diagnóstico convencionales no superan el 70% de efectividad.

## **2.4 Discusión**

En los trabajos citados anteriormente se puede ver como el procesamiento digital de imágenes ha sido una herramienta de gran ayuda para resolver problemas en diversas áreas. En el caso

particular del sector salud hay una gran cantidad de enfermedades donde las imágenes son vitales para las tareas de diagnóstico y toma de decisiones.

En el estado del arte, se pueden resaltar los siguientes hallazgos:

- Se pueden identificar dos análisis de minería de datos comúnmente utilizados en las aplicaciones de minería de imágenes en la salud: clasificación y clustering. En el caso de la clasificación se observa cierta preferencia por técnicas como las máquinas de soporte vectorial, los árboles de decisión y las redes neuronales. En las técnicas mencionadas, la evaluación se realiza frecuentemente con la matriz de confusión y las medidas “*precision*” y “*recall*”. En el caso del Clustering se observa que K-means es una técnica de uso frecuente y ha presentado buenos resultados. En Clustering la evaluación se realiza, por lo general, con el índice de la silueta que es una medida utilizada en otras herramientas de software como es el caso del ToolBox de MatLab, R y SPSS modeler.
- Las metodologías de minería en su mayoría incluyen fases de selección de los datos, preparación, modelamiento y evaluación. Sin embargo, para analizar datos multimedia es necesario incluir etapas de pre procesamiento e indexamiento de los datos, donde la información multimedia es representada en estructuras de datos que son analizadas por las técnicas de minería. Como un acercamiento a esta etapa de pre procesamiento, algunos autores han modificado las metodologías CRISP-DM y KDD para realizar minería de texto (Santana, Costaguta, & Missio, 2014) (Tapia, Ruiz, & Chirinos, 2014). En el presente trabajo, se opta por modificar la metodología KDD para incluir dicha etapa, la nueva metodología se denomina KDM.
- En general para realizar minería de datos, se observa cierta preferencia por la plataforma WEKA debido, entre otros factores, a que permite utilizar una amplia variedad de métodos. Esta herramienta ha sido ampliamente utilizada en aplicaciones de la salud, sin embargo no cuenta con soporte para minería de imágenes, de hecho no se encontró ninguna plataforma con soporte a imágenes, las plataformas tradicionales solo soportan datos estructurados y algunas de ellas están empezando a soportar texto. Esta situación evidencia la necesidad de realizar una plataforma de minería que cuente con soporte para tratamiento de imágenes.

# Parte III: Solución Propuesta

## 1. MINERÍA DE IMÁGENES DEL SECTOR SALUD

En la revisión literaria presentada anteriormente se puede ver como el procesamiento digital de imágenes ha sido una herramienta de gran ayuda para resolver problemas en diversas áreas. En el caso particular del sector salud hay una gran cantidad de enfermedades donde las imágenes son vitales para las tareas de diagnóstico y toma de decisiones.

### 1.1 Requisitos de la Minería de Imágenes del Sector Salud

A continuación se describen los requisitos de minería de imágenes según los hallazgos del estado del arte.

- **Caracterización de las imágenes**

En el campo de la medicina se ha establecido un estándar con reconocimiento a nivel mundial para el almacenamiento y transmisión de pruebas médicas. Se trata de DICOM (Digital Imaging and Communication in Medicine). Este estándar (Bidgood, Horii, Prior, & Syckle, 1997) establece el tipo de fichero DICOM, que permite a entidades de salud compartir exámenes médicos facilitando su visualización, almacenamiento, transmisión e impresión. En este formato de archivos se distinguen dos componentes, un encabezado con campos informativos, y un contenido donde se encuentra la imagen en cuestión.

En la investigación realizada se encontró evidencia de la importancia del estándar DICOM para el manejo de imágenes médicas. Sin embargo aún se utilizan otros tipos de formato para el almacenamiento de imágenes médicas, es el caso de los formatos: IMG y .PGM. Este último permite almacenar una imagen en escala de grises, muy conveniente para algunos tipos de exámenes médicos.

Los formatos mencionados permiten almacenar imágenes de diferentes tamaños, por lo que las características a extraer de cada imagen no pueden depender del tamaño de la misma.

En la revisión literaria se encontró que las características más utilizadas son:

- **Características de color:** Permiten extraer aspectos relacionados con la distribución del color en la imagen. En este tipo de características se destacan descriptores como: **AuctoColorCorrelogram**, que entrega información sobre la correlación del color en la imagen

basado en el espacio de color HSV y Color Layout que consiste en dividir la imagen en varios bloques e identificar el color representativo de cada uno de los bloques.

- **Características de textura:** Permite identificar aspectos relacionados con el tipo de superficies presentes en una imagen. En este tipo de características se destaca el descriptor Tamura, que trabaja sobre seis diferentes tipos de textura relacionados con la percepción visual del ojo humano.
- **Características de forma:** Este tipo de características permite identificar la forma de los objetos contenidos en una imagen. Uno de los descriptores utilizados para ello es la pirámide de histograma de orientación de gradiente (PHOG)

- **Análisis y técnicas usadas en la minería de imágenes del sector salud**

Dentro de los estudios realizados sobre el procesamiento digital de imágenes de la salud, se pueden identificar dos necesidades básicas: la clasificación y el clustering. En el caso de la clasificación se observa cierta preferencia por las MÁQUINAS DE SOPORTE VECTORIAL y las REDES NEURONALES para el análisis de imágenes. Aunque en los estudios de minería de datos de la salud es bastante común el uso de ÁRBOLES DE DECISIÓN por ser un método de fácil entendimiento y ayuda visual para su interpretación. En cuanto al clustering la técnica utilizada tanto en los proyectos de minería de datos de la salud como en minería de imágenes es K-MEANS.

- **Evaluación de modelos de minería**

La evaluación de los modelos de clasificación se realiza con la matriz de confusión, a partir de la cual se calculan las medidas de PRECISION, RECALL y AREA ROC. Estas medidas se pueden obtener en todas las herramientas de minería de datos. En cuanto a la evaluación de los modelos de clustering, las herramientas de minería suelen utilizar medidas diferentes. Aunque recientemente se identifica una tendencia a utilizar el ÍNDICE DE SILUETA, el cual ya está incluido como medida de evaluación de clustering en Matlab, R y SPSS Modeler.

## **1.2 Metodología de Minería Multimedia KDM (Knowledge Discovery in Multimedia)**

Aunque las metodologías en su mayoría incluyen fases de preparación de los datos, dicha preparación sólo incluye análisis estadísticos y transformaciones. Para analizar datos multimedia es necesario incluir etapas de pre procesamiento e indexamiento de la multimedia, donde se pueda representar la información multimedia en vectores de características que puedan ser procesados por las técnicas de minería de datos. Como un acercamiento a esta etapa de pre procesamiento, algunos autores han modificado la metodología CRISP-DM para realizar minería multimedia. En (Santana, Costaguta, & Missio, 2014) se presenta una aplicación de algoritmos de

clasificación de minería de textos, a pesar de tratarse de datos no estructurados, se utiliza la metodología CRISP-DM. Un caso similar se presenta en (Tapia, Ruiz, & Chirinos, 2014) donde se utiliza la minería de texto en el diseño de un modelo de clasificación de opiniones subjetivas utilizando la metodología CRISP-DM.

Aunque en el estado del arte se observa una tendencia fuerte hacia la utilización de la metodología CRISP\_DM, recientemente se ha observado una disminución de su aplicación debido a que no se adapta a aplicaciones en *big data* y *data science*. Esto ha sido publicado en *Kdnuggets*, sitio web que publica constantemente noticias y encuestas sobre minería de datos, ciencia de datos y *big data*<sup>3</sup>.

Teniendo en cuenta que KDD propone una metodología de descubrimiento de conocimiento más estándar, en este trabajo se opta por adecuar la metodología planteada en KDD para analizar información multimedia.

KDD plantea el desarrollo de 5 etapas para datos estructurados: selección de los datos, pre procesamiento, transformación, modelamiento (minería) y evaluación/interpretación. Sin embargo la información multimedia no está estructurada, sino que debe ser preparada para ser analizada por un computador. Por este motivo se propone la metodología KDM (Knowledge Discovery in Multimedia), la cual se presenta en la Ilustración 1.

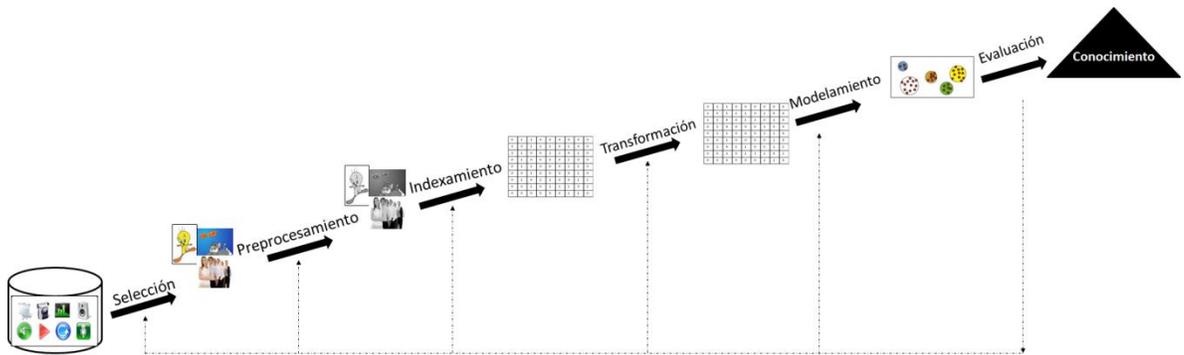


Ilustración 1: Metodología propuesta KDM

En la metodología se pueden presentar reprocesos en cualquiera de las etapas. A continuación se describen las etapas de la metodología KDM:

- Mediante el proceso de SELECCIÓN se extrae una muestra de la multimedia que será analizada, en este caso se seleccionan las imágenes a ser estudiadas.

<sup>3</sup> <http://www.kdnuggets.com/2014/10/new-poll-methodology-analytics-data-mining-data-science.html>

- La etapa de PREPROCESAMIENTO realiza procesos de limpieza, dependiendo del tipo de multimedia a analizar. En la Tabla 4 se presenta una descripción del pre procesamiento necesario para diferentes tipos de multimedia, aunque en este trabajo le compete sólo el pre procesamiento de imágenes.

<b>Tipo de multimedia</b>	<b>Pre procesamiento</b>
Textos	En el enfoque conocido como bolsa de palabras, se eliminan stopwords y se realiza reducción de raíces (stemming)
Imágenes	Se pueden realizar diferentes actividades dependiendo de las imágenes a analizar: <ul style="list-style-type: none"> <li>• Cambiar el formato de la imagen a png, jpg, etc</li> <li>• Segmentación para extraer el segmento de la imagen de nuestro interés y así eliminar ruido</li> <li>• Cambiar el espacio de color a binario, escala de grises, HSV, YUV, etc</li> <li>• Remuestreo de las imágenes al mismo tamaño</li> </ul>
Audios	Remuestreo de los audios a la misma resolución

**Tabla 4: Pre procesamiento de Multimedia**

- En la etapa de INDEXAMIENTO, se representa la multimedia en vectores numéricos que puedan ser analizados por las técnicas de minería de datos. Esta etapa también depende del tipo de recurso multimedia analizado como se presenta en la Tabla 5.

<b>Tipo de multimedia</b>	<b>Indexamiento</b>
Textos	En el enfoque conocido como bolsa de palabras, se calcula la frecuencia de las palabras.
Imágenes	Se pueden calcular diferentes características como: <ul style="list-style-type: none"> <li>• Características de color</li> <li>• Características de textura</li> <li>• Características de forma</li> </ul>
Audios	Se divide la señal de audio en ventanas y se extraen: <ul style="list-style-type: none"> <li>• Características en el dominio del tiempo</li> <li>• Características en el dominio de la frecuencia</li> </ul>

**Tabla 5: Indexamiento de multimedia**

- En la etapa de TRANSFORMACIÓN, es común realizar algunas modificaciones a las características como normalización, o conversión a categorías.

- En el MODELAMIENTO se aplican las técnicas de minería de datos.
- Finalmente en la etapa de EVALUACIÓN, se interpretan los resultados obtenidos para encontrar nuevo conocimiento.

## 2. DESARROLLO DE SOFTWARE

Para el desarrollo del prototipo de software de minería de imágenes se deben tener en cuenta asuntos como la metodología de desarrollo, el lenguaje de programación a utilizar, la base de datos, entre otros. Esta sección está dedicada a estos aspectos.

### 2.1 Metodología de desarrollo de software

Una metodología de desarrollo de software consiste en un marco de trabajo (framework) que establece condiciones para estructurar, planificar y controlar un proceso de desarrollo, brindando herramientas, modelos y métodos que se utilizan como apoyo a la hora de desarrollar.

En general, se puede identificar dos grandes paradigmas de desarrollo de software, el paradigma tradicional y el paradigma ágil.

El paradigma tradicional se caracteriza por tener mayor control en la programación del desarrollo, lo que permite reducir el riesgo de incurrir en gastos excesivos. Sin embargo este paradigma presenta una serie de desventajas en tanto que no se incluye la participación del usuario durante el proceso de desarrollo y si no se hace una aclaración adecuada de los requerimientos se pueden generar sobrecostos y exceder el tiempo planeado para el desarrollo. En este paradigma se suelen detectar posibles fallos en etapas tardías, lo que genera costos exhaustivos.

Por otra parte el paradigma de desarrollo ágil se caracteriza por incluir al usuario dentro del proceso de desarrollo lo que permite realizar una evaluación rápida y continua del desarrollo, que a su vez permite una detección temprana de los inconvenientes que se puedan presentar en el desarrollo. En este paradigma el desarrollo es iterativo, utilizando la retroalimentación del usuario en vez de un proceso extendido de planificación.

En la Tabla 6 se presenta una breve comparación entre las metodologías ágiles y tradicionales. Esta tabla ha sido tomada de (Calderón, Dámaris, Rebaza, & Carlos, 2007) donde se mencionan algunos aspectos a tener en cuenta para decidir qué tipo de metodología se debe utilizar dependiendo del tipo de proyecto.

<b>Metodologías Ágiles</b>	<b>Metodologías tradicionales</b>
Basadas en heurísticas provenientes de prácticas de producción de código.	Basadas en normas provenientes de estándares seguidos por el entorno de desarrollo.
Especialmente preparadas para cambios durante el proyecto.	Cierta resistencia a los cambios.
Impuestas internamente (por el equipo).	Impuestas externamente.
Proceso menos controlado, con pocos principios.	Proceso mucho más controlado, con numerosas políticas/normas.
No existe contrato tradicional o al menos es bastante flexible.	Existe un contrato prefijado.
El cliente es parte del equipo de desarrollo.	El cliente interactúa con el equipo de desarrollo mediante reuniones.
Grupos pequeños (<10 integrantes) y trabajando en el mismo sitio	Grupos grandes y posiblemente distribuidos.
Pocos artefactos.	Más artefactos.
Pocos roles.	Más roles.
Menos énfasis en la arquitectura de software.	La arquitectura del software es esencial y se expresa mediante modelos.

**Tabla 6: Diferencias entre metodologías ágiles y tradicionales**

De acuerdo a las características del proyecto enunciadas en los capítulos de levantamiento de requisitos (III) y análisis y diseño (IV) se opta por utilizar para el desarrollo del prototipo de software de minería de imágenes, una de las metodologías del paradigma ágil.

Las metodologías ágiles de desarrollo de software se basan en el manifiesto ágil, estas son un conjunto de postulados y principios a los que se deben ceñir las metodologías ágiles de desarrollo de software. Los postulados del manifiesto ágil son (Mendes C., Estevez, & Fillotrani, 2010):

- Valorar al individuo y a las interacciones del equipo de desarrollo por encima del proceso y las herramientas
- Valorar el desarrollo de software que funcione por sobre una documentación exhaustiva
- Valorar la colaboración con el cliente por sobre la negociación contractual
- Valorar la respuesta al cambio por sobre el seguimiento de un plan

Entre las diversas metodologías ágiles de desarrollo de software, se destaca la metodología Scrum por la popularidad y el acogimiento que ha tenido en los últimos años, en los que ha sido utilizada por empresas como Yahoo y Google (Amaya B., 2013). Para el desarrollo del prototipo de software de minería de imágenes se tendrá en cuenta algunos elementos de esta metodología.

## **2.2 Lenguaje de programación**

En el capítulo correspondiente al documento de análisis y diseño (IV), se planteó una solución orientada a objetos. Existe un gran número de lenguajes de programación orientados a objetos, se caracterizan principalmente porque permiten definir tipos de datos, realizar operaciones nuevas sobre dichos tipos de datos e instanciar los tipos de datos creados.

Para el desarrollo del prototipo de software de minería de imágenes se escoge trabajar sobre el lenguaje de programación Java. La filosofía de este lenguaje permite que los desarrolladores realicen los programas una sola vez y lo puedan ejecutar en cualquier dispositivo ("write once, run anywhere"), esto se logra gracias a la máquina virtual de Java. Esta característica permite cumplir con las especificaciones establecidas en el capítulo de levantamiento de requisitos (III).

## **2.3 Base de datos**

Una base de datos consiste en una serie de datos que tienen características en común y han sido almacenados en un mismo sitio para su posterior uso, la administración de las bases de datos se realiza mediante gestores de bases de datos que proporcionan el acceso a los datos.

Para el desarrollo del prototipo de software de minería de imágenes se optó por utilizar el gestor de base de datos Apache Derby, el cual puede ser embebido dentro de una aplicación Java, lo que implica la utilización de un menor espacio en disco. Esta característica es acorde con lo expresado en el capítulo de levantamiento de requisitos (III), permite que la aplicación sea altamente portable.

## **2.4 Librerías utilizadas**

Este proyecto requiere la utilización de técnicas elaboradas para el procesamiento de las imágenes y la aplicación de la minería. Para estas dos actividades se realizó una búsqueda con el fin de identificar librerías de terceros que faciliten el desarrollo.

### **Extracción de características**

En el proceso de extracción de características de una imagen se encontró que la librería LIRE (Lucene Image Retrieval) es utilizada frecuentemente en proyectos similares. Se trata de una librería de uso libre bajo licencia GNU-GPL para el lenguaje de programación Java que permite la recuperación de contenido a partir de una imagen (Lux & Chatzichristofis, 2008).

Mediante la librería LIRE es posible extraer de una imagen los siguientes tipos de descriptores:

- **Auto Color Correlogram:** Este descriptor utiliza el espacio de color HSV para entregar información sobre la correlación del color en una imagen.
- **Color Layout:** El objetivo de este descriptor es entregar información sobre la distribución espacial del color en una imagen. Es ampliamente utilizado gracias a su velocidad de ejecución.
- **Scalable Color:** Se trata de un análisis del histograma de color en el espacio HSV utilizando la transformada Haar.
- **Simple Color Histogram:** Este descriptor realiza un histograma del color de una imagen en el espacio de color RGB.
- **Fuzzy Color Histogram:** Es un histograma de color que se caracteriza por ser poco sensible a la interferencia de ruidos generados por factores como la intensidad de la iluminación.
- **Edge Histogram:** Presenta la distribución espacial de cinco tipos de bordes. Para ello divide la imagen en 16 sub imágenes y para cada una de ellas genera el histograma.
- **Tamura:** En esta técnica se propone la utilización de seis diferentes tipos de textura que puede percibir el ser humano mediante las cuales busca identificar el tipo de superficies presentes en una imagen.
- **PHOG:** El descriptor Pirámide de Histograma de Orientación de Gradiente, está basado en la apariencia global de una imagen representando la forma y distribución espacial en la imagen.
- **Local Binary Patterns:** Este descriptor entrega información sobre la textura de una imagen a partir de la conversión de la misma a escala de grises.
- **Rotation Invariant Local Binary Patterns:** Es una extensión del descriptor Local Binary Pattern en donde se aplica una rotación a la imagen.

Con la extracción de estos descriptores es posible representar las características de color, forma y textura de un conjunto de imágenes, y así generar vectores característicos para la representación de las mismas y aplicar sobre estos vectores, técnicas de la minería de datos.

### **Aplicación de la Minería de datos**

Para la aplicación de la minería de datos se identificó, en la investigación realizada, una alta preferencia por la herramienta WEKA (Waikato Environment for Knowledge Analysis). Se trata de una herramienta para el aprendizaje automático y la minería de datos desarrollada por la Universidad de Waikato de Nueva Zelanda, escrita en lenguaje Java y distribuida bajo la licencia GNU-GPL. Mediante Weka se pueden realizar las tareas de clasificación y clustering, utilizando para ello diversas técnicas de la minería de datos como los árboles de decisión, las máquinas de soporte vectorial y K-Means.

