



UNIVERSIDAD DE CHILE  
FACULTAD DE CIENCIAS FÍSICAS Y MATEMÁTICAS  
DEPARTAMENTO DE CIENCIAS DE LA COMPUTACIÓN

REGLAS DE ASOCIACIÓN PARA LÍNEAS MOLECULARES

MEMORIA PARA OPTAR AL TÍTULO DE INGENIERO CIVIL EN COMPUTACIÓN

NICOLÁS MARTÍN MIRANDA CASTILLO

PROFESOR GUÍA:  
GUILLERMO CABRERA VIVES

MIEMBROS DE LA COMISIÓN:

Este trabajo ha sido parcialmente financiado por .

SANTIAGO DE CHILE  
DICIEMBRE 2014



# Resumen

Hoy en día, las observaciones astronómicas se realizan a distintas bandas de frecuencia del espectro electromagnético. Estas tienen como resultado lecturas de la radiación recibida desde diversos objetos del cielo, tales como estrellas, galaxias y otros. A estas se les llama espectros de frecuencia. Cada uno de estos espectros contiene picos de intensidad llamados líneas espectrales. Estas líneas son el resultado de una emisión electromagnética de mayor intensidad en frecuencias específicas por parte de las especies químicas (e.g. átomos, moléculas, isótopos) que componen el objeto de origen. Una misma línea espectral puede estar presente en más de un espectro de frecuencia; y, a su vez, una línea detectada en el espectro puede estar identificada (cuando se conoce la especie que la originó) o puede ser desconocida.

Asumiendo que se cuenta con los conjuntos de líneas detectadas en cada espectro de frecuencias, el objetivo del presente trabajo es obtener información de cómo se relacionan entre sí las líneas espectrales a lo largo de distintos espectros. Para ello se utilizó algoritmos de Aprendizaje de Reglas de asociación, o *Association Rule Learning (ARL)*; que, aplicado a estos datos, entregó reglas que agrupan de forma lógica a los conjuntos de líneas y que traen asociadas distintas medidas de relevancia.

Los algoritmos se probaron sobre datos de observaciones ópticas obtenidas del *Sloan Digital Sky Survey (SDSS)*, previo un pre-procesamiento adecuado de estos, y se espera a futuro poder realizar su aplicación a datos en otras frecuencias del espectro electromagnético, como por ejemplo a los datos radioastronómicos del *Atacama Large Millimeter/submillimeter Array (ALMA)*.



*Una dedicatoria corta. Por ejemplo, A los creadores de U-Campus*

# Agradecimientos

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Nam dui ligula, fringilla a, euismod sodales, sollicitudin vel, wisi. Morbi auctor lorem non justo. Nam lacus libero, pretium at, lobortis vitae, ultricies et, tellus. Donec aliquet, tortor sed accumsan bibendum, erat ligula aliquet magna, vitae ornare odio metus a mi. Morbi ac orci et nisl hendrerit mollis. Suspendisse ut massa. Cras nec ante. Pellentesque a nulla. Cum sociis natoque penatibus et magnis dis parturient montes, nascetur ridiculus mus. Aliquam tincidunt urna. Nulla ullamcorper vestibulum turpis. Pellentesque cursus luctus mauris.

# Índice general

<b>Introducción</b>	<b>1</b>
Contexto y motivación . . . . .	1
Objetivos . . . . .	2
Objetivo General . . . . .	2
Objetivos Específicos . . . . .	2
Descripción de la solución . . . . .	3
 <b>1. Marco Teórico</b>	 <b>5</b>
1.1. Antecedentes Astronómicos . . . . .	5
1.1.1. Espectroscopía Astronómica . . . . .	5
1.1.2. Sloan Digital Sky Survey (SDSS) . . . . .	7
1.1.3. Atacama Large Millimeter/submillimeter Array (ALMA) . . . . .	8
1.2. Reglas de asociación . . . . .	9
1.2.1. Definición formal . . . . .	9
1.2.2. Algoritmos, implementaciones y aplicaciones . . . . .	11
 <b>2. Especificación del Problema</b>	 <b>13</b>
2.1. Descripción del problema . . . . .	13
2.2. Requisitos de la solución y casos de uso . . . . .	14
 <b>3. Descripción de la Solución</b>	 <b>15</b>
3.1. Arquitectura de software . . . . .	15
3.1.1. Paquete de Association Rule Learning (ARL) . . . . .	15
3.1.2. Paquete de procesamiento de datos . . . . .	18
3.2. Diseño de clases . . . . .	19
3.2.1. Clase <i>ItemSet</i> . . . . .	19
3.2.2. Clase <i>AssociationRule</i> . . . . .	20
3.2.3. Clase <i>FrequentItemSetMiner</i> . . . . .	20
3.2.4. Clase <i>RuleMiner</i> . . . . .	20
3.3. Algoritmos y estructuras de datos . . . . .	20
3.3.1. Algoritmo <i>Apriori</i> . . . . .	21
3.3.2. Algoritmo <i>FP-Growth</i> . . . . .	22
3.3.3. Extracción de reglas de asociación . . . . .	22
3.3.4. Desempeño comparativo . . . . .	22
3.4. Interfaz de usuario . . . . .	22

<b>4. Validación de la Solución</b>	<b>23</b>
4.1. Antecedentes de datos de prueba . . . . .	23
4.2. Pre-procesamiento de datos . . . . .	23
4.3. Resultados . . . . .	23
<b>Conclusión</b>	<b>23</b>



# Índice de figuras



# Introducción

*[...] we may in time ascertain the mean temperature of heavenly bodies, but I regard this order of facts as for ever excluded from our recognition. We can never learn their internal constitution [...]*

Auguste Comte, *Astronomy, Ch. I: General View*, 1835

## Contexto y motivación

En los últimos tiempos, y en gran parte debido al explosivo desarrollo tecnológico, han surgido numerosos campos en los cuales se ha requerido el uso de procesamiento masivo de datos e inteligencia computacional con el fin de automatizar y auxiliar el proceso de generación de nuevo conocimiento. La astronomía es, sin lugar a dudas, uno de ellos. Esto se debe, en parte, al explosivo desarrollo de nuevas tecnologías que ponen al alcance de la comunidad científica una cantidad nunca antes vista de datos; los cuales tienen el potencial de contener invaluable información sobre el universo, su composición, estructura, origen y destino.

Un claro ejemplo de esto lo constituye el *Atacama Large Millimeter/sub-millimeter Array (ALMA)*, un interferómetro radio-astronómico que consiste en un arreglo de 66 antenas que observan el espacio en las bandas milimétricas y submilimétricas del espectro electromagnético. Ubicado en el desierto de Atacama, en el norte de Chile, es parte de uno de los proyectos científicos más importantes del último tiempo a nivel nacional; en el cual se ha hecho uso de tecnologías de punta por parte de investigadores, ingenieros y técnicos expertos en computación de alto rendimiento, redes de fibra óptica, Machine Learning, minería de datos, entre otros.

La tecnología involucrada en el proyecto *ALMA* ha permitido, entre otras cosas, obtener datos de alta resolución provenientes de distintas fuentes u objetos del espacio observable desde la tierra, cuya posición en el cielo es cuantificada mediante coordenadas celestes. La radiación electromagnética emitida por estos objetos, en bandas de frecuencia de radio, son captadas por el arreglo de antenas y posteriormente procesadas por equipos de alta capacidad

con el fin de obtener los espectros electromagnéticos correspondientes. Estos, a su vez pueden ser analizados directamente o utilizarse para generar imágenes de alta calidad de regiones del espacio en diversos rangos de frecuencia.

Si bien el caso de *ALMA* es notable por las características particulares de las observaciones en frecuencias de radio y por las oportunidades tecnológicas que involucra, el analizar y extraer información a partir de radiación electromagnética, en diversos rangos de su espectro, es parte primordial de la labor astronómica en todos sus campos.

Parte principal de la importancia de estos espectros de radiación electromagnética es que dan información valiosa sobre la composición química de los objetos de los que esta proviene; ya sean estrellas, galaxias u otros muchos tipos de estructuras celestiales. Esto se debe a que los átomos que componen estos objetos emiten o absorben una mayor cantidad de energía en frecuencias muy específicas, tal y como lo predice la teoría cuántica subyacente. Por lo tanto, un espectro en particular tendrá puntos más altos (picos) en ciertas frecuencias dependiendo de los elementos químicos de los que está compuesto el objeto del que proviene.

