

## FASE 4: TRANSICIÓN DE LA FORMULACIÓN DE IDEAS A LOS DISEÑOS PRELIMINARES

### Descarte de ideas no factibles

<i>Idea no factible</i>	Justificación
<i>1.b. Técnica CRISP-DM (Cross Industry Standard Process for Data Mining).</i>	Debido a la naturaleza del proyecto, esta metodología para el desarrollo de proyectos excede los alcances del mismo. Exigiendo un mayor esfuerzo para obtener resultados que se pueden alcanzar con otra metodología más simple. Es por esto, que se descarta la idea de utilizar la técnica CRISP-DM.
<i>2.a. Modelo de Identificación.</i>	No va acorde a los objetivos que busca lograr la solución, ya que sólo evidenciar la existencia de objetos, eventos o actividades en un conjunto de datos no es suficiente para generar una solución al problema planteado.
<i>3.a Uso de grafos para ver qué producto está asociado con cuál en el entorno de ventas.</i>	Esta idea de presentación se descarta debido a su complejidad al momento de implementar la visualización gráfica de un grafo. Además, no es una alternativa que sea muy clara para representar varias asociaciones.

Después de realizar el descarte de ideas no factibles, se tienen las siguientes alternativas:

1. Metodologías para el desarrollo de proyectos en Minería de Datos:
  - a. Técnica Proceso de Generación de Conocimiento o KDD (*Knowledge Discovery in DataBases*).
2. Modelos (y algoritmos respectivos) para la Minería de Datos:
  - a. Agrupación
    - i. Algoritmo de Clústeres K-Means.
    - ii. Algoritmo de Clústeres K-Medoids.
    - iii. Autómatas finitos.
  - b. Asociación
    - i. Algoritmo de Fuerza Bruta.
    - ii. Algoritmo Apriori
    - iii. Algoritmo Partition
  - c. Predicción
    - i. Algoritmo de Bayes naive.
    - ii. Algoritmo de Regresión logística.
3. Presentación (gráfica) de resultados:
  - a. Uso de histogramas para mostrar la frecuencia de los itemsets.
  - b. Uso de tablas para mostrar las asociaciones encontradas.

## Técnica Proceso de Generación de Conocimiento o KDD (*Knowledge Discovery in DataBases*).

Este proceso surge cuatro pasos para la generación de conocimiento. Estas etapas pueden ser recursivas, es decir, que se retorna a ellas una y otra vez (proceso iterativo) a medida que se obtienen resultados preliminares que requieren replantear las variables iniciales.