Weka se ajusta perfectamente a las necesidades de este desarrollo, por tanto se decide incorporar sus librerías para la aplicación de las técnicas de minería de datos.

**Parte IV: Documento de  
Requisitos de Software  
IEEE830**

## 1. PRESENTACIÓN

Las ciudades inteligentes se han convertido en una temática de interés común en los últimos años, donde se plantea el objetivo de ayudar a mejorar la calidad de vida de los ciudadanos, teniendo como enfoque áreas como: gobierno, construcción, movilidad, energía, medio ambiente y servicios. Los servicios de salud se enmarcan dentro de las áreas de gobierno y servicios, presentando como objetivo principal prevenir enfermedades y mejorar la salud de los ciudadanos.

Un gran porcentaje de los datos requeridos para el análisis de enfermedades se encuentra representado en texto e imágenes. En el caso de la información representada en texto, hoy en día existen algunas plataformas que permiten realizar labores de minería con el fin de analizar ese tipo de información, sin embargo, en el caso de las imágenes, no se ha encontrado reporte de plataformas similares. En este sentido existe una oportunidad de desarrollar un prototipo de plataforma de minería de imágenes para el sector salud, donde se busca extraer información no trivial a partir de un conjunto de imágenes y de esta manera realizar tareas de diagnóstico y prevención de enfermedades

### 1.1 Propósito

El objetivo de esta sección consiste en realizar el levantamiento de requisitos para el desarrollo de un prototipo de software que permita realizar minería, sobre imágenes provenientes del sector salud como apoyo a los servicios de salud en el marco de las ciudades inteligentes.

El levantamiento de requisitos servirá como apoyo durante el desarrollo y la evaluación del proyecto, en esta última etapa se realizará un caso de estudio con imágenes médicas donde se evaluará el funcionamiento del prototipo de software a desarrollar.

### 1.2 Alcance

El alcance del proyecto consiste en el desarrollo de un prototipo de software que permita realizar minería de imágenes, incluyendo análisis predictivo y descriptivo, se contempla una posterior validación sobre un conjunto de imágenes médicas.

Durante el desarrollo del prototipo se contempla la evaluación e implementación de técnicas de minería de imágenes que han presentado un buen desempeño en trabajos similares, teniendo en cuenta específicamente la realización de análisis predictivo (clasificación) y análisis descriptivo (clustering).

Cabe anotar que el objetivo del proyecto no es hacer un estudio profundo de minería de imágenes, se trata de un proyecto de desarrollo de software donde se aprovechan las técnicas existentes para realizar minería sobre imágenes del campo de la salud.

### 1.3 Personal involucrado

Nombre	Efraín Alberto Oviedo
Rol	Diseño, desarrollo, validación
Categoría profesional	Estudiante de maestría en Tics
Responsabilidades	Realizar las tareas de diseño, desarrollo y validación del prototipo de software
Información de contacto	eaoc46@gmail.com

Nombre	Ana Isabel Oviedo
Rol	Experta en minería de imágenes, Directora del proyecto de maestría.
Categoría profesional	Doctora en Ingeniería Electrónica
Responsabilidades	Apoyo y revisión de las labores de diseño, desarrollo y validación del prototipo de software
Información de contacto	ana.oviedo@upb.edu.co

Nombre	Gloria Liliana Vélez
Rol	Experta en gestión de proyectos
Categoría profesional	Doctora en Ingeniería Electrónica
Responsabilidades	Revisión de resultados del proyecto
Información de contacto	gloria.velez@upb.edu.co

Tabla 7: Personal involucrado en el proyecto

## 1.4 Resumen

El prototipo de software a realizar debe permitirle al usuario aplicar minería sobre un conjunto de imágenes. Para ello inicialmente se le solicita al usuario indicar la ruta donde se encuentran almacenadas las imágenes a analizar, luego se debe indicar el tipo de características que se desea extraer de dichas imágenes, posteriormente se debe indicar el tipo de análisis y la técnica particular a utilizar, y finalmente se presentan los resultados del análisis.

En esta primera sección se realizó una presentación general del proyecto al cual se le realiza la especificación de requisitos, específicamente se presentó información relacionada con el propósito del documento, el alcance del producto y las personas involucradas en el desarrollo del mismo.

En la segunda sección de este capítulo se presenta una descripción general del producto a desarrollar, donde se podrá encontrar información relacionada con la perspectiva que se tiene del producto, su funcionalidad, el tipo de usuarios a quienes va dirigido, las restricciones que presenta, suposiciones y dependencias e información relativa a posibles desarrollos futuros.

Posteriormente en la tercera sección se detallarán los requisitos específicos del producto, clasificados en las siguientes categorías: requisitos comunes de las interfaces, requisitos funcionales y requisitos no funcionales.

## 2. DESCRIPCIÓN GENERAL

### 2.1 Perspectiva del producto

El producto a desarrollar hace parte de un proyecto del grupo de Investigación GIDATI de la Universidad Pontificia Bolivariana, que lleva por título “Plataforma de Minería de Datos Estructurados y No Estructurados – Caso de Estudio Salud Pública” y tiene como objetivo apoyar los servicios de salud en las ciudades inteligentes mediante el desarrollo de una plataforma de minería de datos estructurados y no estructurados (específicamente texto e imágenes).

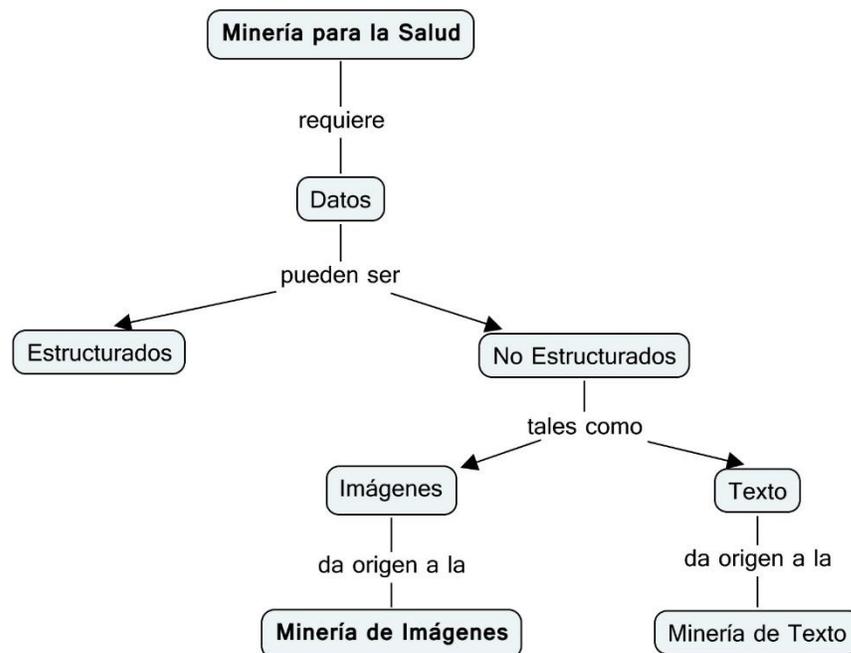


Ilustración 2: Mapa conceptual: minería de datos para la salud

En la Ilustración 2 se presenta un mapa conceptual de la minería de datos para la salud, se destaca que los datos pueden ser de tipo estructurados o no estructurados y estos últimos pueden referirse a texto o imágenes. El producto a desarrollar involucra exclusivamente el tratamiento de imágenes conocido como minería de imágenes.

## 2.2 Funcionalidad del producto

El prototipo de software a desarrollar debe contar con un sistema de control de acceso que valide el ingreso del personal autorizado. Se dispone inicialmente de dos tipos de usuarios, los administradores que tienen la facilidad de gestionar usuarios y modificar los parámetros de las técnicas de minería de imágenes, y los operadores que solamente pueden acceder a las funciones que aplican dichas técnicas sobre un conjunto de imágenes.

El sistema, debe permitirle al usuario cargar un conjunto de imágenes desde un directorio en el equipo local, posteriormente el usuario debe seleccionar el tipo de características que desea extraer de las imágenes y finalmente se le ofrecen opciones para realizar análisis predictivo o descriptivo. En ambos casos el usuario debe seleccionar la técnica que desea utilizar para la minería, estas técnicas se mencionan en la sección de requisitos funcionales y fueron escogidas en base a un estudio previo sobre técnicas comúnmente utilizadas para minería de imágenes en la salud.

## 2.3 Características de los usuarios

Se contempla que el prototipo de software a desarrollar será utilizado por los siguientes tipos de usuario.

Tipo de usuario	Profesionales
Formación	Ingenieros
Habilidades	Profesionales de Ingeniería con conocimientos en inteligencia artificial y tecnologías de la información y la comunicación
Actividades	Validación del proceso de minería de datos, ajuste de parámetros de las técnicas de minería.

Tipo de usuario	Profesionales
Formación	Doctores en área de la salud
Habilidades	Personal médico con amplios conocimientos en áreas de la

	salud
Actividades	Agrupación y clasificación de imágenes, verificación y análisis de los resultados

**Tabla 8: Características de los usuarios de prototipo de software**

## 2.4 Restricciones

Debido a que no se tiene información sobre el sistema operativo de los equipos de cómputo utilizados por los usuarios de la aplicación, es importante que esta permita la ejecución en los sistemas operativos más comunes a la fecha (Windows y Linux).

Adicionalmente se debe tener en cuenta que por lo general el personal que utiliza la aplicación no posee amplios conocimientos en sistemas, por tanto la instalación y ejecución del prototipo de software deben ser procedimientos sencillos e intuitivos de modo que dicho personal pueda acceder a la aplicación sin mayor complejidad. Así mismo es importante mencionar que la aplicación a desarrollar se ejecutará en paralelo con las demás aplicaciones usadas hoy en día por los usuarios, es decir que no se tendrán equipos dedicados para su ejecución, por tanto se recomienda evitar componentes de software que requieran un uso excesivo de la memoria de los equipos de cómputo.

En cuanto al desarrollo del prototipo de software cabe mencionar que se trata de un trabajo de grado que tiene una duración promedio de seis meses, lo que se convierte en una limitante de tiempo a considerar.

## 2.5 Suposiciones y dependencias

Se espera que el conjunto de imágenes suministrado para la validación del prototipo de software se encuentre en buen estado de conservación y sea suficiente para el desarrollo del proyecto. La calidad de dichas imágenes puede llegar a tener una influencia significativa sobre el resultado del proceso de minería de imágenes.

## 2.6 Evolución previsible del sistema

En el desarrollo de aplicaciones para las ciudades inteligentes se vienen presentando una serie de tendencias en cuanto a tecnologías o conceptos, entre ellas cabe destacar el internet de las cosas, la computación en la nube y las aplicaciones móviles.

Dentro del desarrollo del proyecto, el internet de las cosas (IoT) puede apoyar las labores de adquisición de las imágenes médicas, la computación en la nube puede aportar en la capa de Software como servicio (SaaS) e Infraestructura como servicio (IaaS) y las aplicaciones móviles pueden facilitar el acceso a la aplicación.

Tener en cuenta estas tendencias dentro del desarrollo del prototipo de software implicaría esfuerzos importantes en cuanto a costo y dedicación, por tanto, teniendo en cuenta las limitantes de tiempo expuestas en la sección 2.4 de este documento, se propone posponer la incorporación de estas tendencias como un trabajo a futuro. Esto significa que dichos conceptos y sus beneficios no se tendrán en cuenta dentro del desarrollo de prototipo de software.

### 3. REQUISITOS ESPECÍFICOS

#### 3.1 Requisitos comunes de los interfaces

Para el diseño de este prototipo de software se han considerado dos tipos de perfiles, a saber perfil de operador y perfil de administrador. Ambos perfiles, deben de autenticarse ante el sistema mediante un mecanismo de usuario y contraseña. Una vez se ha iniciado sesión, el sistema debe presentar en forma de menú las opciones disponibles. La aplicación de minería de datos es común a ambos perfiles de usuario, mientras que las opciones de gestión de usuarios y configuración de los métodos empleados son exclusivas de usuarios con el perfil de administrador.

La aplicación de la minería de imágenes, requiere que el usuario realice las siguientes tareas:

- Ingresar la ruta donde se encuentran almacenadas las imágenes que desea analizar
- Seleccionar el tipo de características que desea extraer
- Seleccionar el tipo de análisis a realizar
- Seleccionar la técnica que desea utilizar

Estas son las entradas principales del sistema, con esta información el sistema realiza la minería de imágenes y presenta al usuario los resultados, siendo estos la salida del sistema.

#### 3.2 Interfaces de usuario

De acuerdo a lo comentado en la sección 2.1 de este capítulo, el prototipo de software a desarrollar hace parte de un proyecto del grupo de investigación GIDATI de la universidad Pontificia Bolivariana, por esta razón se deben utilizar los logos indicados por el grupo de investigación o en su defecto el logo de la universidad.

Para el prototipo de software se requiere del diseño de las siguientes pantallas:

- **Inicio de sesión:** Es la pantalla de acceso al sistema, en esta se debe permitir el ingreso de usuario y contraseña, y se debe contar con un botón para iniciar sesión. (RF1)
- **Pantalla principal:** En esta pantalla se permite el acceso a las funcionalidades de la aplicación dependiendo del perfil con el que se ha iniciado sesión.

La pantalla debe contar con una serie de botones que permitan el acceso a cada una de las funcionalidades.

- **Gestión de usuarios:** En esta pantalla los administradores pueden agregar, modificar o eliminar usuarios de la aplicación.

La pantalla debe listar en forma de tabla todos los usuarios de la aplicación y al lado de cada usuario se debe agregar botones para eliminar y modificar. Adicionalmente se dispone de un botón que permite adicionar un nuevo usuario.

Las opciones de modificar y agregar deben presentar una pantalla adicional donde se permite el ingreso o modificación de los datos y se debe contar con botones para aceptar o cancelar. (RF2)

- **Modificar información personal:** En esta pantalla tanto operadores como administradores pueden actualizar sus datos personales como: nombre, usuario y contraseña.

La pantalla debe contar con campos de texto editable para el ingreso de esta información y con botones para aceptar o cancelar la operación. (RF3)

- **Cargar imágenes:** En esta pantalla el usuario le entrega al sistema el conjunto de imágenes que desea utilizar para realizar la minería. (RF5)

Se debe contar con un botón de cargar mediante el cual el sistema le presenta al usuario una opción para indicar la ruta donde están almacenadas las imágenes.

El sistema deberá mostrar datos relacionados con el resultado de dicha operación, tales como número de imágenes cargadas.

- **Extracción de características:** En esta pantalla el usuario podrá escoger el tipo de características que desea extraer de las imágenes. Una vez seleccionado el tipo de características, se debe contar con un botón mediante el cual el sistema proceda a realizar el proceso de extracción de cada una de las imágenes. A esta funcionalidad se le denomina el procesamiento de las imágenes.(RF6)

Para realizar esta labor el usuario debe haber cargado las imágenes previamente.

- **Visualización de datos:** En esta pantalla el sistema le permite al usuario: (1) visualizar las imágenes cargadas, (2) visualizar las características extraídas, (3) ver algunos estadísticos de las características como media, desviación estándar, máximo y mínimo. (RF7)

Para ingresar a esta opción el usuario debe haber realizado la extracción de características previamente.

- **Análisis predictivo:** Es la pantalla en la cual el usuario realiza un análisis predictivo, de las imágenes cargadas en la pantalla anterior, en caso de que el usuario aún no tenga

imágenes cargadas o no realice previamente el procesamiento de las mismas, no se debe permitir el acceso a esta funcionalidad. (RF8)

Adicionalmente se deben contemplar labores de entrenamiento y prueba.

- **Análisis descriptivo:** Es la pantalla en la cual el usuario realiza un análisis descriptivo de las imágenes cargadas, en caso de que el usuario aún no tenga imágenes cargadas, o no realice previamente el procesamiento de las imágenes, no se debe permitir el acceso a esta funcionalidad. (RF9)
- **Configuración de las técnicas de minería:** Esta pantalla es para el uso exclusivo de los usuarios con perfil de administrador, en ella se permite ajustar parámetros de ejecución de las técnicas implementadas para la minería de imágenes.  
La pantalla listará todos aquellos parámetros que puedan ser modificados por el usuario, para que este los configure de acuerdo a su necesidad, se debe contar con un botón de aceptar mediante el cual el sistema almacena estos parámetros para la operación. (RF4)

### 3.3 Interfaces de hardware

En esta etapa del proyecto se está desarrollando un prototipo de software que aún no contempla la interacción con algún hardware específico. Por tal motivo no se consideran dentro del alcance del proyecto requisitos de interfaces de hardware.

Sin embargo, de acuerdo a lo comentado en la sección 2.6 de este capítulo, en el contexto de las ciudades inteligentes se están presentando tendencias relacionadas con áreas como el Internet de las cosas. La inclusión de este concepto puede implicar en un desarrollo a futuro que involucre requerimientos en las interfaces de hardware.

### 3.4 Interfaces de software

Para este prototipo de software no se contempla la interacción con otros sistemas de software. Por tal motivo no se consideran dentro del alcance del prototipo requisitos de interfaces de software.

En posteriores etapas del proyecto es posible que esta integración sea necesaria con el fin de compartir bases de datos con software similares.

### 3.5 Interfaces de comunicación

En el desarrollo del prototipo de software no se está contemplando la posibilidad de establecer comunicación con sistemas similares, por tanto no se consideran dentro del alcance requisitos de interfaces de comunicación.

Así como se comentó en la sección 3.4 de este capítulo, un trabajo a futuro en el que se contemple integrar con otros sistemas de software, puede implicar requisitos de comunicación.

### 3.6 Requisitos funcionales

#### Requisito funcional 1

Número de requisito	RF1
Nombre de requisito	Sistema de control de acceso
Tipo	<input type="checkbox"/> Requisito
Prioridad del requisito	<input type="checkbox"/> Media/Deseado

Tabla 9: Requisito funcional 1

Cuando un usuario desee acceder al sistema se le solicitará la siguiente información:

- Usuario
- Contraseña

Ambos campos serán alfanuméricos de máximo 16 caracteres cada uno.

En caso de que la autenticación sea exitosa se procede a mostrar la pantalla principal que brinda acceso a las diferentes funciones del prototipo de software, en caso contrario se sigue presentando esta pantalla de inicio de sesión y se le indica al usuario que la autenticación no se ha realizado correctamente, en este caso el usuario debe revisar los datos ingresados.

## Requisito funcional 2

Número de requisito	RF2
Nombre de requisito	Gestión de usuarios
Tipo	<input type="checkbox"/> Requisito
Prioridad del requisito	<input type="checkbox"/> Media/Deseado

Tabla 10: Requisito funcional 2

Solo los usuarios con perfil de administrador pueden acceder a esta funcionalidad.

En esta opción se pueden realizar las siguientes acciones:

- **Crear usuario:** En esta opción se crea un nuevo usuario que tendrá acceso a la aplicación, para la creación de un usuario se solicitan los siguientes datos:
  - Perfil (Selección entre operador y administrador)
  - Nombre (Alfanumérico)
  - Usuario (Alfanumérico)
  - Contraseña (Alfanumérico)

La aplicación debe validar los datos ingresados, si se presenta algún tipo de anomalía se le informa al usuario y se le pide que corrija los datos que están generando conflicto.

- **Eliminar usuario:** Accediendo a esta opción se elimina un usuario existente de la aplicación. Una vez eliminado un usuario no será posible recuperarlo nuevamente.

El sistema muestra una lista con todos los usuarios actuales de la aplicación, para eliminar uno de ellos se debe hacer clic sobre el botón eliminar situado al frente del usuario que se desea borrar del sistema.

- **Modificar usuario:** Esta opción permite modificar la información de un usuario previamente creado en la aplicación, todos los campos del usuario pueden ser modificados.

El sistema muestra una lista con todos los usuarios actuales de la aplicación, para modificar uno de ellos se debe hacer clic sobre el botón modificar situado al frente del usuario que se desea modificar.

### Requisito funcional 3

Número de requisito	RF3
Nombre de requisito	Modificar Información personal
Tipo	<input type="checkbox"/> Requisito
Prioridad del requisito	<input type="checkbox"/> Media/Deseado

Tabla 11: Requisito funcional 3

Cada usuario, independientemente de su perfil, puede proceder a modificar su información personal. A continuación se presentan los campos que pueden ser modificados:

- Nombre (Alfanumérico)
- Usuario (Alfanumérico)
- Contraseña (Alfanumérico)

Una vez realizadas las modificaciones se debe dar clic en el botón aceptar para que el sistema revise los nuevos valores de cada campo y almacene la información en la base de datos.

En caso de presentarse alguna inconsistencia, el sistema indicará mediante un mensaje de error.