La detección de líneas espectrales es un problema de interés en sí, y que puede llegar a ser muy complejo dependiendo de en qué bandas de frecuencia se esté trabajando. Sin embargo, se puede seguir obteniendo información valiosa de los objetos observados a partir de estas líneas ya detectadas. Esto incluye potencialmente respuestas a preguntas como: ¿de qué forma se relacionan ciertos tipos de líneas entre sí? ¿Existe una mayor correlación de presencia de líneas de ciertos isótopos o moléculas en particular? ¿Hay una mayor presencia de líneas de ciertas especies en algunos objetos que en otros? ¿Qué nos dice esto de la composición de los objetos y de su química subyacente?

## Objetivos

Los objetivos del presente trabajo se enumeran a continuación:

### Objetivo General

- Asociar líneas de transición cuántica presentes en espectros de frecuencia (obtenidos a partir de observaciones astronómicas) entre sí; en particular, tanto aquellas líneas de las cuales se sabe actualmente a qué especie (e.g. átomo, ion, molécula, isótopo) corresponden como aquellas que aun no han sido identificadas o asociadas a alguna especie conocida.

### Objetivos Específicos

- Inferir reglas de asociación para líneas moleculares detectadas en espectros dentro de un mismo rango de frecuencias.

- Inferir reglas de asociación para líneas moleculares detectadas en espectros dentro de distintos rangos de frecuencias.
- Evaluación de los algoritmos sobre simulaciones y datos reales.
- Desarrollar una herramienta para asociación de líneas moleculares que funcione en un ambiente de computación de alto rendimiento (HPC) y que sea compatible con plataformas de observatorios virtuales (VO).

## Descripción de la solución

Si bien existen diversas técnicas de clasificación y caracterización de puntos en un espacio multidimensional (en nuestro caso objetos descritos por parámetros), para resolver las preguntas anteriores se requiere más bien de una herramienta que permita encontrar relaciones explícitas entre los parámetros en sí, y que permita asignar medidas de relevancia a estas relaciones.

En el presente trabajo se plantea el uso de *Association Rule Learning (ARL)*, o Aprendizaje de Reglas de Asociación, como una herramienta que puede dar respuesta directa a algunas de las interrogantes mencionadas anteriormente, y ayudar a obtener información clave para el proceso de utilizar otras técnicas en el largo plazo.

El Aprendizaje de Reglas de Asociación como técnica se ubica dentro del área de la minería de base de datos, y su concepción original fue el ser aplicada a sistemas de puntos de venta con el fin de encontrar las relaciones más comunes entre artículos comprados por los clientes. Sin embargo, con el tiempo se ha convertido en una de las herramientas más utilizadas de su área, en una diversa gama de contextos.

En el presente trabajo se llevó a cabo el uso de esta técnica con el fin de encontrar relaciones comunes entre líneas espectrales a través de distintos espectros de frecuencia. Ahora bien, la naturaleza innata de estos es más bien continua y las líneas en sí mismas poseen diversos parámetros que las caracterizan. Por lo tanto, este caso dista mucho de la binaridad del problema original para el cual se pensó ARL. Sin embargo, como se muestra a lo largo de este informe, si se asume que se realizó con anterioridad un buen trabajo de detección de líneas y se efectúa un pre-procesamiento adecuado de los datos, el algoritmo de ARL arroja resultados de interés.

En particular, se utilizó una implementación de dos de los algoritmos más utilizados de ARL: *Apriori* y *FP-Growth*. Luego, se obtuvo una base de datos de líneas espectrales ya detectadas (pero no necesariamente asociadas a alguna especie [átomo, isótopo, etc.]) correspondientes a observaciones del *Sloan Digital Sky Survey (SDSS)*, un sondeo espectroscópico del espacio realizado con un telescopio óptico. Sobre este conjunto de datos se procedió a realizar un pre-procesamiento que, entre otros, consta de filtrar las líneas según su brillo o razón señal a ruido. Luego, se efectuaron particiones según las características de los objetos de procedencia (como tipo de objeto o estructura estelar, cercanía, etc.). Finalmente, sobre estas se procedió a aplicar los algoritmos de ARL.

Los resultados obtenidos fueron efectivamente reglas de asociación entre líneas espectrales que resultaron tener mayor relevancia sobre el conjunto de datos bajo distintas medidas. Quedó para su estudio en trabajos posteriores el crear una forma eficiente e intuitiva para un usuario de recorrer y visualizar estas reglas, filtrar el conjunto inicial por parámetros y encontrar diferencias entre las reglas generadas por distintos conjuntos. Esto con el fin de facilitar aun más el descubrimiento de información valiosa sobre la química y composición de los objetos estudiados. Junto con esto, queda para desarrollo a futuro una implementación más general del procedimiento para así aplicar los algoritmos a datos obtenidos en otras bandas de frecuencia, como por ejemplo, las observaciones radioastronómicas de ALMA; que por sus características, promete un mayor número de datos sobre los cuales obtener reglas de asociación para líneas espectrales.

# Capítulo 1

## Marco Teórico

### 1.1. Antecedentes Astronómicos

#### 1.1.1. Espectroscopía Astronómica

En los comienzos del siglo XIX, los astrónomos comenzaron a realizar medidas que, por primera vez, revelaron con exactitud cuán lejanas se encuentran de la tierra incluso las estrellas más cercanas. Y, al igual que entonces, actualmente sigue siendo técnicamente imposible, con la tecnología que contamos, viajar a las cercanías de estos objetos estelares. Sin embargo, hoy en día es bastante bien conocida la composición química de las estrellas y del material difuso presente en los vastos espacios que las separan. El estudio de los espectros de los objetos astronómicos, o *espectroscopía astronómica*, es lo que ha hecho esto posible.

En el año 1814, el científico Joseph von Fraunhofer (1787 - 1826), mediante el uso de prismas de alta calidad contruidos por él mismo, logró difractar un rayo de luz solar y proyectarlo hacia un muro blanco. Además de los colores característicos del arcoíris, observados de esta manera desde los tiempos de Newton, vio en la proyección resultante muchas líneas oscuras. Procedió, luego, a catalogar meticulosamente la longitud de onda exacta de cada una de estas líneas, que hasta el día de hoy se conocen como líneas de Fraunhofer, y asignó letras a las más notorias. De esta forma, Fraunhofer registró el primer espectro astronómico de alta resolución.

Ahora bien, él no sabía cuál era la causa de que estas líneas oscuras estuvieran presentes en el espectro. Sin embargo, posteriormente procedió a realizar el mismo experimento, pero esta vez utilizando un rayo de luz proveniente de la estrella roja cercana Betelgeuse, y observó que el patrón de líneas oscuras cambiaba considerablemente. Fraunhofer concluyó correctamente que estas se encuentran de cierta forma relacionadas con la composición del objeto observado. En efecto, algunas de las líneas observadas por Fraunhofer se deben a las especies (e.g átomos, iones, moléculas) que componen la atmósfera terrestre.

Sin embargo, el gran paso en la comprensión general de las observaciones de Fraunhofer

llegó a mediados del siglo XIX de la mano del trabajo de los científicos Gustav Kirchhoff (1824 - 1887) y Robert Bunsen (1811 - 1899), quienes estudiaron el color de la luz emitida al poner distintos metales en llamas. Al hacer esto, descubrieron que, en ciertos casos, la longitud de onda de la luz emitida coincidía exactamente con las líneas observadas por Fraunhofer. Estos experimentos demostraron que las líneas de Fraunhofer son una consecuencia directa de la composición atómica del sol.

En el siglo XX se llegó a comprender de manera más profunda la razón de la existencia de estas líneas, denominadas *líneas espectrales*, gracias a la revolución que significó la llegada de la mecánica cuántica. Los desarrollos en materia de espectroscopía han estado, desde entonces, estrechamente ligados a los de aquel campo de la física.

Hoy en día, esencialmente toda la información con la que se cuenta sobre objetos astronómicos que residen fuera del sistema solar se ha obtenido mediante el estudio de la radiación electromagnética que estos emiten. Esta radiación contiene mucha información detallada, la cual puede ser obtenida solo mediante un análisis cuidadoso. En términos generales, se puede clasificar la información obtenida a partir de esta radiación según la resolución espectral; esto es, el grado de sensibilidad a las distintas longitudes de onda utilizado para realizar la observación.