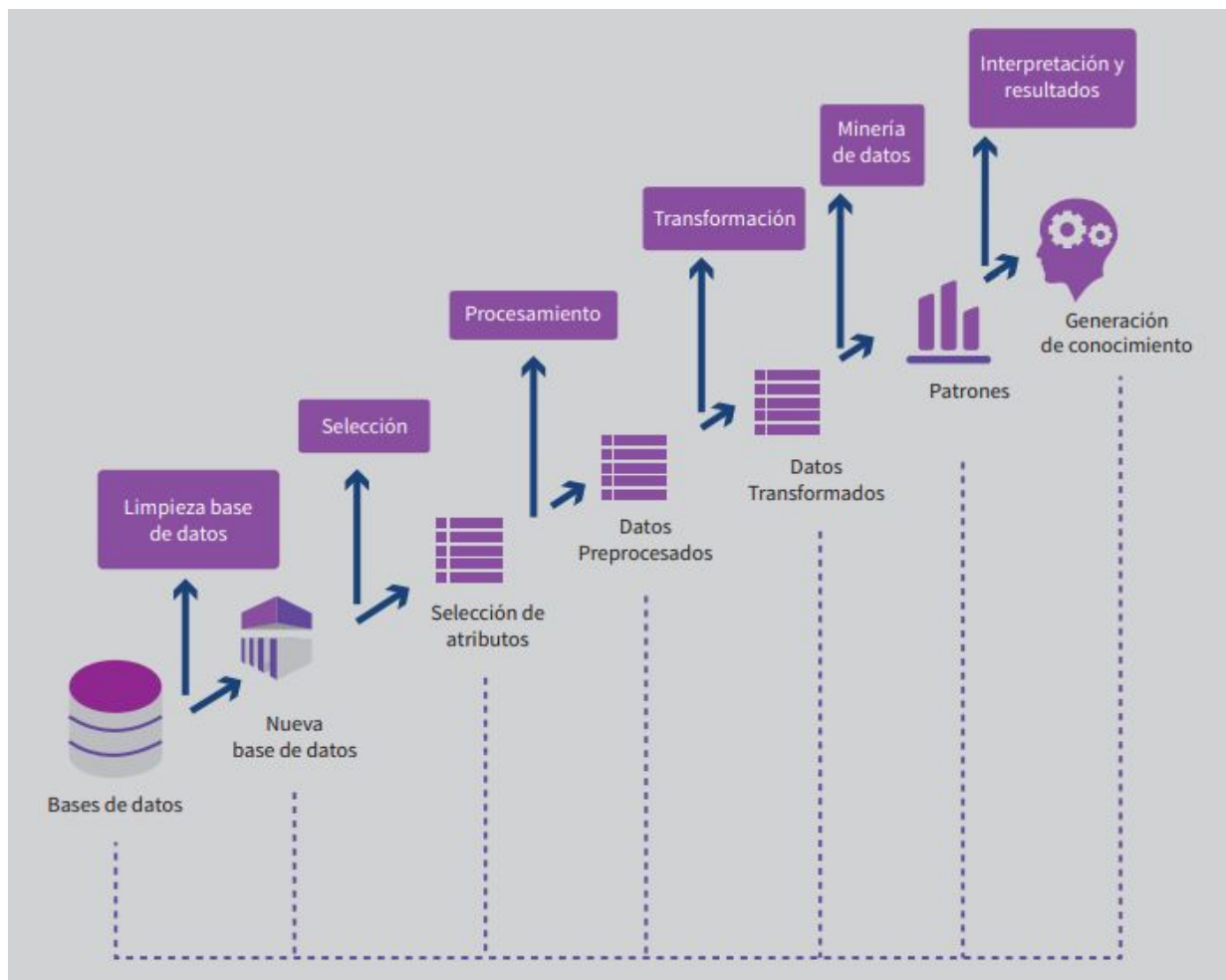


Figura . Proceso de KDD

Las etapas del proceso son:

1. Selección de los datos
2. Pre procesamiento de datos
3. Selección de características
4. Minería de Datos
5. Interpretación y Resultados

Etapas	Descripción
--------	-------------

<b>Selección de los datos</b>		Consiste en la recolección y preparación de los datos. En este proceso se comprende la problemática asociada a la base de datos y se establecen objetivos. Y se identifican las variables que serán consideradas para la construcción del modelo de minería de datos
<b>Pre-procesamiento de datos</b>	Integración de datos	Se analiza si la base de datos requiere incluir o integrar información o variables que reposan en otras bases de datos, y que será relevante para el modelo de minería de datos
	Reconocimiento y limpieza	Se depura el conjunto de datos respecto a valores atípicos, faltantes y erróneos (eliminación de ruido e inconsistencias).
<b>Selección de características</b>	Exploración y limpieza de datos	Aplicando técnicas de análisis exploratorio de datos se busca identificar la distribución de los datos, simetría, pruebas de normalidad y correlaciones existentes entre los datos. En esta etapa es útil el análisis descriptivo del conjunto de datos (clustering y segmentación, escalamiento, reglas de asociación y dependencia, reducción de la dimensión), identificación de datos nulos, ruido y outliers, así como el uso de matrices de correlación (si las variables son numéricas), diagramas (barras, histogramas, caja y bigotes), entre otras técnicas adecuadas de muestreo
	Transformación	Se estandariza o normaliza la información (colocarla en los mismos términos de formato y forma).
	Reducción de datos	Se disminuye el tamaño de los datos mediante la eliminación de características redundantes.
<b>Minería de Datos</b>		Se puede definir como un proceso no trivial de identificación válida, novedosa, potencialmente útil y entendible de patrones comprensibles que se encuentran ocultos en los datos, que a su vez, facilita la toma de decisiones y emplea técnicas de aprendizaje supervisado y no-supervisado.
<b>Interpretación y Resultados</b>		Se analizan los resultados de los patrones obtenidos en la fase de MD, mediante técnicas de visualización y de representación, con el fin de generar conocimiento que aporte mayor valor