### Requisito funcional 4

Número de requisito	RF4
Nombre de requisito	Parámetros de funcionamiento
Tipo	<input type="checkbox"/> Requisito
Prioridad del requisito	<input type="checkbox"/> Media/Deseado

Tabla 12: Requisito funcional 4

Esta es una funcionalidad que solo está disponibles para usuarios con perfil de administrador. El objetivo es que se puedan ajustar parámetros de operación de las técnicas de minería de imágenes utilizadas, de modo que el sistema las almacene en memoria y cada vez que se inicie, utilice dichos parámetros para el funcionamiento.

En este capítulo no se detalla cuales son dichos parámetros, estos surgirán en la medida en que se realice la aplicación a criterio del desarrollador.

El sistema mostrará cada uno de esos parámetros y permitirá su edición, en el momento en que se tengan todos los parámetros configurados, se debe presionar el botón aceptar con el fin de que los parámetros sean almacenados en memoria, esta acción de inmediato elimina los parámetros anteriores.

#### Requisito funcional 5

Número de requisito	RF5
Nombre de requisito	Cargar imágenes
Tipo	<input type="checkbox"/> Requisito
Prioridad del requisito	<input type="checkbox"/> Alta/Esencial

Tabla 13: Requisito funcional 5

Esta es la primera acción que se debe realizar para aplicar la minería de imágenes, en ella el usuario le indica al sistema cuál es la ruta donde se han almacenado las imágenes que se pretende analizar.

El sistema cuenta con un botón de examinar que al ser presionado le permite al usuario indicar dicha ruta, una vez cargadas las imágenes el sistema debe indicar el número de imágenes que se encontraron en dicha ruta.

En caso de que en la ruta indicada no se encuentre ninguna imagen, el sistema debe indicárselo al usuario, en este caso no se habilita aún el acceso a las demás funcionalidades del sistema.

#### Requisito funcional 6

Número de requisito	RF6
Nombre de requisito	Procesamiento de las imágenes
Tipo	<input type="checkbox"/> Requisito
Prioridad del requisito	<input type="checkbox"/> Alta/Esencial

Tabla 14: Requisito funcional 6

Una vez cargadas as imágenes y antes de realizar los análisis permitidos (predictivo y descriptivo), se debe realizar un procedimiento de preparación de las imágenes el cual consiste en extraer de cada una de ellas las características que se desean analizar. Se deberá contar con mecanismos para extraer las siguientes características:

- Color
- Textura
- Forma

Se le debe brindar al usuario un mecanismo para indicar el tipo de características que se desea analizar y posteriormente se debe proceder a extraer dichas características de cada imagen. Luego de realizar este proceso se deben habilitar las opciones para realizar los análisis permitidos por la aplicación.

#### Requisito funcional 7

Número de requisito	RF7
Nombre de requisito	Visualización de datos
Tipo	<input type="checkbox"/> Requisito
Prioridad del requisito	<input type="checkbox"/> Alta/Eencial

Tabla 15: Requisito funcional 7

El sistema debe contar con una opción donde se puedan visualizar lo siguiente:

- Imágenes cargadas
- Características extraídas

#### Requisito funcional 8

Número de requisito	RF8
Nombre de requisito	Análisis predictivo
Tipo	<input type="checkbox"/> Requisito
Prioridad del requisito	<input type="checkbox"/> Alta/Eencial

Tabla 16: Requisito funcional 8

Esta funcionalidad permite realizar un análisis predictivo a partir de las imágenes cargadas, solo se habilitará si se han cargado previamente las imágenes, y se ha realizado la respectiva preparación de las mismas, en caso contrario la funcionalidad permanece deshabilitada.

El análisis predictivo se basa en algoritmos supervisados donde se dispone de un conjunto de datos históricos (datos de entrenamiento) y un conjunto de datos de prueba, utilizados para validar el resultado de la minería.

Este análisis puede realizar dos tipos de predicciones, a saber, predicción de categorías y predicción de números. En el área de la salud lo más común es realizar predicción de categorías, usualmente conocido como clasificación. Para este prototipo se debe considerar exclusivamente la predicción de categorías, no se tendrán en cuenta labores de predicción de números.

Se espera que el prototipo de software a partir del conjunto de imágenes de entrenamiento realice el entrenamiento del clasificador y posteriormente permita realizar la clasificación de una nueva imagen. Durante la etapa de entrenamiento se debe indicar la clase a la que pertenece cada una de las imágenes cargadas, mientras que en la etapa de prueba, simplemente se carga una nueva imagen y el sistema se encarga de realizar la respectiva clasificación

Para realizar este análisis, el usuario debe seleccionar una de las siguientes técnicas:

- Árboles de decisión: Son muy utilizados en el campo de la salud por su simplicidad en la interpretación
- Máquinas de soporte vectorial: Cuentan con soporte para una alta dimensión de datos

### Requisito funcional 9

Número de requisito	RF9
Nombre de requisito	Análisis descriptivo
Tipo	<input type="checkbox"/> Requisito
Prioridad del requisito	<input type="checkbox"/> Alta/Esencial

Tabla 17: Requisito funcional 9

Al igual que en el requisito funcional 8, esta funcionalidad solo se activará luego de cargar las imágenes y realizar la respectiva preparación de las mismas.

El análisis descriptivo consiste en la utilización de algoritmos no supervisados para descubrir patrones en los datos, es decir, se analizan las imágenes cargadas para descubrir patrones o tendencias entre ellas.

Este análisis comprende las labores de clustering y de asociación. Esta última solo se puede realizar con datos categóricos pero los datos que se extraen de las imágenes son numéricos, por este motivo solo se contemplará dentro del alcance del prototipo de software la tarea de clustering.

En este caso no se dispone de un conjunto de datos de entrenamiento, el sistema recibe las imágenes cargadas y empieza a clasificarlas automáticamente de acuerdo a la coincidencia de patrones o tendencias en ellas.

Se espera que el prototipo de software a desarrollar permita realizar clustering sobre las imágenes previamente cargadas y procesadas. Para realizar este análisis se utiliza la siguiente técnica:

- K-means: Por lo general, es el método que más se utiliza.

### **3.7 Requisitos no funcionales**

#### **Requisitos de rendimiento**

Este prototipo está enfocado a validar como la minería de imágenes puede apoyar las labores de análisis y predicción sobre imágenes del sector salud, tareas como rendimiento, optimización de recursos y operación masiva se salen del alcance del prototipo. Estas funciones se podrán ejecutar en una etapa posterior.

Por lo anterior no se pretende soportar varios usuarios de manera simultánea, tampoco se requiere un tiempo mínimo de respuesta esperado.

#### **Seguridad**

En el acceso al sistema se debe tener en cuenta las siguientes condiciones de seguridad.

- Todas las contraseñas se deben cifrar, utilizando una llave de seguridad.
- En la base de datos se almacenan la contraseñas después de aplicar el algoritmo de cifrado
- Cuando se va a iniciar sesión, la aplicación toma la contraseña ingresada por el usuario, le aplica la técnica de cifrado y hace una comparación con la contraseña almacenada en base de datos.

En la aplicación se consideran dos tipos de perfiles: administradores y operadores. A continuación se mencionan las acciones específicas que son permitidas solo a usuarios con perfil administradores.

- Gestión de usuarios
- Configuración de las técnicas de minería

Para este prototipo no se contempla la utilización de logs.

### **Fiabilidad**

Para este prototipo no se requieren factores específicos de fiabilidad. En una etapa posterior donde posiblemente se plantee la solución como un software como servicios (Saas) se pueden requerir características especiales.

### **Disponibilidad**

Para este prototipo no se requieren factores específicos de disponibilidad. En una etapa posterior donde posiblemente se plantee la solución como un software como servicios (Saas) se pueden requerir características especiales.

### **Mantenibilidad**

Para el prototipo no se especifican labores de mantenibilidad.

### **Portabilidad**

En principio, se espera que el prototipo de software sea compatible con los sistemas operativos más comunes (Windows, Linux), por tanto se recomienda la utilización de un lenguaje de programación que permita ejecutar la aplicación generada en dichos entornos.

# **Parte V: Documento de Análisis y Diseño con UML**

## 1. DESCRIPCIÓN DE LA HERRAMIENTA SOFTWARE

De acuerdo a lo expresado en la sección de levantamiento de requisitos, se desea implementar un prototipo de plataforma de minería de imágenes para el sector salud, en el cual el personal de la salud pueda cargar un conjunto de imágenes y realizarles análisis predictivo y descriptivo.

Para el análisis predictivo se desea considerar exclusivamente la predicción de categorías, clasificación, utilizando como técnicas los árboles de decisión y las máquinas de soporte vectorial. Este tipo de análisis requiere realizar etapas de entrenamiento y prueba. Para el análisis descriptivo se desea realizar la tarea de clustering, utilizando la técnicas K-means.

La plataforma debe contar con un sistema de autenticación que identifica dos tipos de usuarios, a saber, operadores y administradores. Estos últimos tiene la capacidad de configurar parámetros especiales de las técnicas de minería utilizadas, además pueden gestionar los usuarios de la plataforma, es decir, que pueden crear, modificar y eliminar usuarios. Tanto operadores como administradores disponen de una opción para modificar su información personal.

Una vez se ha accedido a la plataforma, el usuario puede iniciar un proceso de minería de imágenes, para ello es necesario realizar previamente el procesamiento de las mismas que consiste en extraer el tipo de características que se desea analizar. Las características contempladas son: color, textura y forma. Adicionalmente, el procesamiento de las imágenes requiere que previamente el usuario cargue al sistema las imágenes que desea analizar, que deben estar almacenadas en el mismo equipo donde se ejecuta la aplicación.

Durante el proceso de minería de imágenes, se le debe indicar a la plataforma la ruta de almacenamiento de las imágenes, el tipo de características a extraer, el tipo de análisis a realizar y la técnica que se desea utilizar para el análisis.

Se debe tener en cuenta que existe una alta probabilidad de que los usuarios de la plataforma no cuenten con amplios conocimientos en sistemas, por tanto la instalación y ejecución del prototipo de software deben ser procedimientos sencillos e intuitivos. Así mismo, es muy probable que la plataforma sea ejecutada en paralelo con las demás aplicaciones usadas hoy en día por los usuarios en sus computadores, es decir que no se tendrán equipos dedicados para su ejecución, por tanto se debe evitar componentes de software que requieran un uso excesivo de memoria.

## 2. DIAGRAMA DE CASOS DE USO

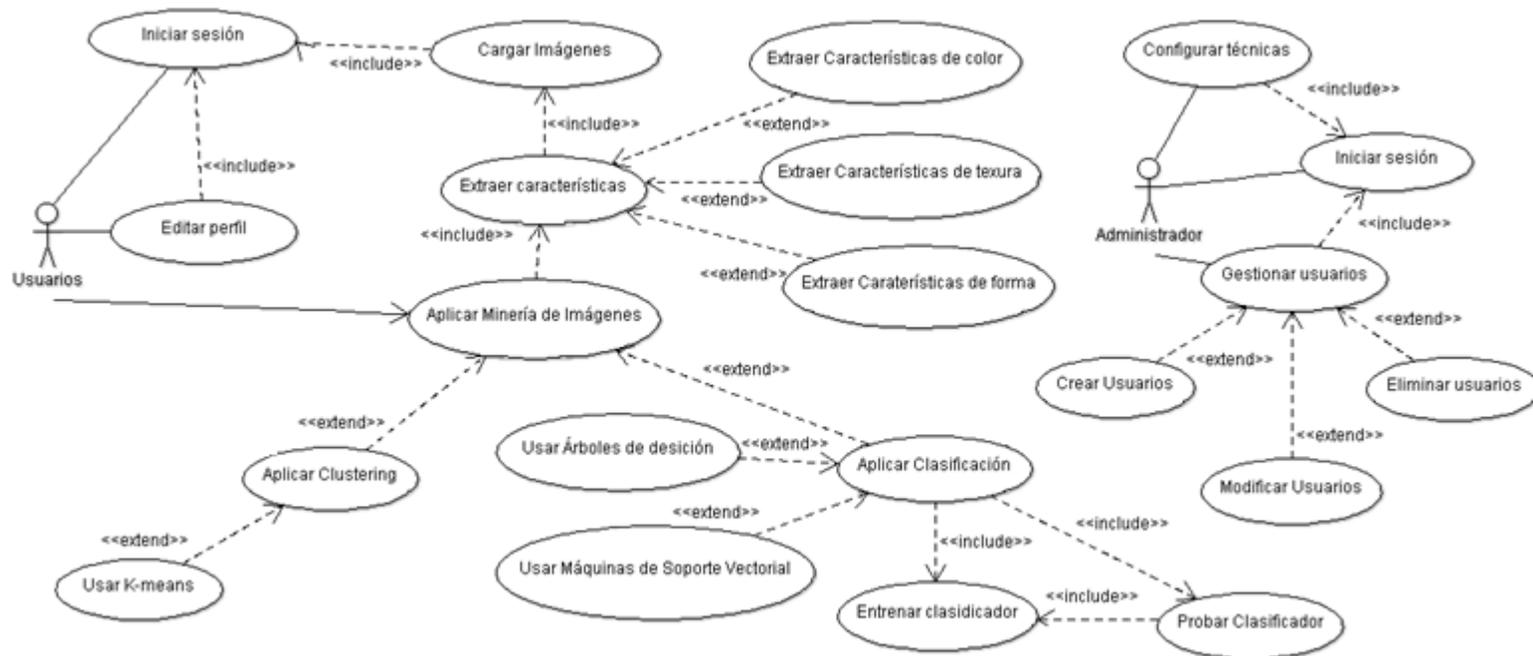


Ilustración 3: Diagrama de casos de uso

En la ilustración 3 se presenta el diagrama de casos de uso identificados para el desarrollo del prototipo de plataforma de minería de imágenes. En él se puede observar las principales funcionalidades del sistema que se detallarán a continuación.

### Escenario parcial 1: Funciones del administrador

En la Ilustración 4 se presenta el escenario parcial que corresponde a lo requerido en el siguiente apartado de la descripción del problema

*“El acceso a la plataforma debe contar con un sistema de autenticación que identifica dos tipos de usuarios, a saber, operadores y administradores. Estos últimos tiene la capacidad para configurar parámetros especiales de las técnicas de minería utilizadas, además pueden gestionar los usuarios de la plataforma, es decir, que pueden crear, modificar y eliminar usuarios.”*

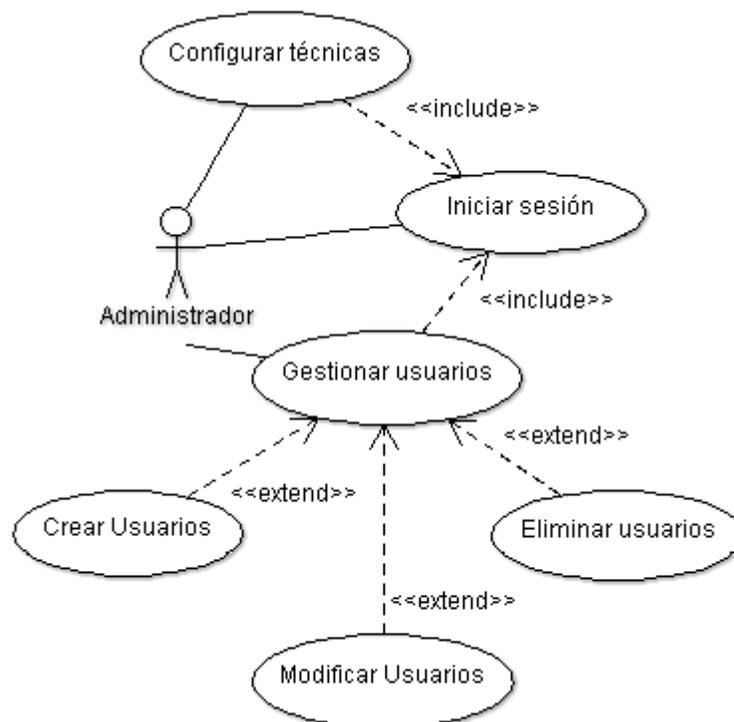


Ilustración 4: Funciones del administrador

### Escenario parcial 2: Editar perfil

En la Ilustración 5 se presenta el escenario parcial que corresponde a lo requerido en el siguiente apartado de la descripción del problema:

*“Tanto operadores como administradores disponen de una opción para modificar su información personal”*

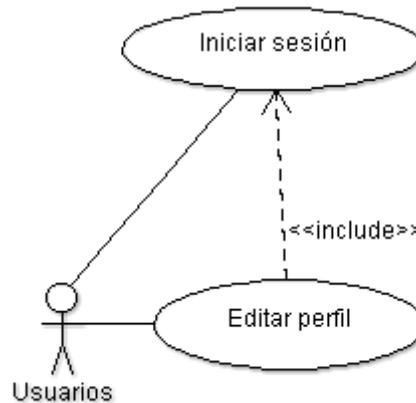


Ilustración 5: Editar perfil

### Escenario parcial 3: Aplicación de la minería de imágenes

En la Ilustración 6 se presenta el escenario parcial que corresponde a lo requerido en el siguiente apartado de la descripción del problema:

*“Una vez se ha accedido a la plataforma, el usuario puede iniciar un proceso de minería de imágenes, para ello es necesario realizar previamente el procesamiento de las mismas que consiste en extraer el tipo de características que se desea analizar. Las características contempladas son: color, textura y forma. Adicionalmente, el procesamiento de las imágenes requiere que previamente el usuario cargue al sistema las imágenes que desea analizar, que deben estar almacenadas en el mismo equipo donde se ejecuta la aplicación”*

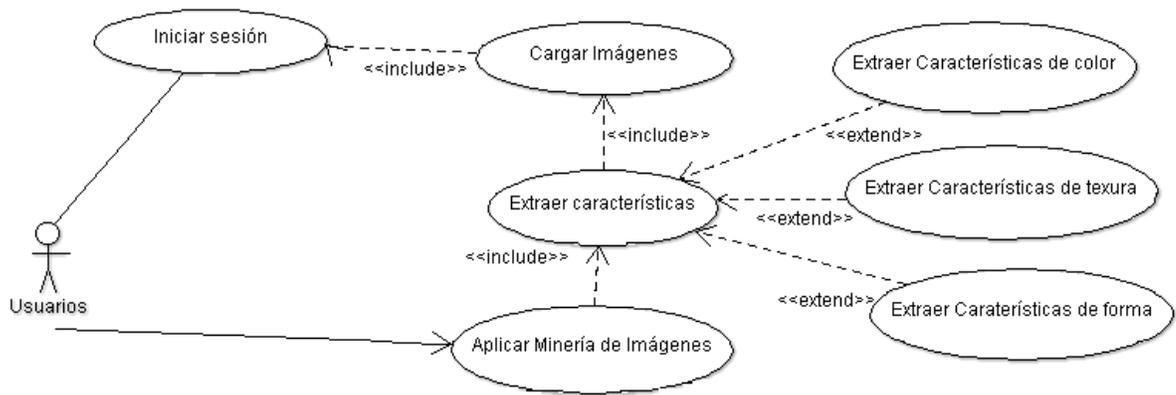


Ilustración 6: Aplicación de la minería de imágenes

#### Escenario Parcial 4: Técnicas aplicadas

En la ilustración 7 se presenta el escenario parcial que corresponde a lo requerido en el siguiente apartado de la descripción del problema:

*“Para el análisis predictivo se desea considerar exclusivamente la predicción de categorías, clasificación, utilizando como técnicas los árboles de decisión y las máquinas de soporte vectorial. Este tipo de análisis requiere realizar etapas de entrenamiento y prueba. Para el análisis descriptivo se desea realizar la tarea de clustering, utilizando la técnica K-means”*