Por ejemplo, cuando se observa el cielo de noche directamente con la vista, la mayoría de los objetos astronómicos se ven blancos. La luz blanca es en realidad luz que consta de muchas longitudes de onda y que no ha sido descompuesta en sus distintos colores. Observando esta luz blanca es posible obtener las posiciones de los objetos en el cielo nocturno, construir mapas de estrellas y galaxias, y registrar el movimiento de cuerpos celestes; como se ha hecho durante siglos hasta nuestros días.

Si se observa cuidadosamente ciertos objetos, tales como los planetas Marte o Júpiter, o estrellas tales como Betelgeuse, se puede apreciar que estos objetos tienden a tener un cierto color. Basta utilizar instrumentos de bajo poder resolutivo para separar la luz que llega desde estos objetos a la tierra en colores de amplio espectro. A su vez, el observar estos colores entrega información sobre la temperatura del objeto. Por ejemplo, las estrellas azules poseen mayor temperatura que las rojas. Objetos que emiten rayos x, como la corona solar, son muy calientes, mientras que objetos fríos emitirán radiación en longitudes de onda mayores; por ejemplo, en forma de ondas de radio.

La única forma de obtener información astrofísica detallada de objetos del cielo es mediante observaciones de alta resolución que involucren el detectar la intensidad de la radiación recibida en función de las longitudes de onda que la componen. Esto se lleva a cabo con equipos de alto poder resolutivo y sensibilidad. Dos ejemplos de estos son, el telescopio óptico SDSS que se encuentra en el Apache Point Observatory (APO, ubicado en Nuevo México, Estados Unidos) y con el cual se lleva a cabo el *Sloan Digital Sky Survey (SDSS)*; y, en mayor medida, el interferómetro radioastronómico *Atacama Large Millimeter/submillimeter Array (ALMA)* ubicado en el norte de Chile.

Observaciones llevados a cabo con estos equipos de alta resolución permiten obtener, no solamente la posición central de una línea dentro del espectro, sino también su forma. A partir de esta información, y con un conocimiento previo de física atómica y molecular,



puede extraerse valioso conocimiento sobre muchas de las propiedades del objeto y de su composición. Dado que existe una relación directa entre los parámetros físicos subyacentes y la información astronómica que se puede extraer a partir de los espectros, es posible utilizar datos generados a partir de observaciones experimentales en un laboratorio y compararlos con las líneas del espectro obtenidos de un objeto del cielo. Mediante este procedimiento se puede inferir propiedades del objeto, tales como su composición química, su temperatura, la abundancia de las especies que lo componen y que se encuentran emitiendo radiación, el movimiento de las especies y del objeto en sí, la presión y densidad local, el campo magnético presente, entre otros.

En síntesis, si se conoce esta información a partir de los datos de laboratorio, entonces una vez que se detecta un conjunto de líneas espectrales y se sabe la longitud de onda a la cual fueron emitidas es posible conocer a qué especies corresponden, y por tanto saber la composición química del objeto observado.

### 1.1.2. Sloan Digital Sky Survey (SDSS)

El *Sloan Digital Sky Survey (SDSS)* es un proyecto de inspección y estudio del espacio llevado a cabo mediante el uso de un telescopio óptico ubicado en el observatorio Apache Point (APO), Nuevo México, Estados Unidos. La recolección de datos comenzó en el año 2000, y las imágenes finales de los datos publicados cubren un 35 % del cielo, con observaciones fotométricas de 500 millones de objetos y espectros de radiación electromagnética de 1 millón de objetos.

El telescopio hace uso de la rotación terrestre para capturar pequeñas franjas del cielo, las cuales son registradas en un circuito integrado llamado dispositivo de carga acoplada o *charge-coupled device (CCD)* que captura las imágenes y permite transmitir las y almacenarlas en formato digital.

Utilizando los datos obtenidos de esta forma, se seleccionan objetos del cielo para su análisis espectroscópico. El espectrógrafo opera mediante asignar una fibra óptica individual a cada objeto que se desea observar y fijándola en su posición correspondiente a través de un agujero en una placa de aluminio. Cada agujero se ubica específicamente para el objeto deseado, por lo tanto, distintas áreas del cielo con distintos objetos requieren distintas placas de aluminio. El espectrógrafo de uso actual es capaz de registrar 1000 espectros a la vez. Cada noche se utilizan entre 6 a 9 placas para registrar espectros.

Los datos de SDSS se hacen disponibles mediante publicaciones regulares o *data releases* a través de internet. La última publicación llevada a cabo fue la correspondiente al data release 10 (DR10), con fecha de julio del 2013. Los datos de todos los data releases se encuentran en un servidor *Microsoft SQL Server* y pueden accederse mediante diversas interfaces o APIs presentes en el sitio web de SDSS. En particular, existe una interfaz web llamada *CasJobs* que permite realizar consultas en lenguaje *SQL* a un servidor que encola la petición, la ejecuta y guarda los resultados en una base de datos asignada al usuario.

### 1.1.3. Atacama Large Millimeter/submillimeter Array (ALMA)

El *Atacama Large Millimeter/submillimeter Array (ALMA)* es un interferómetro astronómico de radiotelescopios ubicados en el desierto de Atacama, en el norte de Chile. Es parte de un proyecto llevado a cabo mediante una asociación de organizaciones de Norteamérica, Europa y el este de Asia. Comenzó sus observaciones científicas en la segunda mitad del año 2011. Se encuentra completamente operacional desde marzo del año 2013; y es el mayor y más caro radiotelescopio construido hasta la fecha.

ALMA realiza observaciones captando radiación electromagnética proveniente del espacio en bandas milimétricas y submilimétricas en sus longitudes de onda, que corresponden a ondas de radio. Debido a que en condiciones normales la humedad del ambiente y del cielo absorbe gran parte de este tipo de radiación, es crucial para el funcionamiento de los telescopios el estar ubicados en uno de los lugares más secos del mundo, el llano de Chajnantor en el desierto de Atacama, a más de 5000 metros de altura.

Gran parte de los datos obtenidos desde ALMA son guardados en estructuras de datos llamadas cubos de datos tipo ALMA (o *ALMA Data Cubes*), que contienen información de distintos puntos de observación del cielo a distintas frecuencias. Los espectros de frecuencia son una forma de representar la intensidad de la radiación electromagnética, recibida desde un punto del espacio, en un cierto rango de frecuencias. Estos contienen puntos altos de intensidad en ciertas frecuencias en las cuales se sabe que una cierta molécula conocida efectúa una transición cuántica. Por lo tanto, mediante reconocer e identificar estos puntos altos, o *peaks*, se puede saber las transiciones moleculares que ocurrieron en el objeto del que proviene la radiación electromagnética. Como, a su vez, se sabe de antemano a qué frecuencia específicamente se realizan las transiciones cuánticas de moléculas conocidas, puede inferirse cuáles son las moléculas presentes en el objeto de origen.

A su vez, los cubos de datos tipo ALMA, como estructura de datos, contienen valores indexados en tres coordenadas. Dos de las coordenadas son espaciales, y corresponden al equivalente a una imagen normal de dos dimensiones, en el sentido que describen puntos del cielo (o del espacio observable desde la tierra). La tercera coordenada corresponde al rango de frecuencias en el que se está detectando radiación electromagnética. Por lo tanto, si se fijan las dos coordenadas espaciales (se fija un punto en el espacio) y se extraen todos los valores en la tercera coordenada de aquel punto, se obtiene el espectro de frecuencias observado en ese punto del espacio.

Los cubos de datos tipo ALMA, por tanto, contienen información de los espectros de frecuencia observados en todos los puntos de un sector dado del espacio. Todos los espectros presentes en un cubo se encuentran en un mismo rango de frecuencia. Sin embargo, distintos cubos de datos pueden tener observaciones hechas en distintos rangos de frecuencia entre sí.

A partir de ALMA se generan enormes cantidades de datos, los cuales, debido a su gran tamaño, necesariamente deben procesarse por parte de sistemas automatizados de extracción y análisis con el fin de facilitar a los investigadores el extraer información útil a partir de estos. La mayoría de estas herramientas se encuentran dentro de las áreas de investigación de la minería de datos y Machine Learning; disciplinas de la computación que han tenido un

gran auge en el último tiempo.