	<p>a los datos. En esta fase se evalúan los resultados con los expertos y, si es necesario, se retorna a las fases anteriores para una nueva iteración</p>
--	--

Tabla . Proceso KDD

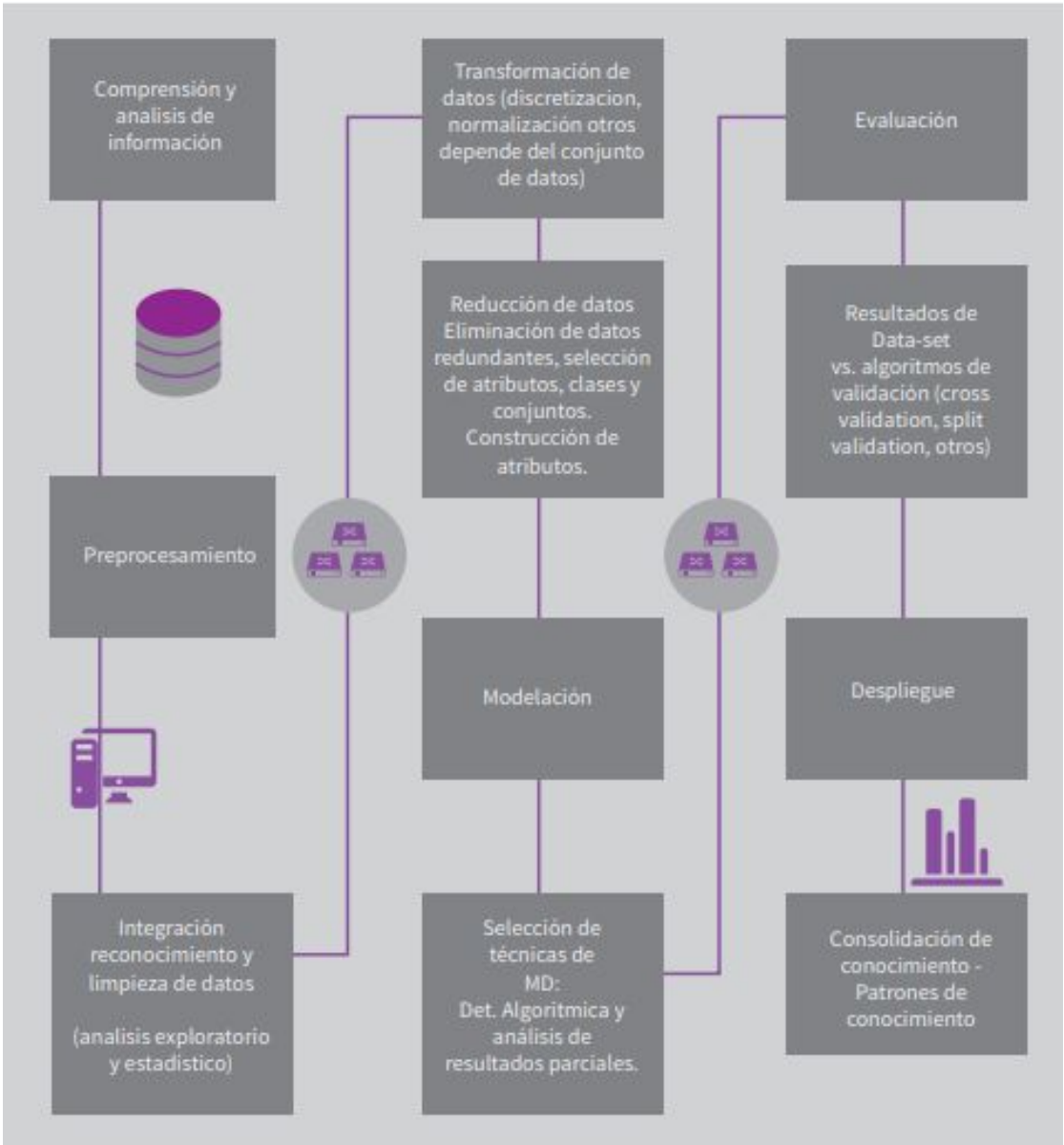
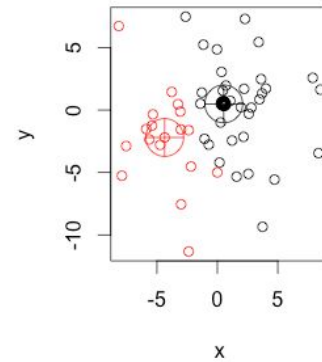
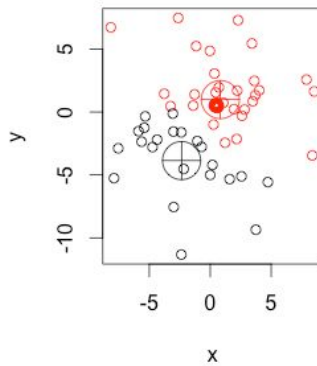