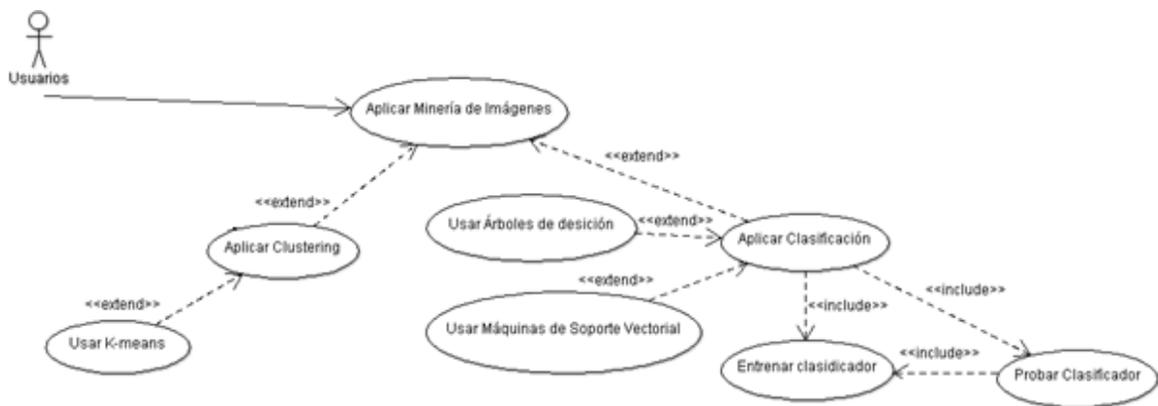


Ilustración 7: Tipos de análisis y técnicas utilizadas

### 3. DIAGRAMA DE CLASES

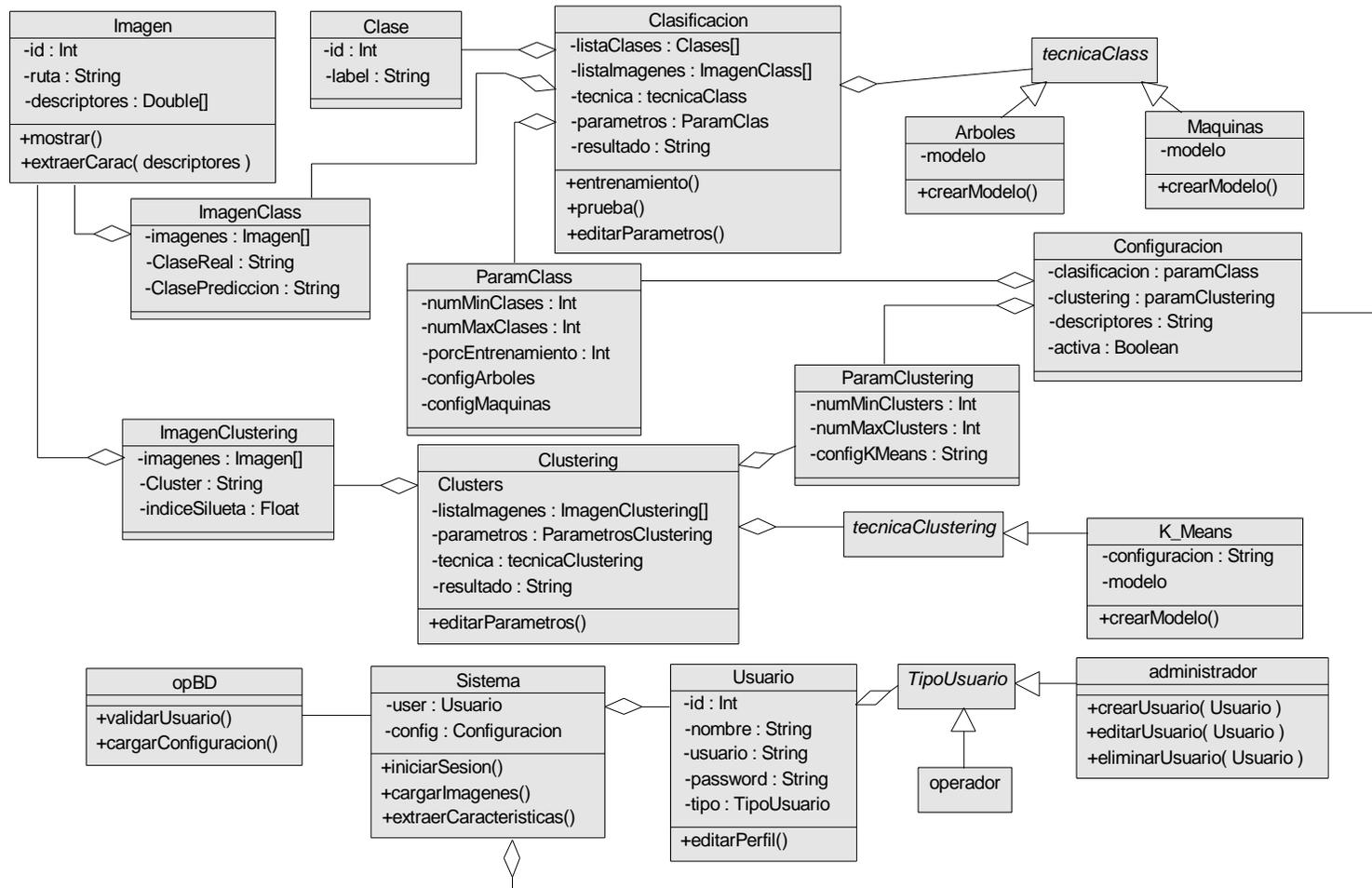


Ilustración 8: Diagrama de clases

En la ilustración 8 se presenta el diagrama de clases propuesto para la elaboración del prototipo de software de minería de imágenes.

En cuanto a los usuarios del sistema, se han contemplado clases por separado para el manejo de los usuarios de tipo operador y administrador. El usuario de tipo administrador cuenta con los métodos necesarios para realizar la gestión de usuario: crear, editar y eliminar usuario. Todos los usuarios: operador y administrador, cuentan con un método para editar su información.

Para controlar la operación del sistema se utilizan los registros de configuración, estos cuentan con parámetros de configuración para análisis descriptivo y predictivo, así como con un campo que indica cuales descriptores se encuentran activos para realizar la extracción de características. Se puede tener varias configuraciones almacenadas en la base de datos, pero solo una de ellas estará activa y es la que se utiliza para la operación del sistema.

La aplicación de la minería de imágenes requiere de dos actividades fundamentales, a saber, cargar las imágenes y extraer las características. Una vez se realizan estas dos actividades se tiene la información que el sistema requiere para crear los objetos de tipo imagen. Posteriormente se utilizará esta información en la aplicación de análisis predictivo y descriptivo.

El análisis predictivo se realiza mediante la clase Clasificación que cuenta con las técnicas árboles de decisión y máquinas de soporte vectorial. Cada una de estas técnicas se encarga de crear un modelo de clasificación, utilizando para ello la configuración asignada. El proceso de clasificación consta de dos etapas, la primera de ellas es el entrenamiento para el cual se tiene en cuenta un segmento de los datos dado por el valor de la variable porcentaje de entrenamiento, la segunda etapa consiste en la evaluación del modelo con el resto de los datos. La clase clasificación añade dos elementos importantes para cada imagen, la clase real a la que pertenece y la clase indicada por el modelo de clasificación, con estos elementos se puede realizar la evaluación del modelo.

Por su parte el análisis descriptivo se realiza mediante la clase Clustering y cuenta con la técnica K-Means que genera un modelo de datos para Clustering, utilizando la configuración asignada. Para realizar el proceso de Clustering solo se requiere indicar el número de cluster, crear el modelo y aplicarlo sobre cada una de las imágenes. Cada imagen de Clustering cuenta con un campo para indicar el cluster asignado a la imagen, y otro campo para indicar el índice la silueta obtenido, que proporciona una medida de evaluación del agrupamiento.

#### 4. DIAGRAMA DE SECUENCIA

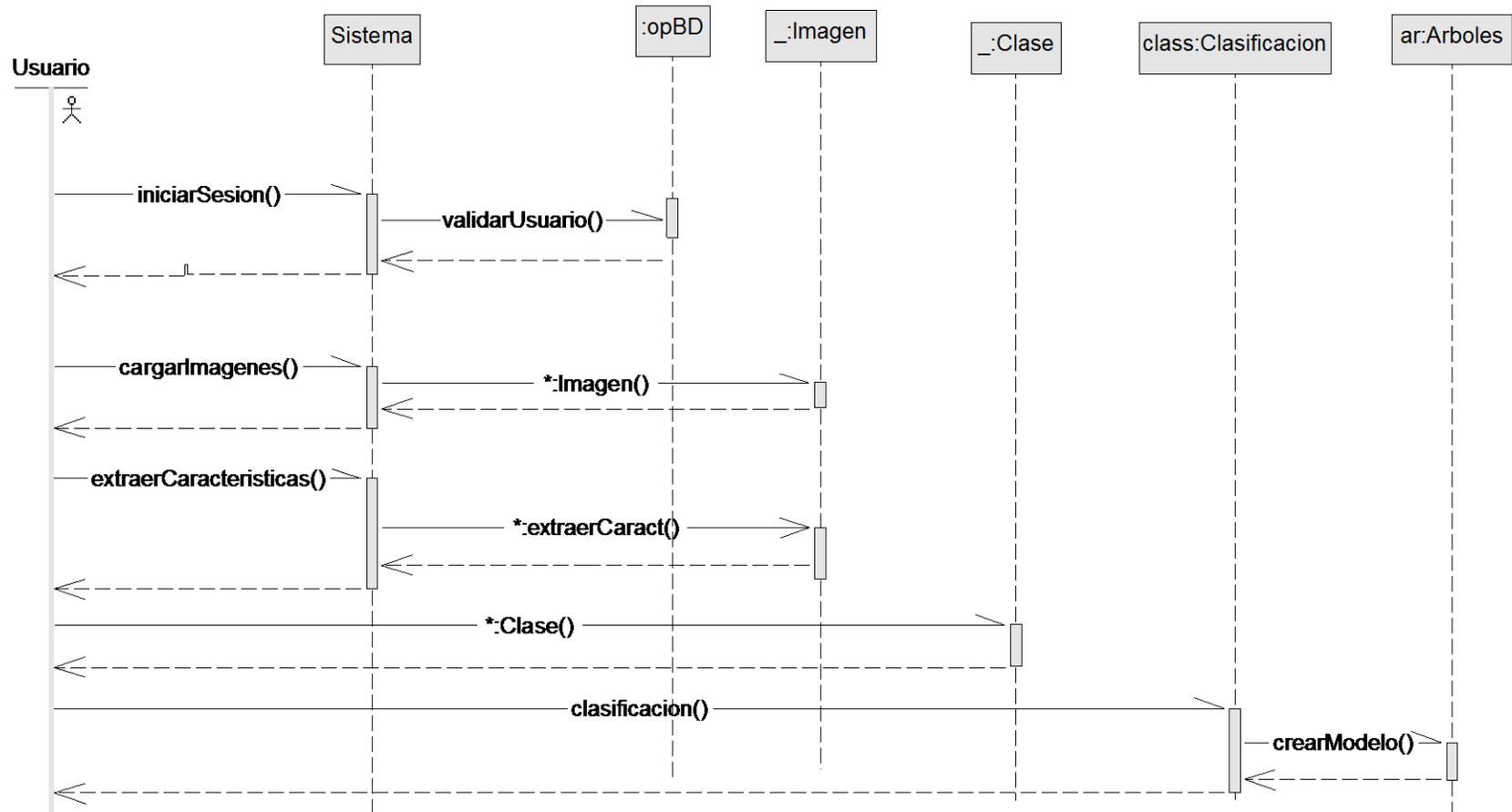


Ilustración 9: Diagrama de secuencia

En la Ilustración 9 se presenta un diagrama de secuencia donde se analiza el proceso de creación de un modelo de clasificación utilizando la técnica árboles de decisión.

En primer lugar, el usuario debe iniciar sesión indicando los valores para los campos nombre de usuario y password, con los que se realizará el proceso de login. Este proceso se realiza mediante una validación contra la base de datos. En caso de que la validación sea exitosa el usuario continúa con la carga de imágenes donde debe indicar la ruta que utilizará el sistema para buscar todos los archivos con extensión JPG, BMP o PNG que pertenecen a dicho directorio. En este momento se crea para cada archivo, con una de las extensiones especificadas, un objeto de la clase Imagen. Este proceso termina cuando se analicen todos los archivos de la ruta seleccionada por el usuario para la carga de imágenes.

Una vez se han cargado las imágenes se procede a realizar el proceso de extracción de características. Esto se hace en cada uno de los objetos de tipo Imagen creados en el paso anterior. Cuando se han extraído las características de cada una de las imágenes, el usuario debe proceder a indicar el número de clases. En este momento se crean objetos de tipo Clase indicando un label para su posterior identificación.

Al tener identificadas las clases, las imágenes y sus descriptores, se procede a crear un objeto de tipo clasificación. En este caso el usuario debe indicar para cada una de las imágenes, la clase real a la que pertenece. Con esta información, se puede crear el modelo para la clasificación, que en este caso se hace utilizando un objeto de la clase árbol, ya que esta fue la técnica seleccionada para el proceso de clasificación en este ejercicio.

Con la creación del modelo el sistema está preparado para evaluarlo y mostrar el resultado obtenido.

## 5. DIAGRAMA DE ACTIVIDADES

En la ilustración 10 se presenta un diagrama de actividades donde se analiza el proceso de aplicar un tipo de análisis.

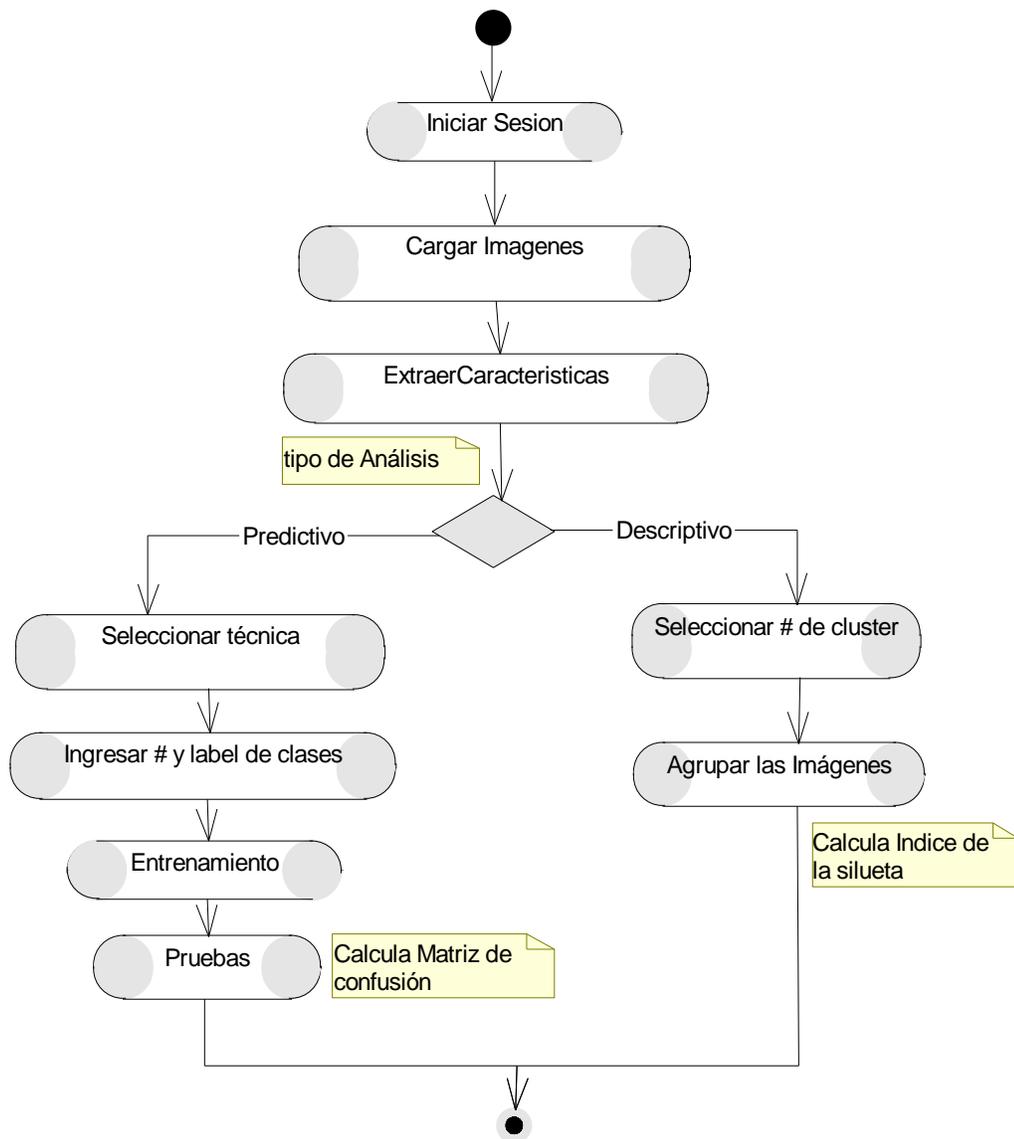


Ilustración 10: Diagrama de actividades

Lo primero que se debe hacer es iniciar sesión, luego se debe realizar el procedimiento de carga de imágenes. Una vez cargadas las imágenes se procede a extraer las características de cada una de ellas.

Posteriormente se debe tomar una decisión sobre el tipo de análisis que se desea realizar.

En caso de que se escoja el análisis predictivo se debe seleccionar la técnica a utilizar (árboles de decisión o máquinas de soporte vectorial), ingresar el número de clases indicado un label para cada una de ellas y realizar la etapas de entrenamiento y prueba. Como resultado del proceso, el sistema presenta la matriz de confusión obtenida de la evaluación del modelo de clasificación creado.

En caso de que se escoja el análisis de tipo descriptivo, solo se debe indicar el número de clusters esperados. Esta información es suficiente para que el sistema realice el proceso de agrupamiento e indique para cada una de las imágenes el cluster asignado y el índice de la silueta respectivo.

# **Parte VI: Documento de Pruebas**

## 1. PLAN DE PRUEBAS

En esta sección se realiza una serie de pruebas para validar la funcionalidad del prototipo de software, haciendo énfasis en los aspectos más relevantes del levantamiento de requisitos presentado en la sección IV. Para ello se realiza este plan de pruebas donde se determinan las acciones a realizar.

### **Ingreso al sistema**

Esta actividad tiene como entradas el usuario y el password ingresados por el usuario, y como salida debe permitir o denegar el acceso a la plataforma. Se debe verificar que el sistema permita el acceso cuando se ingresan usuario y contraseña correcta y lo deniegue cuando una de ellas o ambas sea incorrecta. Validar que sucede si alguno de los campos no es ingresado.

### **Funciones exclusivas del administrador**

Se identificaron algunas funciones que son exclusivas para los usuarios con perfil de administración, es decir que los usuarios con perfil de operador no deben poder acceder a dichas funciones. Para esta prueba la entrada será el tipo de usuario con el cuál se acceda al sistema y la salida esperada es que se habiliten las funciones exclusivas para administradores y se deshabiliten para operadores.

### **Administración de usuarios**

Ingresar al sistema con un usuario que tenga perfil de administración y validar que sea posible crear, editar y eliminar usuarios. Se espera que el sistema presente pantallas para la realización de estas funciones

### **Editar perfil**

Validar que los usuarios tengan una opción para visualizar y editar la configuración de su cuenta. Se espera que el sistema presente en una pantalla la información de la cuenta del usuario y permita su edición.

### **Configuraciones**

Ingresar al sistema con un usuario que tenga perfil de administración y validar que sea posible acceder a las configuraciones. Estas deben permitir la configuración de las técnicas de minería utilizadas.

### **Cargar imágenes**

Como entrada se le indicará al sistema la ruta donde se encuentran almacenadas unas imágenes de prueba, como salida se espera que el sistema detecte las imágenes de la ruta especificada y habilite la función de extracción de características.

### **Extracción de características**

El sistema debe presentar una opción para que se le indique el tipo de características que se desean utilizar. Seleccionar un tipo de característica y verificar que se realice el proceso de extracción de características de cada una de las imágenes cargadas previamente.

### **Aplicación de análisis predictivo**

El sistema debe brindar una opción para seleccionar una de las técnicas del análisis predictivo, también debe permitir la selección del número de clases y la clase a la que pertenece cada una de las imágenes. Luego de esto debe realizar el proceso de entrenamiento y mostrar los resultados obtenidos.

### **Aplicación de análisis descriptivo**

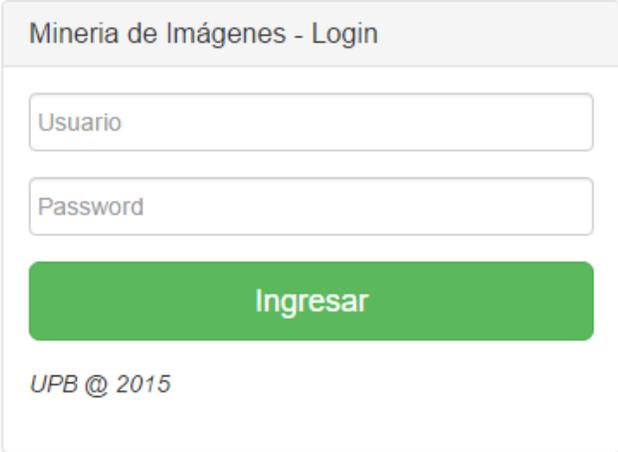
Al seleccionar la técnica K-Means y el número de clusters, el sistema debe realizar el Clustering y mostrar los resultados obtenidos.

## 2. PRUEBAS DEL PROTOTIPO DE SOFTWARE

A continuación se presenta el resultado de las pruebas realizadas de acuerdo al plan de pruebas establecido.

### Ingreso al sistema

En la Ilustración 11 se presenta la pantalla de acceso al sistema. Se evidencia la presencia de campos para ingreso de usuario y contraseña, y un botón para realizar el acceso.