## 1.2. Reglas de asociación

El aprendizaje mediante reglas de asociación, o *Association Rule learning (ARL)*, es sin lugar a dudas uno de los métodos más populares y mejor estudiados dentro de la minería de datos. Basta para ello ver que el artículo seminal de Agrawal et al.[5], donde se sentaron las bases de la teoría subyacente, es uno de los más citados del área; según el catálogo y herramienta de búsqueda de publicaciones científicas *Google Scholar*.

La motivación principal de ARL en su concepción fue el encontrar relaciones lógicas entre los artículos adquiridos por usuarios en puntos de venta del tipo "Si un cliente compra los artículos  $A$  y  $B$ , entonces es muy probable que también compre el artículo  $C$ ". Sin embargo, la teoría de fondo que se desarrolló con el tiempo tiene una gran cantidad de aplicaciones en los más diversos ámbitos.

### 1.2.1. Definición formal

Sea  $\mathcal{I} = \{i_1, i_2, i_3, \dots, i_n\}$  un universo de ítemes posibles. Se denomina, entonces a un conjunto  $X \subseteq \mathcal{I}$  como *conjunto de ítemes* o *itemset*. Se tiene, además un conjunto de transacciones  $\mathcal{T} = \{T_1, T_2, \dots, T_m\}$ , donde  $T_i \subseteq \mathcal{I}$ ,  $\forall i \in [1, m]$ . Dados un conjunto de ítemes  $X$  y una transacción  $T_i$ , se dice que la transacción  $T_i$  *satisface*  $X$  si y solo si  $X \subseteq T_i$ .

Una *regla de asociación* es, entonces, una relación (más específicamente, una implicancia) entre dos conjuntos de la forma  $X \Rightarrow Y$ , donde  $X \subset \mathcal{I}$ ,  $Y \subset \mathcal{I}$ , y  $X \cap Y = \emptyset$ . A  $X$  se denomina el *antecedente* de la regla y a  $Y$  se denomina el *consecuente* de la regla.

Existen, también, una serie de medidas para cuantificar la relevancia de una regla de asociación. A continuación se define algunas de ellas.

El *soporte* de un conjunto de ítemes  $X$ , o  $supp(X)$ , se define como

$$supp(X) = \frac{|\mathcal{T}_X|}{|\mathcal{T}|} \quad , \text{ tal que } \quad \mathcal{T}_X = \{T \in \mathcal{T} : X \subset T\},$$

vale decir, corresponde a la fracción del total de transacciones en la que está presente el conjunto.

A su vez, el soporte de una regla de asociación  $X \Rightarrow Y$ , o  $supp(X \Rightarrow Y)$ , se define como

$$supp(X \Rightarrow Y) = supp(X \cup Y),$$

vale decir, corresponde a la fracción del total de transacciones en las cuales está presente tanto el antecedente como el consecuente de la regla simultáneamente<sup>1</sup>.

---

<sup>1</sup>Debe tenerse en mente que la expresión *mathsupp*( $X \cup Y$ ) indica la fracción del total de transacciones

La *confianza* de una regla de asociación  $X \Rightarrow Y$ , denotada por  $conf(X \Rightarrow Y)$ , se define como

$$conf(X \Rightarrow Y) = \frac{supp(X \cup Y)}{supp(X)},$$

es decir, indica en qué fracción de las transacciones en las cuales está presente el antecedente la regla se cumple (i.e. está presente también el consecuente de la regla). Debido al uso frecuente de esta medida de relevancia, resulta usual el expresar una regla de asociación mediante la notación

$$X \Rightarrow Y \Big|_c$$

donde  $c = conf(X \Rightarrow Y)$ .

El *lift* de una regla de asociación  $X \Rightarrow Y$ , denotado por  $lift(X \Rightarrow Y)$ , se define como

$$lift(X \Rightarrow Y) = \frac{conf(X \Rightarrow Y)}{supp(Y)} = \frac{supp(X \cup Y)}{supp(X) \times supp(Y)}.$$

La intuición detrás del concepto de lift tiene lugar al interpretar las medidas descritas anteriormente desde un punto de vista probabilístico. Tomando el conjunto  $\mathcal{T}$  como un universo de posibles resultados, o espacio muestral, se tiene que

$$supp(X) = P(X) \quad \text{y} \quad conf(X \Rightarrow Y) = P(Y|X).$$

Desde este punto de vista, la medida de lift indica qué tan bien la presencia del antecedente de una regla lograría predecir la presencia del consecuente. Por lo tanto, si la presencia del antecedente y del consecuente en una transacción cualquiera son eventos estadísticamente independientes (i.e. la ocurrencia de uno no afecta la probabilidad de que el otro ocurra), se tendrá que  $lift(X \Rightarrow Y) = 1$ ; y este valor irá variando en la medida que ambos eventos sean más dependientes entre sí.

Por ejemplo, supongamos que se tiene el siguiente conjunto de transacciones

TID	Items
1	$a, c$
2	$a, d$
3	$b, c$
4	$b, d$

donde *TID* es el número identificador de la transacción. Luego, para este caso, se tiene que

$$lift(\{a\} \Rightarrow \{c\}) = \frac{supp(\{a\} \cup \{c\})}{supp(\{a\}) \times supp(\{c\})} = \frac{1/4}{1/2 \times 1/2} = 1,$$

lo cual indica que la que la ocurrencia de que una transacción cualquiera satisfaga  $\{a\}$  es estadísticamente independiente de que una transacción cualquiera satisfaga  $\{b\}$ .

---

en las cuales está presente **tanto** el antecedente como el consecuente de la regla **simultáneamente**, y **no** de aquellas en las cuales está presente el antecedente **o** el consecuente. El argumento del soporte *supp* es un conjunto de “pre-condiciones”, y, por lo tanto, se vuelve más restrictivo en la medida que su cardinalidad aumenta.

En cambio, en el siguiente conjunto de transacciones

<b>TID</b>	<b>Items</b>
1	<i>a, c</i>
2	<i>a, d</i>
3	<i>b, c</i>
4	<i>b, c</i>

se tiene que

$$lift(\{a\} \Rightarrow \{c\}) = \frac{supp(\{a\} \cup \{c\})}{supp(\{a\}) \times supp(\{c\})} = \frac{1/4}{1/2 \times 3/4} = 2/3 < 1,$$

lo cual quiere decir que hay una mayor razón de transacciones que satisfacen  $\{c\}$  dentro del total de transacciones que dentro del conjunto de transacciones que satisfacen  $\{a\}$ .

Finalmente, en el conjunto de transacciones

<b>TID</b>	<b>Items</b>
1	<i>a, c</i>
2	<i>a, d</i>
3	<i>b, d</i>
4	<i>b, d</i>

se cumple que

$$lift(\{a\} \Rightarrow \{c\}) = \frac{supp(\{a\} \cup \{c\})}{supp(\{a\}) \times supp(\{c\})} = \frac{1/4}{1/2 \times 1/4} = 2 > 1,$$

lo cual indica que hay una mayor razón de transacciones que satisfacen  $\{c\}$  dentro del conjunto de transacciones que satisfacen  $\{a\}$  que dentro del total de transacciones.

### 1.2.2. Algoritmos, implementaciones y aplicaciones

En el mismo artículo seminal de ARL por Agrawal et al.[5], se presentó el algoritmo *Apriori*. Este algoritmo hace uso de las propiedades de clausura descendiente de la frecuencia de los conjuntos con respecto a sus subconjuntos con el fin de optimizar el proceso de generación de conjuntos de ítemes frecuentes.

Posteriormente, Agrawal et al. presentaron el algoritmo *AprioriTid*, cuyas mejores características fueron combinadas con el algoritmo *Apriori* para crear el algoritmo *AprioriHybrid*, de orden de complejidad lineal en el número de transacciones[7]. Luego se han realizado más desarrollos en ARL orientado a transacciones secuenciales de clientes de puntos de ventas[6].

Savasere et al. introdujeron el algoritmo *Partition*[35] con el fin de extraer reglas de asociación en base de datos, el cual presenta reducciones en las operaciones de la CPU y de entrada/salida, y que además facilita la paralelización. Posteriormente se creó el algoritmo *Dynamic Itemset Counting (DIC)*[10], que realiza menos lecturas sobre los datos que los

algoritmos previos, y que utiliza la métrica de *Convicción* a la hora de generar reglas de asociación. Luego, Park et al. presentaron un algoritmo que hace uso de funciones de Hashing con el fin de generar reglas candidatas[30]. Se han realizado, también, adaptaciones de los algoritmos previos con el fin de realizar ARL en datos de tipo cuantitativo[38].