Imagen. Proceso KDD

Técnica de Agrupación utilizando el Algoritmo de Clústeres K-Means.	Técnica de Agrupación utilizando el Algoritmo de Clústeres K-Medoids.
El análisis de conglomerados o Clustering, es una técnica que permite analizar y examinar datos que no se encuentran etiquetados, formando conjuntos de grupos a partir de su similitud.	
Pseudocódigo	
<p><b>Algorithm 1:</b> K-Means Algorithm</p> <hr/> <p><b>Input:</b> <math>E = \{e_1, e_2, \dots, e_n\}</math> (set of entities to be clustered)  <math>k</math> (number of clusters)  <math>MaxIters</math> (limit of iterations)</p> <p><b>Output:</b> <math>C = \{c_1, c_2, \dots, c_k\}</math> (set of cluster centroids)  <math>L = \{l(e) \mid e = 1, 2, \dots, n\}</math> (set of cluster labels of E)</p> <pre> foreach <math>c_i \in C</math> do     <math>c_i \leftarrow e_j \in E</math> (e.g. random selection) end foreach <math>e_i \in E</math> do     <math>l(e_i) \leftarrow \operatorname{argmin}_{j \in \{1 \dots k\}} \operatorname{Distance}(e_i, c_j)</math> end  changed <math>\leftarrow</math> false; iter <math>\leftarrow</math> 0; repeat   foreach <math>c_i \in C</math> do       <math>\operatorname{UpdateCluster}(c_i)</math>;   end   foreach <math>e_i \in E</math> do       <math>\minDist \leftarrow \operatorname{argmin}_{j \in \{1 \dots k\}} \operatorname{Distance}(e_i, c_j)</math>;       if <math>\minDist \neq l(e_i)</math> then           <math>l(e_i) \leftarrow \minDist</math>;           changed <math>\leftarrow</math> true;       end   end   iter ++; until changed = true and iter <math>\leq</math> MaxIters ; </pre>	<p><b>Algorithm 1</b> The <math>k</math>-Medoid Algorithm</p> <hr/> <p><b>Input:</b> <math>k</math> : The total number of clusters to generate,  <b>Input:</b> <math>D</math> : The collection of input documents.  <b>Input:</b> <math>T</math> : Total number of iterations allowed.  <b>Output:</b> <math>P</math> : The array of <math>k</math> clusters of the partition.</p> <pre> 1: Function k-Medoids(<math>k, D</math>) 2:   iterCount <math>\leftarrow</math> 0 3:   repeat 4:     <math>\{m_1, m_2, \dots, m_k\} \leftarrow \operatorname{InitializeMedoids}(k, D)</math> 5:     <math>P \leftarrow [\{m_1\}, \{m_2\}, \dots, \{m_k\}]</math> 6:     repeat 7:       <math>\tilde{D} \leftarrow \{d_i : d_i \in D \wedge d_i \notin \{m_1, m_2, \dots, m_k\}\}</math> 8:       for each: <math>d_i \in \tilde{D}</math> do 9:         for <math>j = 1</math> to <math>k</math> do 10:          <math>\operatorname{sim}_{i,j} \leftarrow \operatorname{CosineSimilarity}(d_i, m_j)</math> 11:        end for 12:        <math>i_{\max} \leftarrow \operatorname{argmax}_{1 \leq j \leq k} (\operatorname{sim}_{i,j})</math> 13:        <math>P[i_{\max}] \leftarrow P[i_{\max}] \cup \{d_i\}</math> 14:      end for 15:      <math>\operatorname{oldMedoids} \leftarrow \{m_1, m_2, \dots, m_k\}</math> 16:      <math>\{m_1, m_2, \dots, m_k\} \leftarrow \operatorname{UpdateMedoids}(P)</math> 17:      iterCount <math>\leftarrow</math> iterCount + 1 18:    until (<math>\operatorname{oldMedoids} = \{m_1, m_2, \dots, m_k\}</math>) <math>\vee</math> (iterCount = <math>T</math>) 19:    NoTinyClusters <math>\leftarrow</math> CheckSizeOfClusters(<math>C</math>) 20:  until (NoTinyClusters = True) <math>\vee</math> (iterCount = <math>T</math>) 21:  return <math>P</math> 22: end function </pre>
Clústeres K-Means.	Clústeres K-Medoids.