Minería de Imágenes - Login

Usuario

Password

Ingresar

UPB @ 2015

Ilustración 11: Pantalla de acceso al sistema

Para validar el acceso al sistema, se realizaron tres pruebas. La primera prueba consistió en presionar el botón “Ingresar” sin haber digitado el password, la segunda prueba es similar solo que esta vez se digita el password pero no el usuario y en la tercera prueba se digita un usuario correcto con un password incorrecto. El resultado de esta prueba se presenta en la Ilustración 12.



**Ilustración 12: Pruebas de acceso al sistema**

En la Ilustración 12 la imagen de la izquierda indica que no se ha ingresado el password, la imagen del centro indica que no se ha ingresado el usuario y la imagen de la derecha indica que el usuario o la contraseña son incorrectos.

Posteriormente se ingresó un usuario y contraseña correcta y el sistema presentó la pantalla de inicio presentada en la Ilustración 13.



**Ilustración 13: Pantalla de inicio**

## Funciones exclusivas del administrador

En la Ilustración 13 se puede ver el menú de operación donde se permite el acceso a las diversas funciones de la plataforma. La imagen presentada corresponde a un usuario con perfil de administrador, se evidencia que tiene acceso a las opciones de “Configuración” y “Usuarios”. Iniciando sesión con un usuario de tipo operador se observa una pantalla similar pero no se muestran las dos funciones mencionadas. Esto garantiza que solo los usuarios con perfil de administración pueden acceder a dichas funcionalidades.

## Administración de usuarios

Efectivamente solo fue posible ingresar a esta opción con usuarios de tipo administrador. Inicialmente se presenta un listado de los usuarios registrados en el sistema donde se puede observar el nombre, usuario para acceder a la plataforma y tipo. También se evidencia la presencia de opciones para editar o eliminar usuarios actuales y crear un nuevo usuario. Todo esto se puede observar en la Ilustración 14.

# Administrar Usuarios



Lista de usuarios

[Crear](#)

Id	Nombre	Usuario	Tipo	Opciones
1	Administrador por defecto	admin	Administrador	 
2	Operador por defecto	operador	Operador	 
3	Ana Isabel Oviedo	anaisaoviedo	Administrador	 
4	Efraín Alberto Oviedo	eaoc46	Administrador	 

Ilustración 14: Pantalla de administración de usuarios

Al acceder a la opción “Crear”, se presentó una ventana modal con los campos necesarios para la creación de un nuevo usuario: nombre, usuario, tipo de usuario y contraseña. Este último campo tiene una confirmación que solicita de nuevo la contraseña para asegurarse que se ha digitado correctamente. Se crearon varios usuarios de tipo administrador y operador, y el sistema se comportó según lo esperado.

En una prueba posterior se seleccionó la opción para editar uno de los usuarios creados en el sistema, en este caso apareció una ventana similar a la de crear usuario donde los campos estaban diligenciados con la información del usuario seleccionado y se permitía la modificación de los mismos. En la Ilustración 15 se presentan las ventanas de creación y edición de usuarios.

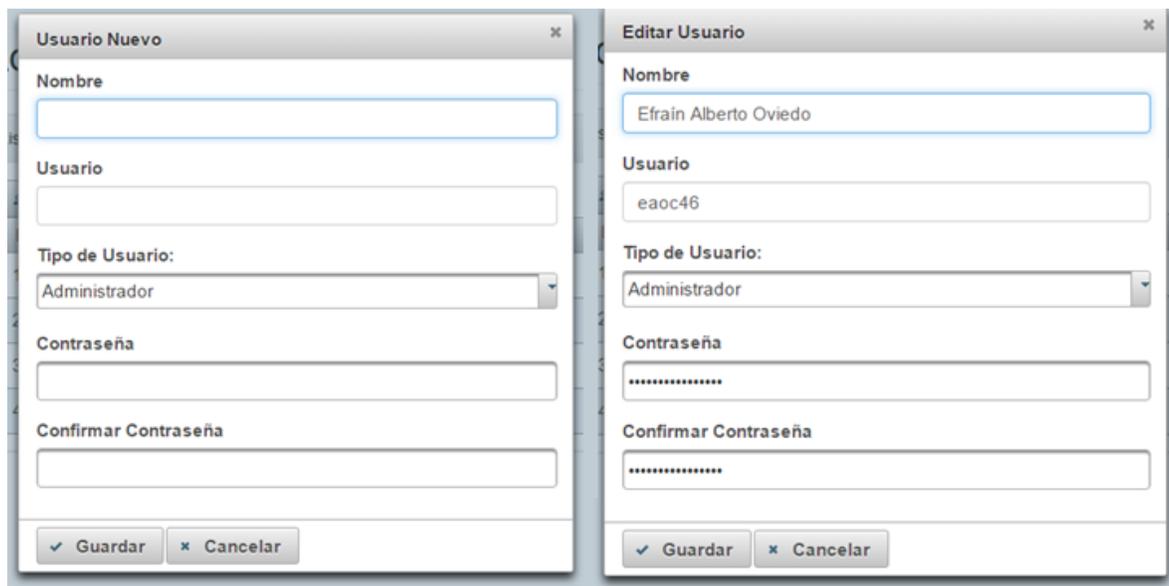


Ilustración 15: Crear y editar usuarios

En la creación de usuarios se realizó una prueba adicional, donde se ingresó de forma errada la confirmación de la contraseña. En este caso el sistema presentó un mensaje donde informaba que las contraseñas no coinciden.

También se validó la opción de eliminar usuario, esta presenta un cuadro de confirmación para eliminar el usuario seleccionado. Se observó un comportamiento adecuado.

### Editar perfil

Se ingresó al sistema con usuarios de tipo administrador y operador, y se observó que ambos tenían acceso a editar su información personal: nombre, usuario, tipo de usuario y contraseña. Cuando el usuario es de tipo administrador se permite seleccionar tipo de usuario: operador o administrador, pero si el usuario es operador en el campo tipo de usuario solo se visualiza la opción operador. Esto garantiza que un usuario de tipo operador no puede cambiarse por sí mismo a administrador. En la Ilustración 16 se presenta la pantalla de editar perfil.

## Editar Perfil

**Nombre**

**Usuario**

**Tipo de Usuario:**

**Contraseña**

**Confirmar Contraseña**

**Ilustración 16: Editar perfil**

Se comprobó que para que la nueva configuración sea almacenada, los campos contraseña y confirmar contraseña deben coincidir, de lo contrario se presenta un mensaje donde se informa que las contraseñas no coinciden.

## Configuraciones

Al ingresar en la opción de configuraciones se presenta un listado con las configuraciones que se han creado en el sistema, donde se observan los siguientes campos: nombre, autor, fecha y hora, activo. Esta pantalla se presenta en la Ilustración 17.

## Configuración

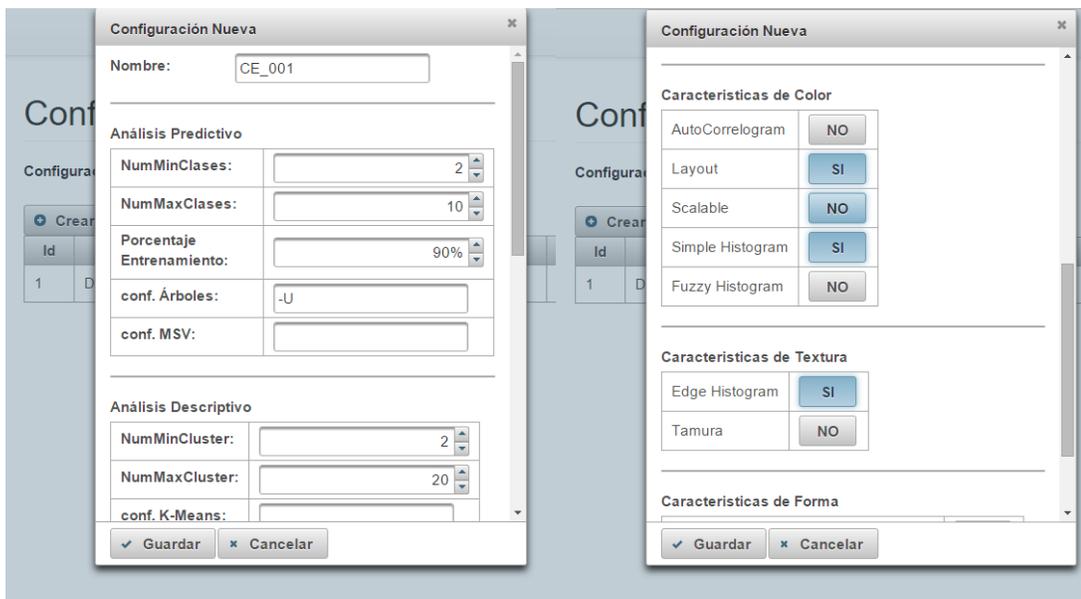
Configuraciones almacenadas

Id	Nombre	Autor	Fecha y hora	Activo	Opciones
1	Defecto admin_1	admin	2015-09-11 04:04:20	true	<input type="button" value="★"/> <input type="button" value="✎"/> <input type="button" value="🗑"/>

**Ilustración 17: Pantalla de configuraciones**

Para cada una de las configuraciones creadas en el sistema se presentan opciones para: asignar como configuración activa, editar configuración y eliminar configuración. Se validó que solo una configuración puede estar activa y esta no puede ser eliminada mientras permanezca como activa. Durante la prueba se observó un comportamiento adecuado de estas opciones.

Mediante el botón “Crear” se desplegó una ventana modal que permite la creación de una nueva configuración. En esta ventana es posible configurar parámetros para el análisis predictivo y descriptivo. En el análisis predictivo se puede configurar el número mínimo y máximo de clases, el porcentaje de datos de entrenamiento y las opciones de los clasificadores (árboles de decisión y máquinas de soporte vectorial). En el análisis descriptivo es posible asignar el número mínimo y máximo de cluster y las opciones de K-means. También es posible habilitar o deshabilitar los tipos de descriptores que se van a extraer de las imágenes para cada tipo de característica. En la Ilustración 18 se puede observar la ventana de creación de una nueva configuración.



**Ilustración 18: Crear configuración**

Lo presentado en la Ilustración 18 corresponde a una sola ventana, se ha duplicado con fines de visualización, nótese el desplazamiento de la barra horizontal.

Para ingresar los parámetros de configuración de las técnicas de minería de datos se ha utilizado un campo de entrada de texto donde se deben ingresar, separados por coma, cada una de las configuraciones deseadas. Estas configuraciones cumplen con el formato establecido por WEKA. Por ejemplo en el caso de árboles de decisión se ha ingresado el comando -U que indica que se quiere el árbol sin podar. En este caso no se añadió coma ya que se trata de un solo parámetro de configuración.

## Cargar imágenes

El sistema presentó una pantalla donde se observa el botón “Cargar Directorio” al presionarlo apareció una ventana para seleccionar la ruta donde se encuentran almacenadas las imágenes, se seleccionó la ruta D:\Image Mining\Clustering2, y al aceptar el sistema detectó un total de 81 imágenes en formato JPG que se encontraban en el directorio seleccionado. En la Ilustración 19 se presenta la pantalla de preparar imágenes donde se encuentra esta funcionalidad.

# Preparar Imágenes

Seleccione la ruta donde se encuentran las imagenes



Ruta seleccionada: D:\Image Mining\Clustering2

Numero de imagenes: 81

[Reiniciar](#)

Id	Ruta	Opciones
1	D:\Image Mining\Clustering2\Mamografias\mdb007.jpg	
2	D:\Image Mining\Clustering2\Mamografias\mdb008.jpg	 Ver
3	D:\Image Mining\Clustering2\Mamografias\mdb009.jpg	
4	D:\Image Mining\Clustering2\Mamografias\mdb011.jpg	
5	D:\Image Mining\Clustering2\Mamografias\mdb014.jpg	
6	D:\Image Mining\Clustering2\Mamografias\mdb016.jpg	

Ilustración 19: Pantalla preparar imágenes

Las imágenes se presentan de forma enumerada en una tabla que contiene los campos: Id, ruta y opciones. Dentro de opciones se encuentra un botón que permite visualizar la imagen deseada. Al hacer presionar este botón en una de las imágenes listadas, se presenta una ventana modal donde se puede visualizar la imagen seleccionada.

Antes de cargar las imágenes, si se trataba de ingresar en la opción “Extraer características” aparecía un mensaje informando que no era posible acceder a esta opción hasta que no se hiciera la carga de imágenes. Una vez cargadas las imágenes, no volvió a salir este mensaje y se permitió el acceso a la opción de extraer características.

## Extraer características

En la Ilustración 20 se presenta la pantalla extraer características. Esta pantalla dispone de un mecanismo para seleccionar el tipo de características a extraer, puede ser seleccionada una, dos o las tres características disponibles. Para la selección de más de una característica se debe utilizar la tecla Ctrl o Shift y el mouse.

Para la prueba se seleccionó Características de color y al presionar el botón “Extraer” se realizó la extracción de características de cada una de las imágenes cargadas. Durante este proceso una barra indicaba el progreso de esta acción y a finalizar se indica que se han extraído 31 descriptores de cada una de las 81 imágenes cargadas. En este caso solo se utilizó el descriptor “Layout” ya que en el momento era el único descriptor de color habilitado en la configuración activa.

# Extraer Características

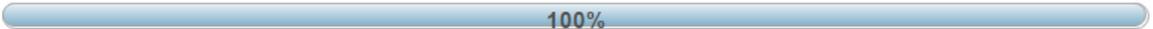
Seleccione el tipo de características que desea extraer

- Características de Color
- Características de Forma
- Características de Textura

Opciones

- Extraer
- Reiniciar

Progreso



Numero de imagenes: 81

Numero de descriptores: 33

Características extraidas

Color				
Id	Nombre	Descriptor Inicial	Descriptor Final	Total
1	Layout	0	32	33

Ilustración 20: Pantalla extraer características

## Aplicación de análisis predictivo

El acceso a esta función solo se habilitó después de haber realizado la extracción de características, antes de esto aparecía un mensaje de error.

Para el entrenamiento se seleccionó el número de clases (3) y la técnica a utilizar (Árboles de decisión). De inmediato apareció una opción para indicar el label que se le quiere dar a cada una de las clases. En este caso la carpeta que se seleccionó para cargar las imágenes, contenía 3 sub carpetas, al coincidir el número de sub carpetas con el número de clases indicadas, el sistema asume que las imágenes se encuentran agrupadas en la carpeta que corresponde con la clasificación esperada, por tanto se propone el nombre de la carpeta como label de cada clase. Así mismo, a cada imagen se le asigna por defecto como clase real, la carpeta a la que pertenece. Esto facilita el hecho de tener que seleccionar manualmente para cada una de las imágenes, la clase a la que pertenece.

En la Ilustración 21 se presenta la pantalla para el entrenamiento en el análisis descriptivo.

## Análisis Predictivo - Entrenamiento

**Seleccione el número de clases**

**Seleccione técnica a Utilizar**

Máquinas de soporte Vectorial

Árboles de decisión

---

**Identificación de las clases**

Id	Label
1	<input type="text" value="Mamografías"/>
2	<input type="text" value="Resonancia Magnética"/>
3	<input type="text" value="Torax"/>

**Opciones**

---

**Indique a cual clase pertenece cada imagen**

Id	Ruta	Clase Real	Clase Pred	Opciones
1	D:\Image Mining\Clustering2\Mamografias\mdb007.jpg	<input type="text" value="Mamografías"/>	NULL	<input type="button" value="📄"/>
2	D:\Image Mining\Clustering2\Mamografias\mdb008.jpg	<input type="text" value="Mamografías"/>		<input type="button" value="📄"/>
3	D:\Image Mining\Clustering2\Mamografias\mdb009.jpg	<input type="text" value="Torax"/>		<input type="button" value="📄"/>
4	D:\Image Mining\Clustering2\Mamografias\mdb014.jpg	<input type="text" value="Mamografías"/>	NULL	<input type="button" value="📄"/>

**Ilustración 21: Pantalla de análisis predictivo- Entrenamiento**

En caso de que el número de clases no coincide con el número de subcarpetas, el sistema propone por defecto labels consecutivos para las clases y asigna a todas las imágenes la clase 1 como clase real. En este caso es necesario seleccionar imagen por imagen la clase real.

Al final de la tabla que presenta las imágenes para el análisis predictivo se encuentra el botón “Entrenar”, al presionarlo se realiza el entrenamiento y se visualiza en un cuadro de texto el resultado obtenido: la matriz de confusión y algunos datos de interés derivados de esta. Adicionalmente en la tabla de las imágenes se llena el campo “Clas Pred” con la clase predicha por el modelo creado. Esta clase aparece en color verde si coincide con la clase real o en color rojo si no coincide.

### Aplicación de análisis descriptivo

Al ingresar en esta opción se seleccionaron 3 Cluster y la técnica K-Means, posteriormente el sistema arrojó el resultado del proceso de Clustering en un cuadro de texto donde se pueden observar los centroides. También se presenta una tabla con la lista de imágenes, donde se presenta el cluster donde quedó agrupada cada imagen y el índice de la silueta respectivo.

En la **¡Error! No se encuentra el origen de la referencia.** Se presenta la pantalla de análisis descriptivo.

## Análisis Descriptivo

**Seleccione el número de clusters**

**Seleccione técnica a Utilizar**

K-Means

---

**Resultado del clustering**

```

kMeans
=====

Number of iterations: 2
Within cluster sum of squared errors: 290.2652137871122
Missing values globally replaced with mean/mode

Cluster centroids:
      Cluster#
Attribute Full Data  0    1    2
          (81) (27) (27) (27)
          =====
    
```

**Cluster identificado para cada imagen**

Id	Ruta	Cluster	ind. Silueta	Opciones
1	D:\Image Mining\Clustering2\Mamografias\mdb007.jpg	Cluster0	0.8194654	
2	D:\Image Mining\Clustering2\Mamografias\mdb008.jpg	Cluster0	0.7608447	

**Ilustración 22: Pantalla de análisis descriptivo**

# **Parte VII: Casos de Estudio, Conclusiones y Trabajo Futuro**

# 1. CASO DE ESTUDIO 1: PREDICCIÓN DE ANORMALIDADES EN MAMOGRAFÍAS

En este caso de estudio, se aplicó la metodología propuesta en este trabajo KDM (Knowledge Discovery in Multimedia), la cual es una modificación de la metodología KDD para incluir etapas de tratamiento de multimedia. En esta sección se presenta la aplicación de las fases de la metodología KDM para la predicción de anormalidad es en mamografías.

## 1.1 Descripción del caso de estudio

La mamografía es un examen médico que consiste en tomar una radiografía de los senos con el fin de detectar signos de cáncer de seno, de hecho es una de las principales fuentes de diagnóstico de esta enfermedad. Con una mamografía se busca hacer una detección temprana de la enfermedad lo que permite aplicar un tratamiento apropiado y así tratar de mejorar la calidad de vida el paciente. Mediante este examen es posible detectar tumores que normalmente no pueden palparse. Así mismo se pueden detectar micro calcificaciones que son pequeños depósitos de calcio que pueden indicar la presencia de cáncer de seno.

Si la mamografía arroja un resultado normal, se recomienda continuar haciendo mamografías periódicamente con el fin de comparar los resultados para detectar si hay cambios en los senos. En caso de que el resultado sea anormal se debe hacer algunos exámenes especializados para determinar si existe presencia de cáncer, en muchos casos las anomalías pueden ser benignas, esto significa que no hay nada de qué preocuparse pero se debe continuar haciendo el examen periódicamente para hacer seguimiento.

## 1.2 Aplicación de la metodología KDM al caso de estudio

A continuación se presenta una descripción de la aplicación de las fases de la metodología KDD en la predicción de anomalías en mamografías.

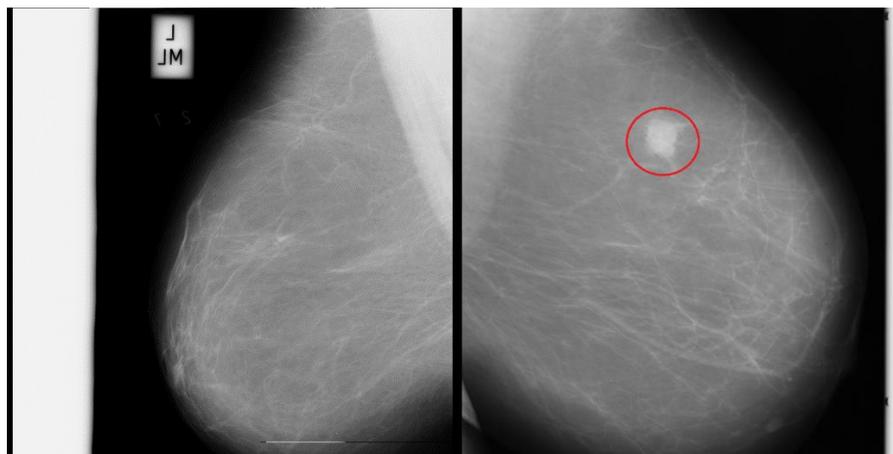
### 1.2.1 Selección de las imágenes

Para este caso de estudio se utilizó la base de datos MIAS (Mamography Image Analyze Society), esta es una sociedad inglesa que se dedica a la investigación de las mamografías. Esta base de datos (Suckling, y otros, 1994) incluye exámenes de ambos senos de 161 pacientes para un total de 322 imágenes. De estas imágenes 208 presentan un diagnóstico de normalidad y 114 presentan

algún tipo de anomalía, estas últimas a su vez se encuentran clasificadas en dos grupos, el primero de ellos consta de 63 imágenes en las que se detecta un cáncer benigno y el otro grupo presenta 51 imágenes que presentan cáncer maligno.