Esfuerzos posteriores se han realizado con el fin de profundizar en los fundamentos teóricos subyacentes en ARL (e.g. definiendo el conjunto de posibles ítemes como una estructura algebraica llamada *retículo*)[41], y con el fin de extender la noción de reglas de asociación a correlaciones[9].

Más recientemente, Han et al. introdujeron el uso de una estructura de datos llamada *Frequent Pattern Tree*[23] en la extracción de reglas de asociación a partir de conjuntos de transacciones. Luego de esto, se han hecho numerosas implementaciones y optimizaciones a los algoritmos más utilizados en ARL, como, por ejemplo, el algoritmo Apriori[8]; así como implementaciones que facilitan el mantener la privacidad de cada una de las fuentes de datos que participan en el proceso[19].

Desde su concepción, el método de ARL ha sido aplicado en numerosas áreas, tales como la detección de intrusiones[28] y anomalías[31][13], educación[33][34], química[17], privacidad de datos[21], búsqueda en la web[20], tráfico en redes[18], computación social[29], búsqueda semántica[15], biología[27][12], salud[25][14], medios de comunicación[16][26], y la investigación forense[24]. Junto con esto, se han realizado numerosas investigaciones sobre el estado actual de ARL y sus posibles desarrollos a futuro dentro del marco de métodos automatizados de generación de conocimiento[22].

Si bien existen numerosos esfuerzos por utilizar minería de datos y Machine Learning en diversos ámbitos de la astronomía (en particular, en detección, clasificación y caracterización de líneas moleculares en espectros de emisión[36]), hasta la fecha no se ha propuesto abiertamente el uso de ARL sobre datos extraídos de espectros de frecuencia.

Sin embargo, se han realizado avances en ampliar los conceptos subyacentes en ARL con el fin de aplicar el método en campos más diversos[9]. Específicamente, una rama de investigación ha desarrollado lo que se denomina *Weighted Association Rule Learning*[39][11]. Este método permite asociar medidas de interés arbitrario a priori a ciertos conjuntos de datos. Si bien esto hace que se pierdan propiedades de clausura que son útiles a la hora de generar algoritmos eficientes, también permite trabajar con distintos conjuntos de transacciones sin que las reglas generadas estos dependan exclusivamente de su soporte u otras medidas estándar.

# Capítulo 2

## Especificación del Problema

### 2.1. Descripción del problema

Explicar que las líneas ya han sido detectadas. Como? Mencionar o citar el trabajo de Karim Pichara y Andrés Riveros?

Se asume que se cuenta con conjuntos de espectros, y que cada uno de ellos posee todas sus líneas espectrales correctamente detectadas y, por lo tanto, se conoce su posición en el espectro. En la práctica eso puede ser muy difícil de lograr, sobre todo en circunstancias donde pueden existir en principio una alta cantidad de líneas espectrales y estas pueden interferir unas con otras en la señal final, lo que se conoce como *blending*.

No sé si poner lo anterior acá o en el capítulo de marco teórico

Por lo tanto, para efectos de lo que sigue, basta con asumir que existe la posibilidad que no todas las líneas hayan sido detectadas. Pero es importante que las que sí fueron detectadas, lo hayan sido con una seguridad suficiente y que se sepa de manera adecuada su posición.

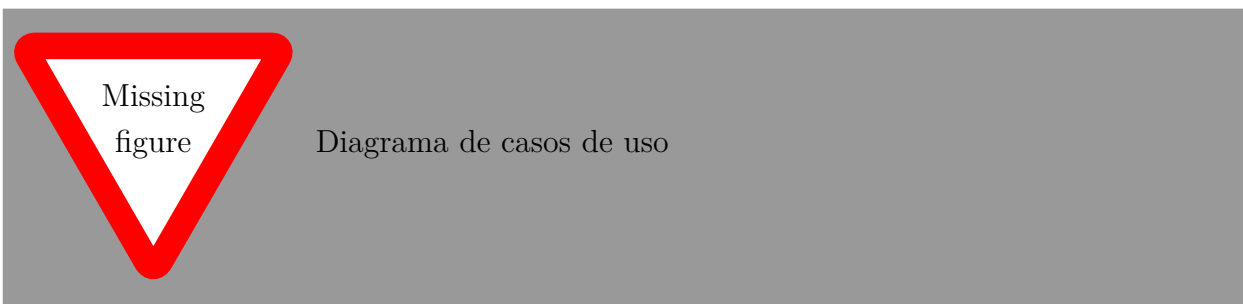
Explicar cómo se llevó a cabo la detección en términos más rigurosos. Cómo?

Teniendo estos conjuntos de espectros con sus respectivas líneas detectadas se desea...

Definir el problema de forma general, quizás con lineamientos de objetivos, que no involucre directamente como solución el uso de reglas de asociación

## 2.2. Requisitos de la solución y casos de uso

Definir los casos de uso





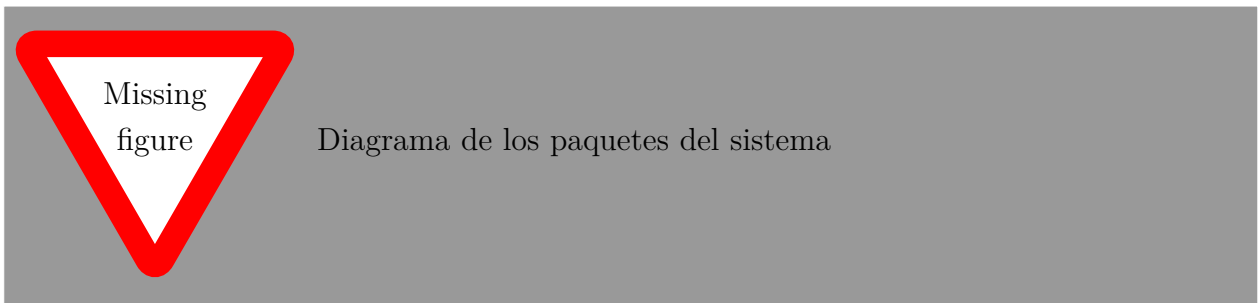
# Capítulo 3

## Descripción de la Solución

A continuación se describe la solución implementada para el presente proyecto. Se detalla aquí la estructura, diseño y funcionamiento del sistema y la aplicación realizados con el fin de cumplir con los requerimientos descritos anteriormente.

### 3.1. Arquitectura de software

Dado que, para fines del proyecto, se requería de una herramienta con la cual se pudiese llevar a cabo una serie de pruebas en distintos contextos, se optó por dividir el sistema en dos paquetes distintos; cada uno con una función específica, e interfaces bien definidas, con el fin de facilitar su posterior extensión y reutilización.



A continuación se detallan los paquetes del sistema, sus módulos, interfaces y funciones:

#### 3.1.1. Paquete de Association Rule Learning (ARL)

El paquete de Association Rule Learning (ARL) es el encargado de realizar el aprendizaje mediante reglas de asociación en sí; vale decir, de recibir un conjunto de datos con transacciones y de retornar reglas de asociación generadas a partir de aquel conjunto.

En las siguientes secciones se especifican los formatos de entrada y salida de este paquete junto con una descripción de los módulos que lo componen.