Ambos algoritmos intentan minimizar la distancia entre los puntos del mismo grupo y un punto particular que es el centro de ese grupo.



### Similitudes

- La variación dentro del clúster disminuye con cada iteración del algoritmo.
- El algoritmo siempre converge, depende en los centros iniciales. Para cualquiera de los algoritmos, uno debe ejecutarlo varias veces con diferentes inicios
- La agrupación final depende de los centros iniciales del clúster. Diferentes comienzos dan como resultado diferentes agrupaciones finales.

### Diferencias

- Devuelve centros que son promedios de puntos de datos.

- Generalmente devuelve un mayor valor de

$$\sum_{k=1}^K \sum_{C(i)=k} \|X_i - c_k\|_2^2$$

- Computacionalmente más difícil (debido al paso 2: calcular el medoid es más difícil que calcular el promedio)
- Tiene la propiedad (potencialmente importante) de que los centros se encuentran entre los puntos de datos
- Se basa en el cálculo de centroides (o medoides) minimizando la distancia absoluta entre los puntos y el centroide seleccionado, en lugar de minimizar la distancia cuadrada.

## Técnica de Asociación utilizando el Algoritmo de Fuerza Bruta.

Es una técnica de búsqueda también llamada búsqueda exhaustiva que consiste en observar todos los posibles candidatos a una solución.

### Pseudocódigo

---

**Algorithm 1:** COMBINATIONS( $c, S$ )

---

**Input:**  $c, S$

**Output:** All subsets of size  $c$  of  $S$

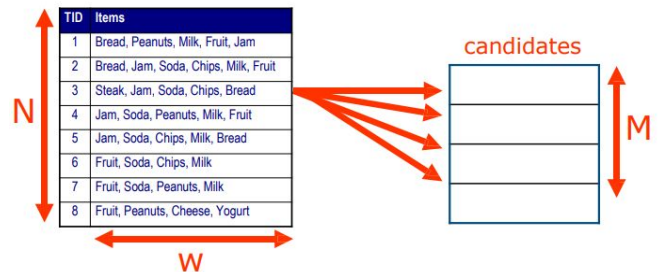
```
1  $root \leftarrow \emptyset$ 
2  $s \leftarrow |S|$ 
3 COMBINATIONSREC( $c, 0, root, s$ )
4  $icombinations \leftarrow \text{EXTRACTFROMTREE}(root)$ 
5  $combinations \leftarrow \text{MAPPOSITIONSTOELEMENTS}(icombinations)$ 
6 return  $combinations$ 
```