En total se seleccionaron 197 imágenes (se seleccionaron de forma aleatoria, corresponde a un valor cercano al 60% del total de las imágenes de la base de datos) de las cuales 108 presentan normalidad y 89 presentan algún tipo de anomalía. Con estas imágenes se realiza un proceso de clasificación con el fin de determinar si se trata de un caso de normalidad o anomalía.

En la Ilustración 23 se presentan dos mamografías seleccionadas de la base de datos MIAS. La figura de la izquierda corresponde a un caso diagnosticado como normal, mientras que la imagen de la derecha presenta una anomalía que ha sido encerrada con fines ilustrativos.



**Ilustración 23: Mamografías de la base de datos MIAS**

La anomalía detectada consiste en una pequeña masa que presenta presencia de cáncer maligno. Los diagnósticos de la base de datos MIAS han sido realizados por radiólogos expertos.

### **1.2.2 Preprocesamiento de las imágenes**

En esta fase se realizan varias operaciones descritas a continuación.

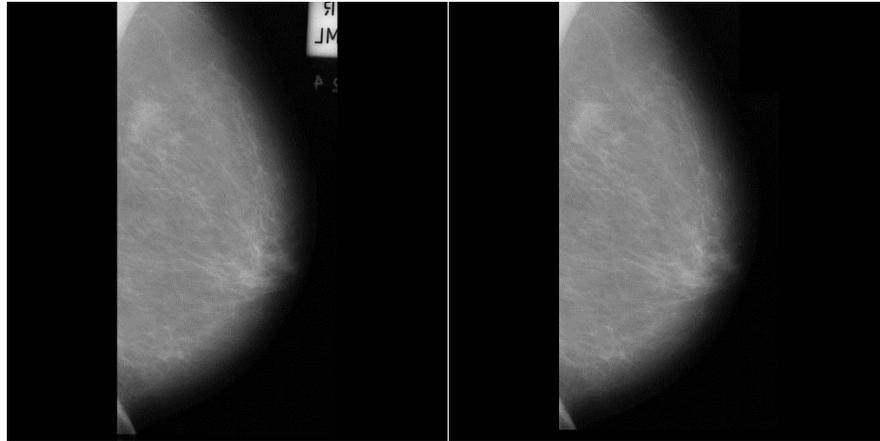
#### **Conversión de las imágenes**

Las imágenes de esta base de datos se encuentran en formato PGM. Este formato no es compatible con el prototipo de software, por este motivo fue necesario realizar una etapa de conversión de las imágenes del formato PGM al formato JPG compatible con el sistema.

## Segmentación

Como se puede observar en la Ilustración 23, las imágenes de la base de datos MIAS presentan algunos elementos visuales como letras, etiquetas y bordos que pueden desviar la atención y presentar una alta influencia a la hora de clasificar las imágenes. Estos son elementos innecesarios que no aportan para el proceso, que debe centrarse en la información contenida dentro del seno.

Un proceso de segmentación permite identificar claramente la región de interés de una imagen, eliminando todos aquellos elementos distractores presentes en la misma. En la Ilustración 24 se presenta un ejemplo de segmentación, la figura de la izquierda corresponde a la imagen sin segmentación, mientras que la figura de la derecha corresponde a la imagen segmentada.



**Ilustración 24: Segmentación de Mamografías**

Puede observarse como la imagen segmentada solo contiene información de interés, se han eliminado los elementos distractores de la imagen.

Existen diversos mecanismos para realizar de forma automática el proceso de segmentación para extraer de una imagen la región correspondiente al seno. En el caso de estudio se quiere resaltar la tarea de clasificación más que la segmentación, por tanto la segmentación de las imágenes se ha realizado de forma manual, por fuera del prototipo de plataforma. Se contempla como un trabajo a futuro la inclusión de mecanismos de segmentación automática dentro del software.

## Cambios en el espacio de color

Las imágenes son tratadas en los espacios de color rgb y hsv.

### 1.2.3 Indexamiento

Para realizar el proceso de clasificación es necesario realizar la extracción de características de cada una de las imágenes. Las características extraídas deben permitir la identificación de las imágenes en las dos clases que se pretende clasificar: diagnóstico normal y diagnóstico anormal.

En las imágenes seleccionadas se observa que las mamografías con diagnóstico anormal presentan algunas masas con una textura, forma y color representativos. Por esta razón se optó por extraer simultáneamente características de color, forma y textura. Las características extraídas son:

- Fuzzy Histogram
- Simple Histogram
- LBP (Local Binary Patterns)
- PHOG (Pyramid Histogram of Oriented Gradients)
- RILBP (Rotation Invariant Local Binary Patterns)
- Edge Histogram
- Tamura

Cada una de estas características cuenta con varios descriptores que se pueden extraer de una imagen, en total se extrae 1209 descriptores de cada imagen.

En esta ocasión no es necesario realizar transformaciones de las características, donde todas son numéricas, así que se pasa a la etapa de modelado.

### 1.2.4 Modelamiento

Para realizar clasificación, el prototipo construido ofrece dos técnicas diferentes: máquinas de soporte vectorial y árboles de decisión. Para cada una de estas técnicas se crearon modelos con los diferentes conjuntos de datos generados en la extracción de características.

En los modelos creados se ha contemplado un 70% de los datos para la etapa de entrenamiento y un 30% para la etapa de evaluación. Estos conjuntos de datos se crean de forma aleatoria al generar el modelo.

La forma de evaluar los modelos es mediante la matriz de confusión, que resume el comportamiento de modelo con el 30% de los datos seleccionados para la evaluación. El modelo a elegir es aquel que brinde un mayor porcentaje de éxito en la clasificación.

### 1.2.5 Evaluación

El modelo que presentó un mejor resultado utiliza la técnica de árboles de decisión. Con un 86% de instancias clasificadas correctamente. En la Tabla 18 se presenta la matriz de confusión generada.

Clasificado como		Diagnóstico real
Normal	Anormal	
34	2	Normal
6	18	Anormal

Tabla 18: Matriz de confusión caso de estudio 1

De la matriz de confusión se calculan las siguientes medidas que indican la calidad del clasificador construido.

TP Rate	FP Rate	Precision	Recall	F-Measure	ROC Area	Class
0.944	0.25	0.85	0.944	0.895	0.856	Normal
0.75	0.056	<b>0.9</b>	<b>0.75</b>	0.818	0.856	Anormal
0.867	0.172	0.87	0.867	0.864	<b>0.856</b>	Weighted Avg

Tabla 19. Medidas de evaluación del clasificador caso de estudio 1

De la tabla 19 es importante resaltar las medidas de precisión y cobertura (recall) de la clase Anormal. La precisión indica que el 90% de anomalías reconocidas estaban correctas, es decir que se tuvieron el 10% de falsos positivos. La cobertura indica que se reconocieron el 75% de anomalías, es decir que 25% de las anomalías no fueron reconocidas. Finalmente el área ROC indica un desempeño general del 85.6%, lo cual lo posiciona en un buen clasificador.

En la Ilustración 25 se presenta el árbol de decisión generado para lograr esta clasificación. En la parte superior del árbol se ubican las características de mayor relevancia para hacer la clasificación y descendiendo por el árbol se pueden alcanzar las clases llamadas “problema” y “ok” que representan las anomalías y normalidades en las imágenes.

Tree View

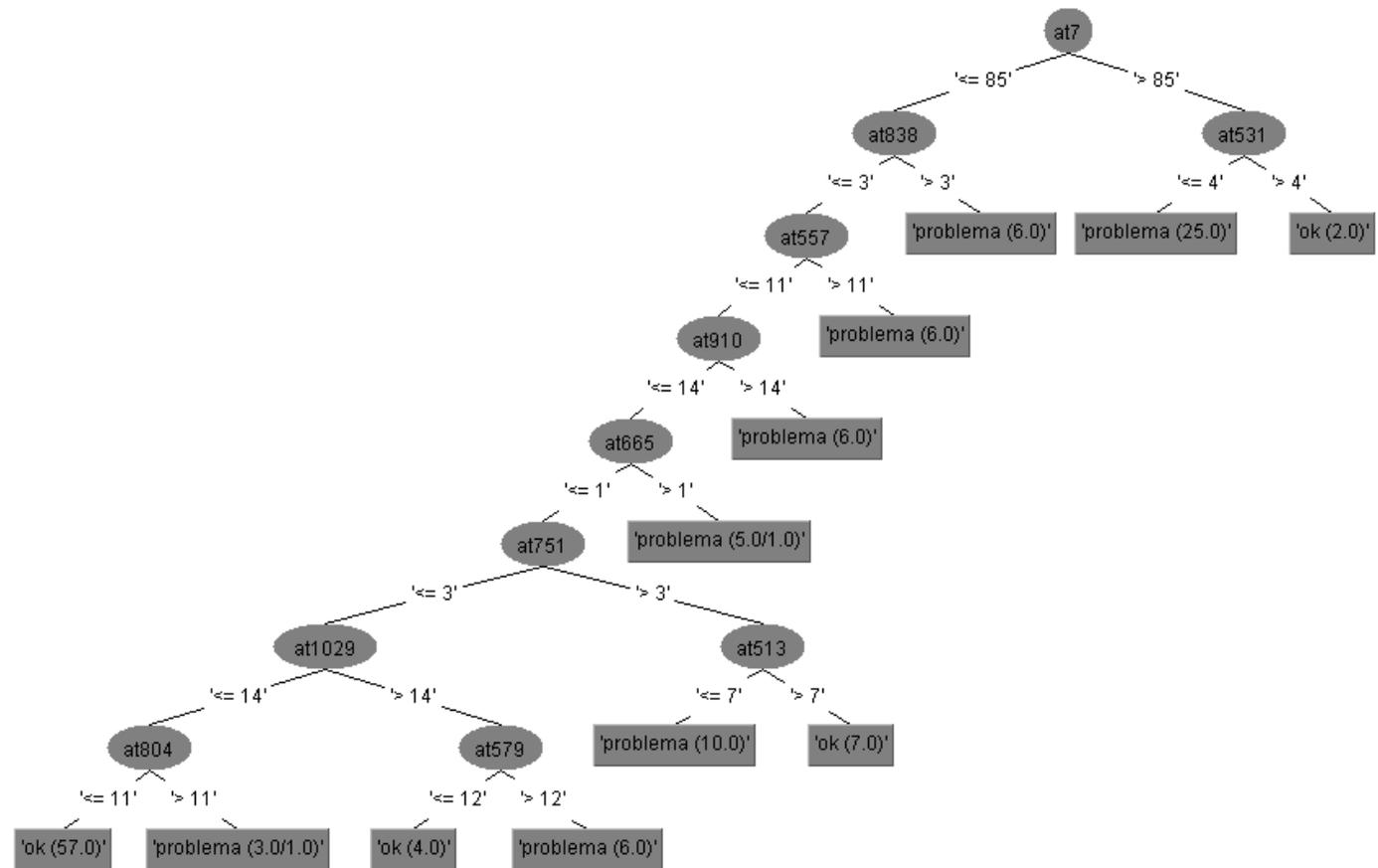


Ilustración 25: Árbol de decisión caso de estudio 1

## 2. CASO DE ESTUDIO 2: CLUSTERING DE TIPOS DE EXÁMENES MÉDICOS

Este segundo caso de estudio se enfoca en el análisis descriptivo, realizando Clustering sobre un conjunto de imágenes que representan distintos tipos de exámenes médicos. Al igual que en el caso de estudio anterior se utiliza la metodología propuesta KDM.

### 2.1 Descripción del caso de estudio

Una gran cantidad de exámenes médicos utilizan imágenes para el diagnóstico y toma de decisiones relacionados con la salud de un paciente. De ahí que este tipo de imágenes sean tan importantes y utilizadas a nivel mundial. Cada tipo de examen presenta unas características propias que lo identifican de los demás tipos de exámenes, relacionadas principalmente con la estructura de la parte del cuerpo a la que se realiza el examen.

Este tipo de imágenes son bastante atractivos para la aplicación de métodos no supervisados como el caso del Clustering. Recientemente ImageCLEF2015, una campaña de evaluación que convoca participantes de todo el mundo para participar en concursos de procesamiento de imágenes, ha incluido dentro de sus actividades una tarea de Clustering en la medicina (Amin & Mohammed, 2015) que consiste en identificar las partes del cuerpo a partir de radiografías.

En este caso de estudio se pretende realizar una actividad similar a la propuesta por ImageClef2015, aplicada a bases de datos de diferentes tipos de exámenes médicos.

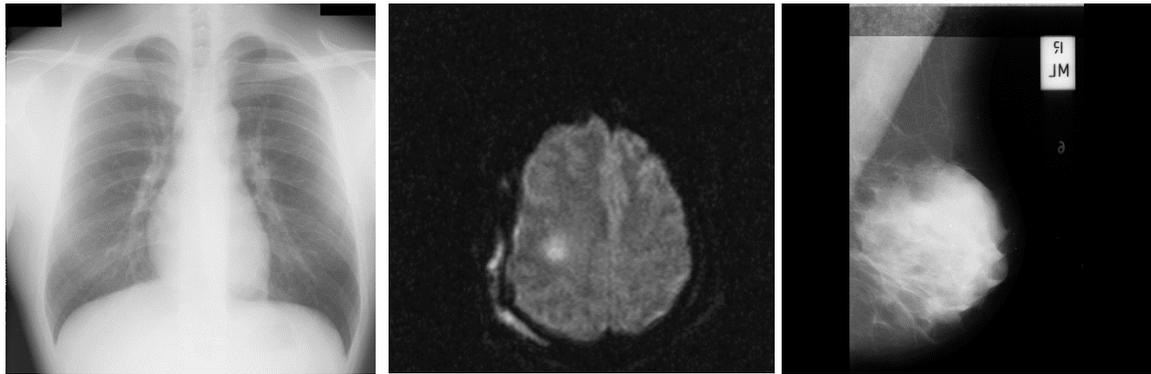
### 2.2 Aplicación de la metodología KDM al caso de estudio

A continuación se presenta una descripción de la aplicación de las fases de la metodología KDD para realizar clustering de exámenes médicos.

#### 2.2.1 Selección de las imágenes

Para este caso de estudio se utilizan imágenes de tres bases de datos distintas. La primera de ellas es conocida como REMBRANDT (REpository for Molecular BRAin Neoplasia DaTa) 900-00-1961, se trata de una serie de imágenes correspondientes a 110 casos de exámenes de resonancia magnética donde se diagnostica la presencia de tumores cerebrales. La segunda base de datos utilizada pertenece a la Sociedad Japonesa de Radiología Tecnológica (JSRT) y tiene como objetivo estudiar la presencia de nódulos pulmonares a través del estudio de radiografías del tórax, mediante el análisis de 154 imágenes donde se diagnostica la presencia de nódulos pulmonares y 93 casos donde no se detectan nódulos pulmonares. La tercer base de datos utilizadas es conocida como MIAS, también utilizada en el caso de estudio 1, tiene como objetivo el análisis de mamografías para detectar la presencia de cáncer de seno.

En la Ilustración 26 se presenta una imagen de cada una de las bases de datos utilizadas. La imagen de la izquierda pertenece a la base de datos JSRT, la del medio pertenece a REMBRANDT y la de la derecha pertenece a MIAS.



**Ilustración 26: Imágenes caso de estudio 2**

Las tres bases de datos mencionadas han sido utilizadas de forma independiente en diversos estudios de procesamiento de imágenes. En este caso de estudio se toman 27 muestras de cada una de ellas para conformar una base de datos de 81 imágenes para realizar clustering de tipos de exámenes médicos.

### **2.2.2 Preprocesamiento de las imágenes**

En esta fase se realizan varias operaciones descritas a continuación.

#### **Conversión de las imágenes**

Las imágenes de la base de datos JSRT se encuentran originalmente en el formato IMG que no es compatible con el prototipo de software, lo mismo sucede con las imágenes de la base de datos REMBRANDT se encuentran en formato DICOM y las imágenes de la base de datos MIAS que están en formato PGM.

Por lo anterior, las 81 imágenes seleccionadas fueron convertidas de su formato original al formato JPG compatible con el prototipo de software. La conversión se realizó con un software externo donde se mantienen las mismas condiciones para todas las imágenes evitando que afecte el proceso de minería.

#### **Segmentación**

A las imágenes extraídas de la base de datos MIAS, se les realizó la segmentación comentada en el caso de estudio 1. Las imágenes de las bases de datos REMBRANDT y JSRT no se consideró necesario realizar segmentación ya que no cuentan con elementos distractores significativos.

## **Cambios en el espacio de color**

Las imágenes son tratadas en los espacios de color rgb y hsv.

### **2.2.3 Indexamiento**

De forma similar al caso de estudio 1, se realizó la extracción de diferentes descriptores de cada una de las características disponibles, color, forma y textura. Con esto se generaron varios modelos de datos para su posterior evaluación. Las características extraídas son:

- Fuzzy Histogram
- Simple Histogram
- LBP (Local Binary Patterns)
- PHOG (Pyramid Histogram of Oriented Gradients)
- RILBP (Rotation Invariant Local Binary Patterns)
- Edge Histogram
- Tamura

Para un total de 1209 descriptores extraídos de cada una de las 81 imágenes de la base de datos construida.

### **2.2.4 Modelamiento**

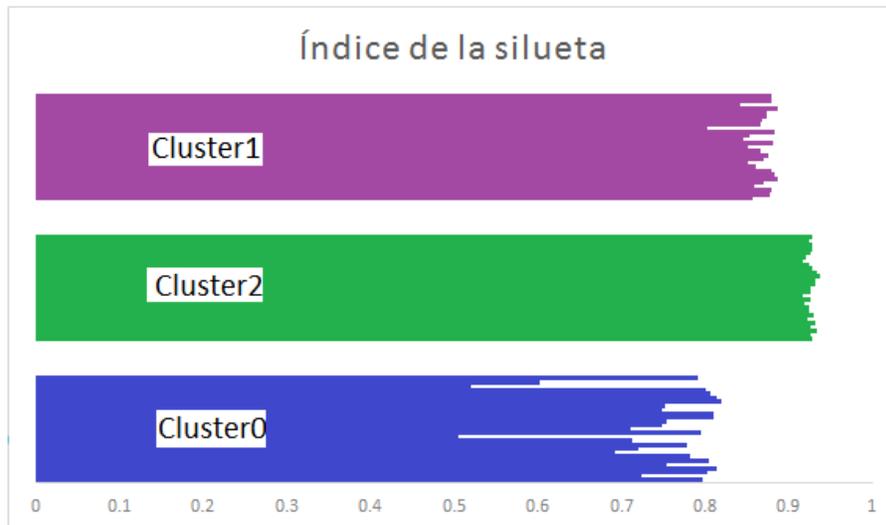
El prototipo de plataforma para minería de imágenes contempla únicamente la técnica K-Means para realizar Clustering. Seleccionando esta técnica se probaron los diferentes modelos de datos generados con el fin de obtener el mejor resultado, que se valida a través del índice de la silueta. Este indicador permite evaluar la efectividad del modelo mediante el análisis de la cohesión de los datos de un mismo cluster y su separación frente a los datos del cluster más cercano.

Adicionalmente se tiene un dato clave para la evaluación del modelo, se conoce con anticipación el número de clusters, que en este caso son tres debido a que se pretende evaluar tres diferentes tipos de imágenes. En algunos casos el número de clusters es desconocido y se puede utilizar el índice la silueta para identificar si se necesita un número de clusters diferente al seleccionado.

### **2.2.5 Evaluación**

En la evaluación de los modelos se observó que varios de ellos realizaron un proceso exitoso al ubicar en el mismo cluster las imágenes del mismo tipo de examen.

En la Ilustración 27 se presenta un gráfico con el índice de la silueta calculado para cada una de las imágenes.