## Entrada y salida

Este paquete recibe como entrada un archivo de tabla en formato de valores separados por coma o *comma separated values (CSV)*. Este archivo debe tener el siguiente formato en cada una de sus filas

```
<TID>,"<ItemList>"
```

donde *<TID>* es el identificador de la presente transacción, e *<ItemList>* es una lista de identificadores únicos de los ítemes presentes en la transacción separados por comas. Tal como se indica, esta lista debe ir rodeada por comillas dobles en el archivo de entrada. A continuación se muestra un ejemplo de archivo de entrada válido.

```
000001,"15,2,44"  
000002,"5,4,23,67,43,234"  
000003,"66,3,53,23"
```

Adicionalmente, se puede especificar para cada transacción un tipo o clase a la que pertenece, o de la cual se origina, con el fin de realizar estadísticas pertinentes con las reglas generadas. De ser así, el archivo de entrada debe tener el siguiente formato en cada una de sus filas,

```
<TID>,<Class>,"<ItemList>"
```

donde, en esta ocasión, se añade en la segunda posición el campo *<Class>*, que consiste en una secuencia de caracteres válidos que identifique de manera unívoca la clase a la cual la transacción pertenece. A continuación un ejemplo de entrada válida en este formato.

```
000001,MORNING,"15,2,44"  
000002,MORNING,"5,4,23,67,43,234"  
000003,NIGHT,"66,3,53,23"
```

Esta lista es leída y procesada dentro del paquete de ARL y luego entregada en una estructura de datos correspondiente al algoritmo indicado, que obtendrá las reglas de asociación presentes en el conjunto de transacciones. Estas reglas, por defecto, serán retornadas en un archivo de texto en formato CSV con la siguiente estructura en cada una de sus líneas.

```
<N>,"<Antecedent>","<Consequent>",<Support>,<Confidence>,<Lift>
```

Donde *N* es un número identificador de la regla de asociación, *<Antecedent>* es una lista de ítemes separados por coma correspondientes al antecedente de la regla, *<Consequent>* es una lista de ítemes separados por coma correspondientes al consecuente de la regla, *<Support>* es un valor de punto flotante entre 0 y 1 correspondiente al soporte de la regla, *<Confidence>* es un valor de punto flotante entre 0 y 1 correspondiente a la confianza de la regla, y *<Lift>*

es un valor de punto flotante entre 0 y 1 correspondiente al lift de la regla. A continuación un ejemplo de este formato de archivo de salida.

```
1,"15,33","2,89,91",0.21,0.85,2.31
2,"12,33,44","5,23,31",0.23,0.81,3.3
```

Si, además, en los datos de entrada se especificó una clase para cada transacción, entonces el archivo de salida tendrá el siguiente formato

```
<N>,"<Antecedent>","<Consequent>",<Support>,<Confidence>,<Lift>,"<ClassCount>
```

en donde *<ClassCount>* es una lista de valores separados por comas con el siguiente formato

```
<Class01>:<Count01>,<Class02>:<Count02>,...
```

donde *<Class01>* es el identificador de la primera clase, *<Count01>* es un número entero que indica cuántas de las transacciones que satisfacen la regla actual pertenecen a esta primera clase, y así sucesivamente con todas las clases posibles. A continuación un ejemplo de archivo de salida con el formato recién descrito.

```
1,"15,33","2,89,91",0.21,0.85,2.31,"MORNING:210,NIGHT:15"
2,"12,33,44","5,23,31",0.23,0.81,3.3,"MORNING:20,NIGHT:91"
```

## Módulo de interfaz de usuario/controlador

El módulo de interfaz de usuario y controlador es el encargado de recibir directamente del usuario los parámetros de entrada correspondientes. Este módulo contiene métodos, clases y funciones que reciben los parámetros del usuario, abren y leen los archivos de entrada adecuados, los procesan de acuerdo al formato especificado, y hacen entrega de los datos al módulo principal de ARL.

Este módulo, es el encargado, además de recibir las reglas de asociación, entregarlas al módulo de formato para luego retornarlas al usuario en un archivo correspondiente.

## Módulo de formato

Es el módulo encargado de analizar los archivos de entrada leídos por el módulo de interfaz de usuario, extraer la información pertinente de ellos según el formato especificado, y retornar los datos en una estructura adecuada para luego ser procesados por el módulo principal de ARL. A su vez, este módulo realiza, además la labor inversa; vale decir, recibe las reglas de asociación en una estructura de datos estándar para luego entregarlas al módulo de interfaz en el formato requerido por el usuario.

Hasta el momento los formatos soportados son CSV para archivos de entrada, y CSV o tabla en formato  $\text{\LaTeX}$  para archivos de salida.

## **Módulo principal de ARL**

El módulo principal de ARL es el encargado de llevar a cabo el algoritmo de aprendizaje mediante reglas de asociación en sí. En su parte lógica, consta de dos sub-módulos principales. El primero es el sub-módulo encargado de extraer los conjuntos de ítemes frecuentes; vale decir, aquellos que cumplen con el requerimiento de soporte mínimo. Y el segundo es el sub-módulo de generación de reglas, que es el encargado de recibir los conjuntos de ítemes frecuentes y generar, a partir de ellos, las reglas de asociación que cumplen con el requerimiento de confianza mínima indicado.

## **Módulo de testeo de ARL**

Se encuentra dentro de este paquete, además, un módulo de testeo de los algoritmos de ARL sobre datos de prueba de pequeña envergadura; con el fin de realizar chequeos periódicos del funcionamiento correcto de estos algoritmos en la medida que se realizan cambios, mejoras o refactorizaciones sobre su código fuente.

## **Módulo de herramientas**

Finalmente, se encuentra el módulo de herramientas generales, que consta de una serie de funciones de uso frecuente por parte de otros módulos del paquete; tales como operaciones sobre listas anidadas, búsqueda de llaves sobre diccionarios específicos, entre otros.

### **3.1.2. Paquete de procesamiento de datos**

Debido a que, en la mayoría de las ocasiones los datos sobre los cuales se desea aplicar los algoritmos de reglas de asociación no se encuentran desde un comienzo en los formatos o estructuras necesarias, se procedió a implementar un paquete de pre-procesamiento. Este contiene una serie de scripts y métodos cuya función principal es extraer los datos desde sus fuentes originales, opcionalmente inferir aquella información que sea relevante, y guardarla en archivos cuyo formato sea comprensible para el paquete de aprendizaje de reglas de asociación.

En su implementación actual, este paquete se encuentra enfocado, en su mayor parte, para trabajar sobre datos extraídos a partir del Sloan Digital Sky Survey (SDSS).

A continuación se enumeran algunos de sus componentes más importantes.

## **Queries SQL**

Una colección de queries relevantes para ejecutar en las bases de datos de SDSS y extraer los datos sobre los cuales obtener las reglas de asociación.

## Módulo de procesamiento de tablas

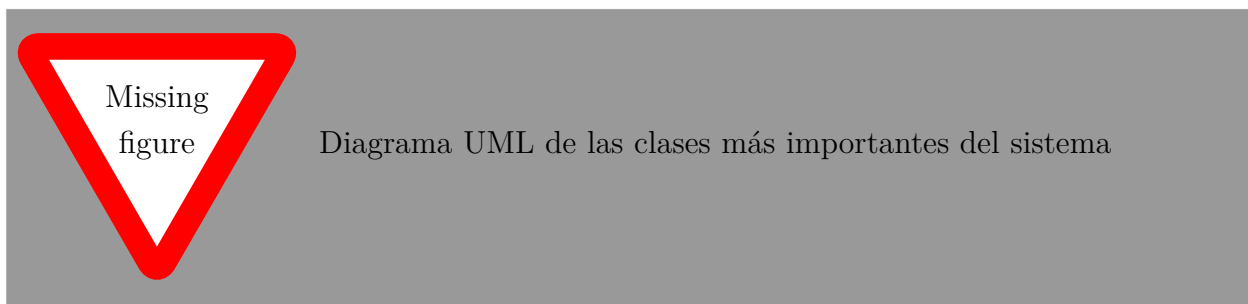
Contiene una serie de scripts cuyo fin es recibir un archivo de tabla de base de datos en formato CSV y procesar los datos que contiene; por ejemplo, eliminando ciertas filas, añadiendo columnas calculadas a partir de las ya existentes, entre otros. Los resultados son guardados en un nuevo archivo de tabla en formato CSV.

## Módulo de extracción de transacciones

Este módulo contiene scripts cuya función es recibir un archivo de tabla de base de datos en formato CSV, y a partir de él generar un archivo CSV que contenga una transacción por cada fila; cada una de estas con una lista de ítems en formato adecuado para ser recibido por el paquete de ARL.

## 3.2. Diseño de clases

A continuación se detallan las clases de objetos más importantes del sistema.



### 3.2.1. Clase *ItemSet*

La clase *ItemSet* es la encargada de mantener información sobre un conjunto de ítems y abstraer su estructura de datos subyacente. Cada instancia de esta clase corresponde a un conjunto de ítems distinto, y contiene campos que guardan la información más reciente sobre su soporte (calculado sobre un cierto conjunto de transacciones) y punteros a meta-datos con información adicional sobre los ítems en sí. Su interfaz asegura que se pueda realizar de forma adecuada, visto desde un punto de vista matemáticamente abstracto, las operaciones más comunes de conjuntos de elementos; como comprobar pertenencia, sumar de conjuntos, diferencia entre conjuntos, entre otros.

### 3.2.2. Clase *AssociationRule*

La clase *AssociationRule* es la que define la estructura y comportamiento de las reglas de asociación. Cada instancia de esta clase corresponde a una regla de asociación en particular, extraída a partir de un cierto conjunto de datos. Cada regla de asociación consta de dos objetos de la clase *ItemSet*; uno para el antecedente y otro para el consecuente de la regla. Además contiene un campo que codifica su soporte, junto con métodos para calcular sus medidas de relevancia, tales como su confianza y lift.

### 3.2.3. Clase *FrequentItemSetMiner*

La clase *FrequentItemSetMiner* es la encargada de abstraer y guardar información sobre el proceso de extraer a partir de las transacciones aquellos conjuntos de ítemes que cumplan con un requisito de soporte mínimo dado. Cada instancia de esta clase corresponde a un proceso de extracción distinto, conteniendo campos y estructuras de datos para los algoritmos involucrados, su estado actual y su resultado.

En su implementación actual, esta clase es heredada por dos sub-clases. Una correspondiente al algoritmo *Apriori*, y otra al algoritmo *FP-Growth*. Cada una contiene su propia implementación de los métodos principales, definidos en su clase padre, junto con sus propias funciones auxiliares y estructuras de datos correspondientes.

### 3.2.4. Clase *RuleMiner*

La clase *RuleMiner* es la que abstrae el proceso de extraer reglas de asociación a partir de conjuntos frecuentes de ítemes. Cada instancia de esta clase corresponde a un proceso de extracción distinto; básicamente el mismo en todo los casos salvo en ciertos detalles, como algunas funciones auxiliares y referencias a estructuras de datos, dependiendo de si los conjuntos fueron extraídos mediante *Apriori* o *FP-Growth*.

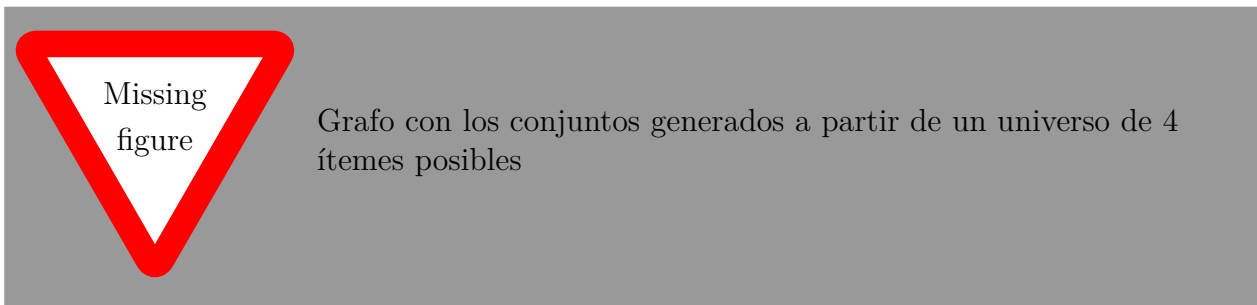
## 3.3. Algoritmos y estructuras de datos

Para el desarrollo de este sistema se llevó a cabo la implementación de dos algoritmos de aprendizaje de reglas de asociación. El primero es el algoritmo *Apriori* y el segundo es una optimización de este, llamada *FP-Growth*. Sus procedimientos generales y estructuras de datos principales se describen a continuación.

### 3.3.1. Algoritmo *Apriori*


El algoritmo *Apriori* recibe como entrada un conjunto de transacciones, y tiene como objetivo encontrar y retornar todos aquellos conjuntos presentes que cumplan con el requisito de soporte mínimo indicado, también llamados *conjuntos frecuentes*.

Por ejemplo, supongamos que se cuenta con el un conjunto de transacciones, y que cada una contiene ítems pertenecientes a un universo de solo 4 elementos posibles,  $\mathcal{I} = \{0, 1, 2, 3\}$ . Luego, en principio, para extraer los conjuntos frecuentes a partir de estas transacciones, por cada uno de los conjuntos que es posible generar con este universo de 4 ítems posibles (llamados *conjuntos candidatos*), se debe recorrer cada una de las transacciones, ver si la transacción satisface este conjunto, y de ser así incrementar un contador. Luego de terminar este proceso para cada uno de los conjuntos posibles, se tendrá el número de veces que cada uno de estos se encuentra dentro del conjunto de transacciones, y teniendo el número total de estas, se puede obtener de forma directa el soporte de estos conjuntos.



El problema radica en que el número de conjuntos candidatos crece de manera exponencial en el número de ítems del universo posible. En efecto, si el número de ítems del universo es  $n$ , entonces a partir de este es posible generar  $2^n + 1$  conjuntos. Por tanto, para un universo de 100 elementos, existen nada menos que  $1,26 \times 10^{30}$  conjuntos candidatos; y debe, por tanto, recorrerse el total de transacciones este número de veces.

No obstante, es posible reducir el número de conjuntos candidatos utilizando la propiedad de *clausura descendiente* de los conjuntos frecuentes, también llamado *principio Apriori*. Esta propiedad asegura que si un conjunto dado es, en efecto, frecuente, entonces necesariamente todos sus subconjuntos también lo son. O, expresado de forma recíproca, si un conjunto dado resulta no ser frecuente, entonces necesariamente todos sus superconjuntos tampoco lo son. Esta última expresión es la que resulta más relevante para nuestro caso. Esto implica que luego de generar un conjunto candidato y verificar si es frecuente verificando el número de transacciones que lo satisfacen, si se comprueba que este conjunto no es frecuente (vale decir, no cumple con el requisito de soporte mínimo), entonces necesariamente ninguno de sus conjuntos posibles que lo contienen será frecuente, y por tanto no será necesario obtener sus soportes correspondientes contando el número de transacciones que los satisfacen.



Missing  
figure

Grafo igual al anterior, pero que muestra cuales de los conjuntos necesariamente no son frecuentes si uno de ellos resulta no serlo.

Esta propiedad permite reducir considerablemente el número de conjuntos candidatos y, por tanto, optimizar el algoritmo final; ya que no será necesario recorrer el total de transacciones tantas veces como se planteó originalmente. Para poder utilizar esta propiedad y beneficiarse de la optimización correspondiente, es necesario generar los conjuntos candidatos comenzando por aquellos que poseen menos elementos, y a partir de estos generar todos los superconjuntos posibles.

El algoritmo *Apriori*, por lo tanto, en terminos generales resulta ser el siguiente

---

**Algoritmo 1:** Algoritmo *Apriori*

---

**Data:** Conjunto de transacciones  $\mathcal{T}$   
**Result:** Conjunto de ítemes frecuentes  $\mathcal{L}$   
 $\mathcal{L}_1 \leftarrow \{\text{conjuntos de 1 solo ítem}\}$   
**for**  $k = 2; L_{k-1} \neq \emptyset; k++$  **do**  
     $C_k = \text{apriori-gen}(L_{k-1})$   
    **for** *transacciones*  $T \in \mathcal{T}$  **do**  
         $C_1 = \text{subset}(C_k, T)$   
        **for** *candidatos*  $C \in C_t$  **do**  
             $C.\text{count} \leftarrow C.\text{count} + 1$   
     $\mathcal{L}_k = \{C \in C_k : C.\text{count} \geq \text{minsup}\}$   
 $\mathcal{L} \leftarrow \bigcup_k \mathcal{L}_k$

---

### 3.3.2. Algoritmo *FP-Growth*

### 3.3.3. Extracción de reglas de asociación

### 3.3.4. Desempeño comparativo

## 3.4. Interfaz de usuario



# Capítulo 4

## Validación de la Solución

### 4.1. Antecedentes de datos de prueba

Explicar datos de sloan

### 4.2. Pre-procesamiento de datos

Filtrado, de objetos, de líneas. Gráficos.

### 4.3. Resultados

Mostrar las reglas de asociación resultantes con tablas, gráficos, grafos, figuritas, pegatinas, etc...