**Ilustración 27: Índice de la silueta caso de estudio 2**

En el caso del cluster 1 se observa que todas las imágenes presentan un índice superior a 0.75, el cluster 2 presenta índices superiores a 0.9, mientras que el cluster 0 presenta índices más bajos pero siempre superiores a 0.5. En promedio el índice de la silueta para este proceso de clustering es de 0.845, un valor que indica una alta probabilidad de éxito al agrupar las imágenes de esta base de datos.

El análisis descriptivo del prototipo de software, cuenta con una función bastante interesante para evaluar la efectividad de un proceso de Clustering. Una vez se han asignado los cluster para cada una de las imágenes, se le permite al usuario almacenar las imágenes separadas en carpetas que indican el cluster asignado, esto con el fin de ser sometidas a un proceso de clasificación donde la clase real de cada imagen corresponde a la carpeta donde se encuentra almacenada. Aplicando esta función se crearon tres carpetas y se sometieron estas imágenes a un proceso de clasificación en el que se obtuvo una efectividad del 100%.

### 3. CONCLUSIONES

Gracias a las investigaciones realizadas para la construcción del estado del arte de este documento, se lograron identificar los requisitos para el desarrollo de un prototipo de software de minería de imágenes de la salud, una herramienta que utiliza la minería de datos no estructurados para apoyar los diagnósticos médicos asistidos por computadora. Estos requisitos permitieron identificar tendencias y necesidades de los servicios de salud en el contexto de las ciudades inteligentes.

El prototipo de software desarrollado permite seleccionar un conjunto de imágenes para aplicar análisis predictivo o descriptivo. El análisis predictivo utiliza técnicas como las máquinas de soporte vectorial y los árboles de decisión para realizar proceso de clasificación. Por su parte, el análisis descriptivo utiliza la técnica K-Means para realizar procesos de Clustering. Estos análisis son posibles gracias al proceso de extracción de características donde se contemplan varios descriptores para extraer información relativa al color, forma y textura presentes en las imágenes. Para el desarrollo de estas actividades se utilizaron librerías como Lire (Lucene Image Retrieval), en la extracción de características y Weka (Waikato Environment for Knowledge Analysis) en la aplicación de la minería.

En el documento de pruebas, presentado en la parte VI de este trabajo, se comprobó el cumplimiento de los requisitos de software establecidos para el desarrollo del prototipo, que están basados en el estándar IEE830 y son presentados en la parte IV. El prototipo cuenta con un sistema de autenticación para el ingreso al sistema, contemplando dos tipos de usuarios: operador y administrador, este último tiene privilegios especiales para creación de nuevos usuarios y configuración de los descriptores y técnicas de minería de datos. Una vez se ingresa al sistema se permite seleccionar la ruta donde se encuentran las imágenes a utilizar y posteriormente se habilitan las opciones para extraer características y aplicar análisis predictivo y descriptivo. En el análisis predictivo se presenta la matriz de confusión como mecanismo de evaluación y se permite almacenar el modelo generado para posteriores usos. En el caso del análisis descriptivo se presenta el índice de la silueta de cada una de las imágenes para medir la efectividad del Clustering, adicionalmente se presenta una opción para utilizar el resultado obtenido en un proceso de clasificación que permite validar la efectividad del agrupamiento.

La precisión de los análisis predictivo y descriptivo fue validada mediante dos casos de estudio. El primero de ellos aplica clasificación sobre una base de datos de 197 mamografías con el fin de identificar la presencia de anomalías que puedan ser consideradas como indicios de cáncer de seno, en este caso se obtuvo un 86% de efectividad extrayendo características de color, forma y textura y aplicando los árboles de decisión como técnica de minería. En el segundo caso de estudio

se aplicó Clustering para realizar la agrupación de 81 imágenes pertenecientes a tres tipos diferentes de exámenes médicos: radiografías del tórax, resonancia magnética y mamografías. Como resultado se logró una agrupación con un 100% de efectividad que presenta un índice de la silueta de 0.85.

## 4. TRABAJO FUTURO

En el desarrollo del presente trabajo se identifican varias líneas en la cuales se puede realizar aportes en el futuro.

### **Ciudades inteligentes**

En la investigación realizada sobre las ciudades inteligentes se identificaron una serie de tendencias en cuanto a tecnologías o conceptos que no se contemplaron en este trabajo. Es el caso del Internet de las cosas (IoT), la computación en la nube y las aplicaciones móviles. En cuanto al Internet de las cosas se puede emplear mecanismos para obtener y compartir las imágenes médicas de forma automática. La aplicación de los conceptos de la computación en la nube como el Software como Servicio (SaaS) y la infraestructura como servicio (IaaS) pueden mejorar significativamente el prototipo de plataforma. Así mismo desde el punto de vista de las aplicaciones móviles se puede facilitar el acceso a la plataforma y permitir la interacción de los usuarios compartiendo imágenes u opiniones.

### **Procesamiento de imágenes**

El procesamiento de imágenes incluye algunas actividades que no se han contemplado dentro de la plataforma, este es el caso de la segmentación de imágenes y remuestreo de imágenes para que tengan el mismo tamaño. La inclusión de estas actividades complementa la plataforma y permite que el proceso sea más automático.

### **Tratamiento de los datos**

Incluir una etapa de transformación de los datos donde se permita realizar actividades como la normalización de los mismos, puede mejorar el desempeño de las técnicas de minería.

### **Minería de datos**

La inclusión de más técnicas de minería de datos para clasificación y Clustering, pueden enriquecer la plataforma en el sentido que presentan un mayor número de alternativas para encontrar un modelo que genere los resultados esperados. Del mismo modo contemplar otros tipos de datos adicionales como texto, audio, y datos estructurados, complementan la plataforma y brindan nuevas alternativas para extraer conocimiento en el área de la salud o incluso en cualquier otra área.

## BIBLIOGRAFÍA

- Achaerandio, R., Curto, J., Bigliani, R., & Gallotti, G. (2012). Análisis de las ciudades inteligentes en España 2012 - El viaje a la ciudad inteligente. *IDC España - Analyze the future*.
- Amaya B., Y. D. (2013). Metodologías ágiles en el desarrollo de aplicaciones para dispositivos móviles. Estado actual. *Revista de Tecnología*, 111-124.
- Amin, A., Takib, R., Raza, S., & Javed, S. (2014). Extract association rules to minimize the effects of dengue by using a text mining technique. 3(4).
- Amin, M. A., & Mohammed, M. K. (2015). *Overview of the ImageCLEF 2015 medical clustering task*. Toulouse, France: CLEF2015 Working Notes. CEUR Workshop Proceedings, CEURWS. org, .
- Apté, C., & Weiss, S. (1997). Data mining with decision trees and decision rules. *Future Generation Computer Systems*, 197-210.
- Azevedo, A., & Rojão, L. (2008). *KDD, SEMMA and CRISP-DM: a parallel overview*.
- Azoumana, K. (2013). Análisis de la deserción estudiantil en la Universidad Simón Bolívar, facultad Ingeniería de Sistemas, con técnicas de minería de datos. *Pensamiento Americano*, 41-51.
- Balu, R., & Devi, T. (2012). Design and Development of Automatic Appendicitis Detection System using Sonographic Image Mining. *International Journal of Engineering and Innovative Technology (IJEIT)*, 67-74.
- Banco Mundial. (22 de 7 de 2015). *Banco mundial*. Obtenido de <http://www.bancomundial.org/temas/cities/datos.htm>
- Beyer, D., Dresler, G., & Wendler, P. (2014). Software Verification in the Google App-Engine Cloud. *Computer Aided Verification*, 327-333.
- Bidgood, W. D., Horii, S. C., Prior, F. W., & Syckle, D. E. (1997). Understanding and using DICOM, the data interchange standard for biomedical imaging. *Journal of the American Medical Informatics Association*, 199-212.
- Cabrera, N., Castro, P., Demeneghi, V., Luque, L., Morales, J., Sainz, L., & Ortiz, M. (2014). mSalUV: un nuevo sistema de mensajería móvil para el control de la diabetes en México. *Rev Panam Salud Publica*, 371-377.
- Calderón, A., Dámaris, S., Rebaza, V., & Carlos, J. (2007). *Metodologías Ágiles*. Trujillo: Escuela de Informatica. Trujillo: Universidad Nacional de Trujillo.

- Chen, C.-C. (2014). The Trend towards "Smart Cities". *International Journal of Automation and Smart Technology*, 4(2), 63-66.
- Chidambaranathan, S. (2015). Image Mining with CBIR. *International Journal of Computer Science and Mobile Computing*, 12-15.
- Choubassi, M., Nefian, A., Kozintsev, I., Bouguet, J., & Wu, Y. (2007). Web image clustering. *Acoustics, Speech and Signal Processing*.
- Corrales, D., Ledesma, A., Peña, A., Hoyos, J., Figueroa, A., & Corrales, J. (2014). A new dataset for coffee rust detection in Colombian crops base on classifiers. *Revista S&T*, 9-23.
- Dávila, F., & Sánchez, Y. (2012). Técnicas de minería de datos aplicadas al diagnóstico de enfermedades clínicas. *Revista Cubana de Informática Médica*.
- Devi, B., Rao, K., Setty, S., & Rao, M. (2013). Disaster Prediction System Using IBM SPSS Data Mining Tool. *International Journal of Engineering Trends and Technology (IJETT)*, 3352-3357.
- Dua, J. (2014). Survey on Image Mining Using Genetic Algorithm. *International Journal of Computer, Electrical, Automation, Control and Information Engineering*, 1822-1827.
- Emam, A. (2014). Intelligent drowsy eye detection using image mining.
- Fernández, J., Miranda, N., Guerrero, R., & Piccoli, F. (2006). Minería de Imágenes. *VIII Workshop de Investigadores en Ciencias de la Computación*.
- Filippone, M., Camastra, F., Masulli, F., & Rovetta, S. (2008). A survey of kernel and spectral methods for clustering. *Pattern recognition*, 41(1), 176-190.
- Filipuzzi, M., Rodrigo, F., Graffigna, J., Isoardi, R., & Noceti, M. (2012). Diseño e implementación de una Base de Datos de Esclerosis Múltiple e Interfaz Gráfica de Usuario. *3º Congreso Argentino de Informática y Salud, CAIS 2012*, 92-103.
- François, O., Ancelet, S., & Guillot, G. (2006). Bayesian clustering using hidden Markov random fields in spatial population genetics. *Genetics*, 174(2), 805-816.
- García, M., Terrones, J., Henarejos, E., & Martínez, M. (2014). Segmentación automática de imágenes de cultivos: estudio comparativo de modelos de color. *I Symposium Nacional de Ingeniería Hortícola*.
- Gómez, J., Huete, J., Hoyos, O., Perez, L., & Grigori, D. (2013). Interaction System Based on Internet of Things as Support for Education. *Procedia Computer Science* 21, 132-139.

- Goodrum, A. (2000). Image Information Retrieval: An Overview of Current Research. *Informing Science*, 63-67.
- Hall, M., Frank, E., Holmes, G., Pfahringer, B., Reutemann, P., & Witten, I. (2009). The WEKA Data Mining Software: An Update. *SIGKDD Explorations*, 10-18.
- Han, J., & Kamber, M. (2006). *Data Mining Concepts and Techniques*. Morgan Kaufmann.
- Hernández, H. (2014). Aplicación de minería de datos a información de pacientes pre-diabéticos. *Revista Iberoamericana de Producción Académica y Gestión Educativa*.
- Hofmann, M., & Ralf, K. (2013). *RapidMiner: Data Mining Use Cases and Business Analytics Applications*. CRC Press.
- Hsu, W., Mong, L., & Ji, Z. (2002). Image Mining: Trends and Developments. *Journal of Intelligent Information Systems*, 19(1), 7-23.
- Ion, A., & Udristoiu, S. (2015). IMAGE MINING FOR ESTABLISHING MEDICAL DIAGNOSIS. *Information Technology And Control*, 123-129.
- Jin, J., Gubbi, J., Marusic, S., & Palaniswami, M. (2014). An Information Framework of Creating a Smart City through Internet of Things. *IEEE Internet of Things Journal*, 112-121.
- Joachims, T., Hofmann, T., Yue, Y., & Yu, C.-N. (2009). Predicting structured objects with support vector machines. *Communications of the ACM*, 52(11), 97-104.
- Kamel, M., & Al-Shorbaji, N. (2014). On the Internet of Things, smart cities and the WHO Healthy Cities. *International Journal of Health Geographics*.
- Kumari, M., & Sunila, G. (2011). Comparative Study of Data Mining Classification Methods in Cardiovascular Disease Prediction. 2.
- Langohr, K. (2010). *Introducción a R*. Barcelona: Departament de Estadística e Investigació Operativa.
- Latorre, A. (2014). Design and implementation of new services for smart cities in Android. *Universitat Politècnica de Catalunya*.
- Leong, L., Toombs, D., Gill, B., Petri, G., & Haynes, T. (01 de 03 de 2015). *Magic Quadrant for Cloud Infrastructure as a Service*. Obtenido de infomall:  
<http://www.infomall.org/I590ABDSSoftware/Resources/Magic%20Quadrant%20for%20Cloud%20Infrastructure%20as%20a%20Service.pdf>
- Lorca, G., Arzola, J., & Pereira, O. (2010). Segmentación de Imágenes Médicas Digitales mediante Técnicas de Clustering. *Rev. Aporte Santiaguino*, 108-116.

- Lovrić, M., Milanović, M., & Stamenković, M. (2014). Algorithmic methods for segmentation of time series: an overview. *1*(1).
- Lux, M., & Chatzichristofis, S. A. (2008). LIRe: Lucene Image Retrieval - An Extensible Java CBIR Library. *Proceedings of the 16th ACM international conference on Multimedia*, 1085-1088.
- Maldonado, J. (2008). *Estudio de métodos de indexación y recuperación en bases de datos de imágenes*. San Sebastiani: Universidad del país Vasco.
- Martínez, J. (2013). La recuperación automatizada de imágenes: retos y soluciones. *Revista General de Información y Documentación* , 423-436.
- Meireles, M., Almeida, P., & Godoy, M. (2003). A comprehensive review for industrial applicability of artificial neural networks. *50*(3).
- Mena, J. (1999). *Data mining your website*. Digital Press.
- Mendes C., K., Estevez, E., & Fillottrani, P. (2010). Evaluación de Metodologías Ágiles para Desarrollo de Software. *WICC 2010 - XII Workshop de Investigadores en Ciencias de la Computación*, 455-459.
- Moine, J. (2013). *Metodologías para el descubrimiento de conocimiento en bases de datos: un estudio comparativo*. Universidad Nacional de la Plata.
- Namiot, D., & Schneps-Schneppe, M. (2013). Smart Cities Software from the developer's point of view. *arXiv preprint arXiv:1303.7115*.
- Ortiz, J. (2013). *Análisis y estudio de viabilidad en la implementación del servicio de imágenes diagnósticas bajo la modalidad de negocio Joint Venture, en el dispensario médico Héroes del Sumapaz*. 2013: Universidad Militar Nueva Granada.
- Palomino, N., & Garcia, N. (2011). Implementación del algoritmo Watershed para el análisis de imágenes médicas. *Revista de Investigación de Sistemas e Informática*, 67-74.
- Pehlivanli, A. (2011). *The comparision of data mining tools*.
- Peláez, G. (2014). *Métodos para la Clasificación Automática de Imágenes de Resonancia Magnética del Cerebro*. Madrid: Escuela Politécnica Superior - Instituto de Investigaciones Biomédicas - Universidad de Madrid.
- Pérez, A. (2012). Aplicación de la red de probabilidad neuronal y escala de framingham para predicción de la hipertension arterial. *Memorias Convención Internacional de Salud Pública*. La Habana.

- Perozo, A., & Boscán, N. (2014). SOFTWARE COMO SERVICIO (SAAS): TENDENCIAS MUNDIALES. *Revista Electrónica Facultad de Ingeniería UVM*, 1127-1138.
- Petrolo, R., Loscri, V., & Mitton, N. (2014). Towards a Smart City based on Cloud of Things. *WiMobCity - International ACM MobiHoc Workshop on Wireless and Mobile*.
- Petrovic, S. (2006). A Comparison Between the Silhouette Index and the Davies-Bouldin Index in Labelling IDS Clusters. *Proceedings of the 11th Nordic Workshop of Secure IT Systems*, 53-64.
- Pyle, D. (2003). *Business modeling and data mining*. Morgan Kaufmann.
- Riquelme, J., Ruiz, R., & Gilbert, K. (2006). Minería de datos: Conceptos y tendencias. *10(29)*, 11-18.
- Rodríguez, C., & Gil, S. (2014). Ciudades amigables con la edad, accesibles e inteligentes. *CEPAT-IMSERSO*.
- Romero, A. (2011). *Nuevos paradigmas para el análisis estadístico de imágenes tomográficas cerebrales*. Granada: Editorial de la Universidad de Granada.
- Santamaría, C. (2014). *Construcción de una base de datos de imágenes de mamografía para la identificación de microcalcificaciones*. Pereira: Universidad tecnológica de Pereira.
- Santana, P., Costaguta, R., & Missio, D. (2014). Aplicación de Algoritmos de Clasificación de Minería de Textos para el Reconocimiento de Habilidades de E-tutores Colaborativos. *Revista Iberoamericana de Inteligencia Artificial*, 57-67.
- Sebastiani, F. (2005). Text Categorization. *Text Mining and its Applications to Intelligence*, 109-129.
- Serrano, N., Gallardo, G., & Hernantes, J. (2015). Infrastructure as a Service and Cloud Technologies. *IEEE SOFTWARE*, 30-36.
- Shaban, M., & Khalaf, A. (2013). Image Information Retrieval Using Wavelet and Curvelet Transform. *International Journal of Soft Computing and Engineering (IJSCE)*, 86-90.
- Shaikh, T. (2014). A Prototype of Parkinson's and Primary Tumor Diseases Prediction Using Data Mining Techniques. *3(4)*.
- Singh, B., Thoke, A., Verma, K., & Chandraka, A. (2011). Image information retrieval from incomplete queries using color and shape features. *Signal & Image Processing : An International Journal (SIPIJ)*, 213-220.
- Slimani, T., & Amor, L. (s.f.). Efficient Analysis of Pattern and Association Rule Mining Approaches. *6(3)*, 70-81.

- Solarte, G. R., & Castro, Y. V. (2012). Modelo híbrido para el diagnóstico de enfermedades cardiovasculares basado en inteligencia artificial. *Tecnura*, 35-52.
- Steinley, D. (2006). K-means clustering : A half-century synthesis. *British Journal of Mathematical and Statistical Psychology*, 59(1), 1-34.
- Suckling, J., Parker, J., Dance, D. R., Astley, S., Hutt, I., Boggis, C., & Savage, J. (1994). The mammographic image analysis society digital mammogram database. *In Exerpta Medica. International Congress Series*, 375-378.
- Suganthira, S., Thamilselvan, P., Sathiaselvan, J., & Lakshmiprabha, M. (2015). A Technical Study on Biomedical image Classification using Mining Algorithms. *National Conference on Recent Advancements in Software Development* .
- Tapia, M., Ruiz, O., & Chirinos, C. (2014). Modelo de clasificación de opiniones subjetivas en redes sociales. *Ingeniería: Ciencia, Tecnología e Innovación*.
- Thamilselvan, P., & Sathiaselvan, J. (2015). A Comparative Study of Data Mining Algorithms for Image Classification. *I.J. Education and Management Engineering*, 1-9.
- Timarán, R., & Yépez, M. (2012). La minería de datos aplicada al descubrimiento de patrones de supervivencia en mujeres con cáncer invasivo de cuello uterino. *Universidad y salud*, 117-129.
- Torres, D. (2013). *Diseño y aplicación de una metodología para análisis de noticias policiales utilizando minería de textos*. Universidad de Chile.
- UIT. (2014). *Una visión general de las ciudades inteligentes sostenibles y el papel de las tecnologías de la información y comunicación*. UIT-T Grupo Temático sobre Ciudades Inteligentes Sostenibles.
- Usama, F., Piatetsky-Shapiro, G., Padhraic, S., & Uthurusamy, R. (1996). Advances in knowledge discovery and data mining.
- Valencia, J., Cruz, J., Caicedo, L., & Chamorro, C. (2014). Extracción de características del iris como mecanismo de identificación biométrica. *Revista Virtual Universidad Católica del Norte*, 182-196.
- Villena, J., Villanueva, S., & Serrano, S. (2013). *Modelado predictivo de la contaminación en la ciudad sostenible*. CIUDAD2020: HACIA UN NUEVO MODELO DE CIUDAD INTELIGENTE SOSTENIBLE.