# Conclusión

Breve resumen del trabajo realizado

Recuento de objetivos alcanzados y no alcanzados

Análisis crítico de por qué los resultados fueron los reportados

Reflexión acerca de la relevancia / impacto del trabajo realizado

Lecciones aprendidas

Posibles trabajos futuros que podrían hacerse a partir de la memoria para mejorar aún más la solución

# Bibliografía

- [1] Atacama large millimeter/submillimeter array (alma). <http://www.almaobservatory.org>, 2014. online, accesed July 2014.
- [2] matplotlib: python plotting. <http://www.matplotlib.org>, 2014. online, accesed July 2014.
- [3] Numpy – numpy. <http://www.sdss.org>, 2014. online, accesed July 2014.
- [4] Welcome to python.org. <http://www.python.org>, 2014. online, accesed July 2014.
- [5] Rakesh Agrawal, Tomasz Imieliński, and Arun Swami. Mining association rules between sets of items in large databases. In *ACM SIGMOD Record*, volume 22, pages 207–216. ACM, 1993.
- [6] Rakesh Agrawal and Ramakrishnan Srikant. Mining sequential patterns. In *Data Engineering, 1995. Proceedings of the Eleventh International Conference on*, pages 3–14. IEEE, 1995.
- [7] Rakesh Agrawal, Ramakrishnan Srikant, et al. Fast algorithms for mining association rules. In *Proc. 20th int. conf. very large data bases, VLDB*, volume 1215, pages 487–499, 1994.
- [8] Ferenc Bodon. A fast apriori implementation. In *Proceedings of the IEEE ICDM workshop on frequent itemset mining implementations (FIMI'03)*, volume 90, 2010.
- [9] Sergey Brin, Rajeev Motwani, and Craig Silverstein. Beyond market baskets: Generalizing association rules to correlations. In *ACM SIGMOD Record*, volume 26, pages 265–276. ACM, 1997.
- [10] Sergey Brin, Rajeev Motwani, Jeffrey D Ullman, and Shalom Tsur. Dynamic itemset counting and implication rules for market basket data. In *ACM SIGMOD Record*, volume 26, pages 255–264. ACM, 1997.
- [11] Chun Hing Cai, Ada Wai-Chee Fu, CH Cheng, and WW Kwong. Mining association rules with weighted items. In *Database Engineering and Applications Symposium, 1998. Proceedings. IDEAS'98. International*, pages 68–77. IEEE, 1998.
- [12] Pedro Carmona-Saez, Monica Chagoyen, Francisco Tirado, Jose M Carazo, and Alberto

- Pascual-Montano. Genecodis: a web-based tool for finding significant concurrent annotations in gene lists. *Genome biology*, 8(1):R3, 2007.
- [13] Varun Chandola, Arindam Banerjee, and Vipin Kumar. Anomaly detection: A survey. *ACM Computing Surveys (CSUR)*, 41(3):15, 2009.
  - [14] R Chaves, JM Górriz, J Ramírez, IA Illán, D Salas-Gonzalez, and M Gómez-Río. Efficient mining of association rules for the early diagnosis of alzheimer’s disease. *Physics in medicine and biology*, 56(18):6047, 2011.
  - [15] Edith Cohen, Amos Fiat, and Haim Kaplan. Associative search in peer to peer networks: Harnessing latent semantics. *Computer Networks*, 51(8):1861–1881, 2007.
  - [16] James Davidson, Benjamin Liebald, Junning Liu, Palash Nandy, Taylor Van Vleet, Ullas Gargi, Sujoy Gupta, Yu He, Mike Lambert, Blake Livingston, et al. The youtube video recommendation system. In *Proceedings of the fourth ACM conference on Recommender systems*, pages 293–296. ACM, 2010.
  - [17] Luc Dehaspe, Hannu Toivonen, and Ross D King. Finding frequent substructures in chemical compounds. In *KDD*, volume 98, page 1998, 1998.
  - [18] Cristian Estan, Stefan Savage, and George Varghese. Automatically inferring patterns of resource consumption in network traffic. In *Proceedings of the 2003 conference on Applications, technologies, architectures, and protocols for computer communications*, pages 137–148. ACM, 2003.
  - [19] Alexandre Evfimievski, Ramakrishnan Srikant, Rakesh Agrawal, and Johannes Gehrke. Privacy preserving mining of association rules. *Information Systems*, 29(4):343–364, 2004.
  - [20] Paolo Ferragina and Antonio Gulli. A personalized search engine based on web-snippet hierarchical clustering. *Software: Practice and Experience*, 38(2):189–225, 2008.
  - [21] Gabriel Ghinita, Panos Kalnis, Ali Khoshgozaran, Cyrus Shahabi, and Kian-Lee Tan. Private queries in location based services: anonymizers are not necessary. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, pages 121–132. ACM, 2008.
  - [22] Jiawei Han, Hong Cheng, Dong Xin, and Xifeng Yan. Frequent pattern mining: current status and future directions. *Data Mining and Knowledge Discovery*, 15(1):55–86, 2007.
  - [23] Jiawei Han, Jian Pei, Yiwen Yin, and Runying Mao. Mining frequent patterns without candidate generation: A frequent-pattern tree approach. *Data mining and knowledge discovery*, 8(1):53–87, 2004.
  - [24] Farkhund Iqbal, Hamad Binsalleeh, Benjamin Fung, and Mourad Debbabi. A unified data mining solution for authorship analysis in anonymous textual communications. *Information Sciences*, 231:98–112, 2013.

- [25] Murat Karabatak and M Cevdet Ince. An expert system for detection of breast cancer based on association rules and neural network. *Expert Systems with Applications*, 36(2):3465–3469, 2009.
- [26] Georgi Kobilarov, Tom Scott, Yves Raimond, Silver Oliver, Chris Sizemore, Michael Smethurst, Christian Bizer, and Robert Lee. Media meets semantic web—how the bbc uses dbpedia and linked data to make connections. In *The semantic web: research and applications*, pages 723–737. Springer, 2009.
- [27] Stefan Kramer, Luc De Raedt, and Christoph Helma. Molecular feature mining in hiv data. In *Proceedings of the seventh ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 136–143. ACM, 2001.
- [28] Wenke Lee and Salvatore J Stolfo. A framework for constructing features and models for intrusion detection systems. *ACM transactions on Information and system security (TiSSEC)*, 3(4):227–261, 2000.
- [29] Xin Li, Lei Guo, and Yihong Eric Zhao. Tag-based social interest discovery. In *Proceedings of the 17th international conference on World Wide Web*, pages 675–684. ACM, 2008.
- [30] Jong Soo Park, Ming-Syan Chen, and Philip S Yu. *An effective hash-based algorithm for mining association rules*, volume 24. ACM, 1995.
- [31] Animesh Patcha and Jung-Min Park. An overview of anomaly detection techniques: Existing solutions and latest technological trends. *Computer Networks*, 51(12):3448–3470, 2007.
- [32] Anthony J Remijan and Andrew J Markwick-Kemper. Splatalogue: Database for astronomical spectroscopy. 2008.
- [33] Cristóbal Romero and Sebastian Ventura. Educational data mining: A survey from 1995 to 2005. *Expert Systems with Applications*, 33(1):135–146, 2007.
- [34] Cristóbal Romero, Sebastián Ventura, and Enrique García. Data mining in course management systems: Moodle case study and tutorial. *Computers & Education*, 51(1):368–384, 2008.
- [35] Ashok Savasere, Edward Robert Omiecinski, and Shamkant B Navathe. An efficient algorithm for mining association rules in large databases. 1995.
- [36] Petr Škoda and Jaroslav Vážný. Searching of new emission-line stars using the astroinformatics approach. *arXiv preprint arXiv:1112.2775*, 2011.
- [37] Ramakrishnan Srikant and Rakesh Agrawal. Mining generalized association rules. In *VLDB*, volume 95, pages 407–419, 1995.
- [38] Ramakrishnan Srikant and Rakesh Agrawal. Mining quantitative association rules in large relational tables. In *ACM SIGMOD Record*, volume 25, pages 1–12. ACM, 1996.

- [39] Wei Wang, Jiong Yang, and Philip S Yu. Efficient mining of weighted association rules (war). In *Proceedings of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 270–274. ACM, 2000.
- [40] Donald G York, J Adelman, John E Anderson Jr, Scott F Anderson, James Annis, Neta A Bahcall, JA Bakken, Robert Barkhouser, Steven Bastian, Eileen Berman, et al. The sloan digital sky survey: Technical summary. *The Astronomical Journal*, 120(3):1579, 2000.
- [41] Mohammed Javeed Zaki and Mitsunori Ogihara. Theoretical foundations of association rules. In *3rd ACM SIGMOD Workshop on Research Issues in Data Mining and Knowledge Discovery*, pages 71–78. Citeseer, 1998.