- Vivas, H., Britos, P., García, N., & Cambarieri, M. (2013). Investigación en Progreso: Estudio y Evaluación de Tecnologías de la Información y la Comunicación para el Desarrollo de Ciudades Inteligentes. *Revista Latinoamericana de Ingeniería de Software*, 146-151.
- Wettig, H., Grünwald, P., Roos, T., Myllymäki, P., & Tirri, H. (2002). *On supervised learning of Bayesian network parameters*. Helsinki Institute for Information Technology (HIIT).
- Wiener, E., Jan, P., & Weigend, A. (1995). A Neural Network Approach to Topic Spotting. *4th Annual Symposium on Document Analysis and Information Retrieval*, (págs. 317-332). Las Vegas.
- Xu, R., & Wunsch, D. (2005). Survey of clustering algorithms. *16*(3).
- Yang, Y., & Liu, X. (1999). A re-examination of text categorization methods. *Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval*, (págs. 42-49).
- Zanella, A., Bui, N., Castellani, A., Vangelista, L., & Zorzi, M. (2014). Internet of Things for Smart Cities. *IEEE Internet of Things Journal*.

## **ANEXOS**

### **1. Manual de usuario del prototipo de software**

# PROTOTIPO DE PLATAFORMA DE MINERÍA DE IMÁGENES PARA EL SECTOR SALUD

MANUAL DE USUARIO

Ver: 1.0

## INTRODUCCIÓN

Las imágenes diagnósticas juegan un papel fundamental para diagnosticar el estado de salud de un paciente, en este tipo de imágenes se encuentra información no trivial, que requiere de técnicas especializadas para ser extraída, procesada y analizada, con el fin de generar nuevo conocimiento. La minería de imágenes puede realizar estas actividades, a partir de un conjunto de imágenes se pueden realizar análisis de tipo predictivo o descriptivo con el fin de extraer la información presente en dichas imágenes.

En este documento se presenta el manual de usuario de un prototipo de plataforma de minería de imágenes del sector salud, que permite realizar análisis de tipo predictivo y descriptivo.

## REQUISITOS

El prototipo se ha desarrollado en el lenguaje de programación Java. Para su funcionamiento se requiere la instalación de Java Runtime Environment (JRE) que puede obtenerse de: <http://www.java.com/es/>

Para almacenar la información se utiliza la base de datos embebida de Apache Derby que no requiere de la instalación de un motor de base de datos.

Para la correcta operación del prototipo de software, el usuario debe proporcionar imágenes de alta calidad donde se presente únicamente información de interés, si el usuario considera necesario un proceso de limpieza de imágenes, este se debe realizar de manera externa antes de ingresar las imágenes al prototipo de software. De igual manera se espera que las imágenes tengan la misma resolución y el mismo tamaño, si estas condiciones no se cumplen, se debe realizar etapas de preprocesamiento de forma externa, en estas etapas se puede cambiar la resolución, redimensionar las imágenes y aplicar cualquier otro tipo de intervención que se considere prudente con el fin de que todas las imágenes, a ingresar al prototipo de software, presenten características similares.

El prototipo de software solo recibe imágenes con alguna de estas extensiones: JPG, BMP, PNG. Si desea utilizar imágenes en otros formatos de archivo, se debe hacer la respectiva conversión de forma externa.

## DESCRIPCIÓN GENERAL

Para utilizar el prototipo de plataforma, se deben realizar las siguientes actividades: login, cargar imágenes, configurar técnicas y descriptores, extraer características y, finalmente, aplicar análisis de tipo predictivo o descriptivo.

A continuación se detallan cada una de estas actividades.

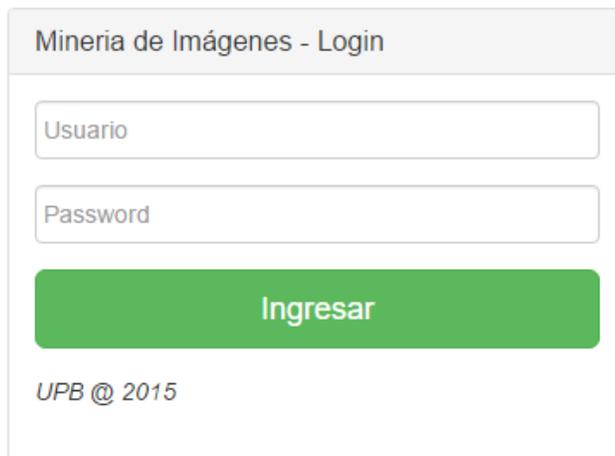
## LOGIN

Para el control de acceso se utiliza un sistema de autenticación de usuario y contraseña. Cuando se ejecuta por primera vez, el sistema crea la base de datos con dos usuarios, que tienen la siguiente información:

Id	Nombre	Usuario	Password	Tipo
1	Administrador por defecto	admin	admin	Administrador
2	Operador por defecto	operador	operador	Operador

Los usuarios de tipo administrador tienen acceso a dos funcionalidades especiales, a saber, administrar usuarios y configuraciones. En la administración de usuarios se permite crear, editar o eliminar usuarios, en configuraciones se pueden habilitar/deshabilitar los descriptores para la extracción de características y asignar parámetros para la aplicación de las técnicas de minería de datos.

En la siguiente figura se presenta la pantalla de login.



Mineria de Imágenes - Login

Usuario

Password

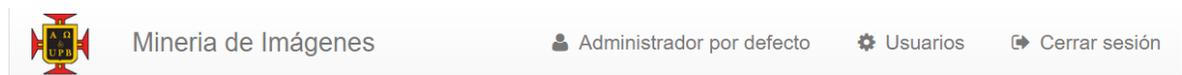
Ingresar

UPB @ 2015

Para acceder al sistema, ingrese usuario y Password, y presione el botón ingresar. Si los datos ingresados son correctos se presentará la pantalla inicial del sistema. En caso de que alguno de los datos o ambos sean incorrectos el sistema permanece en la pantalla de login y presenta un mensaje de error de acceso al sistema.

## ACCESO A LAS FUNCIONES

Para acceder a las distintas funcionalidades, las pantallas del sistema cuentan con una barra superior y un menú. En la siguiente imagen se presenta la barra superior.

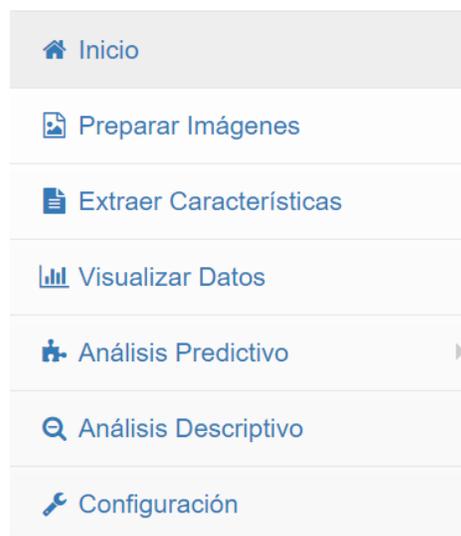


En este caso se ha iniciado sesión con el usuario admin cuyo nombre es Administrador por defecto. El nombre del usuario es presentado como un botón en la barra superior que permite editar la configuración del usuario.

Se cuenta también un botón de usuarios (disponible únicamente para administradores) que permite realizar la administración de los usuarios del sistema.

Por último se cuenta con un botón para cerrar la sesión y volver a la pantalla de login.

En la siguiente figura se presenta el menú que permite el acceso a las principales funciones del sistema



Este menú permite el acceso a las funciones relacionadas con la aplicación de la minería.

## CONFIGURACIÓN

Mediante esta opción, los usuarios de tipo administrador, pueden crear, editar o eliminar registros de configuración del sistema. Para la configuración de las técnicas de minería de datos, se utiliza la nomenclatura propuesta por Weka.

Cada registro de configuración contiene la siguiente información:

- **Análisis Predictivo:**

---

**Análisis Predictivo**

NumMinClases:	<input type="text" value="2"/>
NumMaxClases:	<input type="text" value="5"/>
Porcentaje Entrenamiento:	<input type="text" value="90%"/>
conf. Árboles:	<input type="text"/>
conf. MSV:	<input type="text"/>

**NumMinClases:** Indica cual es el mínimo número de clases para realizar un proceso de clasificación.

**NumMaxClases:** Indica cual es el número máximo de clases para realizar un proceso de clasificación.

**Porcentaje Entrenamiento:** Indica el porcentaje de datos que se utilizarán para el entrenamiento, los datos restantes se utilizarán para evaluar el modelo de clasificación creado en el entrenamiento.

**Conf. Árboles:** Son los parámetros de configuración para la técnica de minería da datos árboles de decisión. Estos parámetros se deben ingresar separados por coma “,”. A continuación se describen los distintos parámetros que se pueden utilizar:

Parámetro	Valor	Descripción
-B		Habilitar binary Splits
-C,X.YZ	X.YZ	Confidence factor
-M,X	X	minNumObj
-S		Subtree Raising
-U		Unpruned
-A		Use Laplace

**ConfMSV:** Son los parámetros de configuración para la técnica de minería da datos máquinas de soporte vectorial. Estos parámetros se deben ingresar separados por coma “,”. A continuación se describen los distintos parámetros que se pueden utilizar:

Parámetro	Valor	Descripción
-S,X	X	Type. X=0. C-SVC Classification X=1. Un-SVC Classification X=2. One-class Classification
-R,X	X	Coef0
-C,X	X	cost
-D,X	X	Degree of the Kernel
-E,X.YYY	X.YYY	Eps (Termination of the tolerance criterion)
-G,X.YYY	X.YYY	gamma
-K,X	X	Type of Kernel X=0. Linear: $u*v$ X=1. Polinomial( $\text{gamma} * u*v + \text{coef0}$ ) <sup>degree</sup> X=3. Radial basis function: $\exp(-\text{gamma} *  u-v ^2)$ X=3. Sigmoid: $\tanh(\text{gamma}*u*v + \text{coef0})$
-Z		normalize
-N	X.Y	nu

- **Análisis Descriptivo**

**Análisis Descriptivo**

<b>NumMinCluster:</b>	<input type="text" value="2"/>
<b>NumMaxCluster:</b>	<input type="text" value="5"/>
<b>conf. K-Means:</b>	<input type="text"/>

**NumMinCluster:** Número mínimo de clusters

**NumMaxCluster:** Número máximo de clusters

**Conf K-Means:** Son los parámetros de configuración para la técnica de minería da datos K-Means. Estos parámetros se deben ingresar separados por coma “,”. A continuación se describen los distintos parámetros que se pueden utilizar:

Parámetro	Valor	Descripción
-A,"s"	s	Distance function s=weka.core.EuclideanDistance –R first-last s=weka.core.ManhattanDistance –R first-last
-I,X	X	Max Iterations
-O	X	Preserve Instances order

- **Descriptores de las características a extraer**

Se presentan los diferentes descriptores que se pueden extraer de cada característica. Haciendo clic sobre cada descriptor se puede habilitar o deshabilitar

#### Características de Color

AutoCorrelogram	<input type="button" value="NO"/>
Layout	<input checked="" type="button" value="SI"/>
Scalable	<input type="button" value="NO"/>
Simple Histogram	<input type="button" value="NO"/>
Fuzzy Histogram	<input type="button" value="NO"/>

---

#### Características de Textura

Edge Histogram	<input checked="" type="button" value="SI"/>
Tamura	<input checked="" type="button" value="SI"/>

---

#### Características de Forma

PHOG	<input checked="" type="button" value="SI"/>
------	--

---

### CARGAR IMÁGENES

Mediante esta función, el usuario le indica al sistema la ruta donde se encuentran almacenadas las imágenes a las que se les desea aplicar la minería. Para cargar las imágenes entre en la opción Preparar imágenes y presione el botón Cargar directorio, al hacer esto aparecerá una ventana de selección donde debe localizar el directorio de las imágenes. El sistema buscará en este directorio todos los archivos con formato: JPG, PNG o BMP, y las listará mostrando la ruta y una opción que permite visualizar cada imagen.

En la siguiente figura se presenta la pantalla para la carga de imágenes.

# Preparar Imágenes

Seleccione la ruta donde se encuentran las imágenes



Ruta seleccionada: D:\Image Mining\Clustering2

Numero de imagenes: 81

[Reiniciar](#)

Id	Ruta	Opciones
1	D:\Image Mining\Clustering2\Mamografias\mdb007.jpg	
2	D:\Image Mining\Clustering2\Mamografias\mdb008.jpg	
3	D:\Image Mining\Clustering2\Mamografias\mdb009.jpg	
4	D:\Image Mining\Clustering2\Mamografias\mdb011.jpg	
5	D:\Image Mining\Clustering2\Mamografias\mdb014.jpg	
6	D:\Image Mining\Clustering2\Mamografias\mdb016.jpg	

## EXTRAER CARACTERÍSTICAS

En esta función se deben seleccionar el tipo de características que se desea extraer. Si desea seleccionar más de un tipo de características se debe utilizar las teclas Ctrl o Shift. Luego de seleccionar características, presione el botón extraer y espere a que la barra de progreso llegue al 100% indicando de que se han extraído las características de todas las imágenes.

Al finalizar se le indicará el número de imágenes extraídas y el número de descriptores obtenido por cada una de las imágenes. Luego de realizar este proceso se habilitan las opciones para aplicar los análisis predictivo y descriptivo.

En la siguiente figura se puede observar el proceso de extracción de características.

# Extraer Características

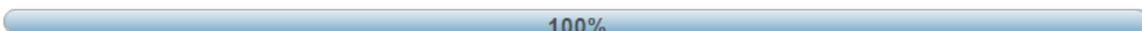
Seleccione el tipo de características que desea extraer

Características de Color  
Características de Forma  
Características de Textura

Opciones

Extraer  
Reiniciar

Progreso



Numero de imagenes: 81

Numero de descriptores: 33

Características extraidas

Color				
Id	Nombre	Descriptor Inicial	Descriptor Final	Total
1	Layout	0	32	33

## ANÁLISIS PREDICTIVO

El análisis predictivo consta de dos funciones, a saber, Entrenamiento y Aplicar Modelo. En la siguiente figura se observa como seleccionar la opción deseada



Para el entrenamiento se debe seleccionar el número de clases y la técnica que se desea utilizar. En caso de que el número de clases coincida con el número de subcarpetas del directorio seleccionado para cargar las imágenes, el sistema asume el nombre de cada subcarpeta como

label de clase y le asigna por defecto a cada imagen la clase real correspondiente a la subcarpeta donde está almacenada. Esta opción se utiliza para facilitar el hecho de tener que seleccionar la clase real de cada una de las imágenes, cuando se trate de muchas imágenes esta puede ser una labor demandante.

En caso de que el número de clases seleccionadas y el número de subcarpetas no coincida, el sistema asigna por defecto a todas las imágenes, la clase 1 como clase real.

En la siguiente figura se presenta la pantalla de entrenamiento para análisis predictivo.

## Análisis Predictivo - Entrenamiento

**Seleccione el número de clases**

**Seleccione técnica a Utilizar**

Máquinas de soporte Vectorial

Árboles de decisión

---

**Identificación de las clases**

Id	Label
1	<input type="text" value="Mamografías"/>
2	<input type="text" value="Resonancia Magnética"/>
3	<input type="text" value="Torax"/>

**Opciones**

Aplicar

Reiniciar

---

**Indique a cual clase pertenece cada imagen**

Id	Ruta	Clase Real	Clase Pred	Opciones
1	D:\Image Mining\Clustering2\Mamografías\mdb007.jpg	<input type="text" value="Mamografías"/>	NULL	
2	D:\Image Mining\Clustering2\Mamografías\mdb008.jpg	<input type="text" value="Mamografías"/>		
3	D:\Image Mining\Clustering2\Mamografías\mdb009.jpg	<input type="text" value="Torax"/>		
4	D:\Image Mining\Clustering2\Mamografías\mdb014.jpg	<input type="text" value="Mamografías"/>	NULL	

Luego de verificar que cada imagen tiene asignada correctamente la clase real, se debe presionar el botón de entrenar. Al hacer esto el sistema toma el porcentaje de entrenamiento asignado en las configuraciones y crea dos grupos de imágenes, el primero de ellos se utiliza para entrenar y el segundo se utiliza para realizar la evaluación. Como mecanismo de evaluación se presenta la matriz de confusión. En la siguiente figura se presenta un ejemplo de ello.

```
=== Detailed Accuracy By Class ===

      TP Rate  FP Rate  Precision  Recall  F-Measure  ROC Area  Class
      1      0      1      1      1      1      Mamografías
      1      0      1      1      1      1      Resonancia
      1      0      1      1      1      1      Torax
Weighted Avg.  1      0      1      1      1      1

=== Confusion Matrix ===

a b c  <-- classified as
2 0 0 | a = Mamografías
0 3 0 | b = Resonancia
0 0 4 | c = Torax
```

Almacenar modelo creado

Nombre:

Nota: El modelo será almacenado en la ruta: \\data\predictivo\modelos

El modelo creado durante la etapa de entrenamiento, puede ser almacenarlo en el equipo para aplicarlo posteriormente sobre un conjunto de imágenes. Para ello debe indicar el nombre que desea asignarle al modelo y presionar el botón Grabar. Se generaran dos archivos, el primero de ellos tiene extensión .model y se encarga de almacenar el modelo creado, el segundo tiene extensión .bin y contiene información relacionada con el tipo de características extraídas y las clases utilizadas en el entrenamiento.

La segunda opción del análisis predictivo consiste en cargar un modelo previamente almacenado en el equipo y aplicarlo sobre un conjunto de imágenes. Al ingresar en la opción Aplicar Modelo, aparecerá un listado con los modelos que han sido almacenados en el equipo, seleccione uno de ellos y presione el botón Cargar Modelo. Una vez se ha cargado el modelo, se debe seleccionar la ruta donde están almacenadas las imágenes a las cuales se les aplicará el modelo, para ello presione el botón Cargar Imágenes y seleccione la ruta.

En ese momento el sistema procede a extraer las características de cada una de las imágenes y aplicarle el modelo, este es un proceso que puede tomar un tiempo apreciable. Al finalizar se presenta una tabla donde se listan las imágenes cargadas indicando su respectiva clasificación. En la siguiente figura se presenta la pantalla de Aplicar Modelo.

# Análisis Predictivo - Aplicar Modelo

Seleccione el modelo que desea cargar

Modelo1.model  
Modelo2.model  
Modelo3.model

Cargar Modelo

Seleccione la ruta donde se encuentran las imágenes

Cargar Imágenes

Ruta Imágenes de Prueba:

D:\Image Mining\Clustering2

Id	Ruta	Clase	Opciones
1	D:\Image Mining\Clustering2\Mamografias\mdb007.jpg	Mamografias	
2	D:\Image Mining\Clustering2\Mamografias\mdb008.jpg	Mamografias	
3	D:\Image Mining\Clustering2\Mamografias\mdb009.jpg	Mamografias	

En esta tabla también cuenta con una opción para visualizar cada imagen.

## ANÁLISIS DESCRIPTIVO

Para aplicar el análisis descriptivo se debe seleccionar el número de clusters y la técnica K-Means. Como resultado se presentan los centroides obtenidos.

Adicionalmente se presenta una tabla donde se enumeran cada una de las imágenes, la ruta, el cluster donde quedó agrupado, una opción para visualizar la imagen y el índice de la silueta calculado. Este es un indicador de la efectividad del proceso de Clustering, se espera que el índice de la silueta sea cercano a 1 para indicar que el agrupamiento de la imagen en cuestión es aceptable. En caso de que el índice de la silueta sea negativo es posible que el agrupamiento no sea el adecuado.

En la siguiente figura se presenta la pantalla de análisis descriptivo

# Análisis Descriptivo

Seleccione el número de clusters

Seleccione técnica a Utilizar

Resultado del clustering

```
kMeans
=====

Number of iterations: 2
Within cluster sum of squared errors: 290.2652137871122
Missing values globally replaced with mean/mode

Cluster centroids:
      Cluster#
Attribute Full Data  0  1  2
      (81) (27) (27) (27)
=====
```

79	D:\Image Mining\Clustering2\Torax\JPCNN034.jpg	Cluster1	0.8504895	
80	D:\Image Mining\Clustering2\Torax\JPCNN035.jpg	Cluster1	0.8995337	
81	D:\Image Mining\Clustering2\Torax\JPCNN036.jpg	Cluster1	0.8900564	

Indice de la silueta promedio: 0.859737

Almacenar imágenes para clasificación

Nombre de la carpeta: \*

Nota: Las imágenes se almacenan en la ruta: ldata\descriptivo\images

Al final se presenta el índice de la silueta promedio de todas las imágenes.

Adicionalmente se dispone de una opción que permite tomar las imágenes usadas para el Clustering y almacenarlas en un nuevo directorio indicado por el usuario, cada imagen se almacena dentro de una subcarpeta que indica el cluster donde fue ubicada. De este modo se puede usar fácilmente estas imágenes dentro de un proceso de clasificación. Para esto indique el nombre del directorio y presione el botón Grabar.