---

### Complejidad

Complejidad temporal

$$O(NMw), M = 2^n$$



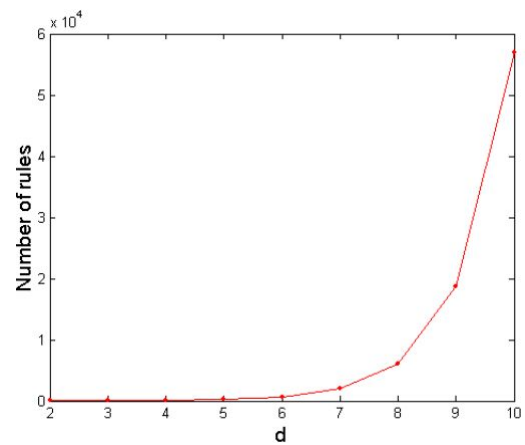
Complejidad computacional

Dado un número único  $d$  de item.

- Número total de itemsets:  $2^d$
- Número total de posibles asociaciones:



$$\sum_{k=1}^{d-1} \left[ \binom{d}{k} \times \sum_{j=1}^{d-k} \binom{d-k}{j} \right]$$



### Estrategias para mejorar la generación de itemsets frecuentes

- Reducir el número de candidatos (M)
  - Búsqueda completa:  $2^n$ .
  - Utilizar técnicas de poda para reducir M.
- Reducir el número transacciones (N)
  - A medida que aumenta el tamaño de los itemsets, reducir N.
- Reducir el número de comparaciones (NM)
  - Usar estructuras de datos eficientes para almacenar los candidatos o transacciones.

## Técnica de Asociación utilizando el Algoritmo Apriori.

### Pseudocódigo

---

**Algorithm 6.1** Frequent itemset generation of the *Apriori* algorithm.

---

```
1:  $k = 1$ .
2:  $F_k = \{ i \mid i \in I \wedge \sigma(\{i\}) \geq N \times \text{minsup} \}$ .    {Find all frequent 1-itemsets}
3: repeat
4:    $k = k + 1$ .
5:    $C_k = \text{apriori-gen}(F_{k-1})$ .    {Generate candidate itemsets}
6:   for each transaction  $t \in T$  do
7:      $C_t = \text{subset}(C_k, t)$ .    {Identify all candidates that belong to  $t$ }
8:     for each candidate itemset  $c \in C_t$  do
9:        $\sigma(c) = \sigma(c) + 1$ .    {Increment support count}
10:    end for
11:  end for
12:   $F_k = \{ c \mid c \in C_k \wedge \sigma(c) \geq N \times \text{minsup} \}$ .    {Extract the frequent  $k$ -itemsets}
13: until  $F_k = \emptyset$ 
14:  $\text{Result} = \bigcup F_k$ .
```

---

---

**Algorithm 6.2** Rule generation of the *Apriori* algorithm.

---

```
1: for each frequent  $k$ -itemset  $f_k, k \geq 2$  do
2:    $H_1 = \{ i \mid i \in f_k \}$     {1-item consequents of the rule.}
3:   call  $\text{ap-genrules}(f_k, H_1)$ 
4: end for
```

---

---

**Algorithm 6.3** Procedure  $\text{ap-genrules}(f_k, H_m)$ .

---

```
1:  $k = |f_k|$     {size of frequent itemset.}
2:  $m = |H_m|$     {size of rule consequent.}
3: if  $k > m + 1$  then
4:    $H_{m+1} = \text{apriori-gen}(H_m)$ .
5:   for each  $h_{m+1} \in H_{m+1}$  do
6:      $\text{conf} = \sigma(f_k) / \sigma(f_k - h_{m+1})$ .
7:     if  $\text{conf} \geq \text{minconf}$  then
8:       output the rule  $(f_k - h_{m+1}) \longrightarrow h_{m+1}$ .
9:     else
10:      delete  $h_{m+1}$  from  $H_{m+1}$ .
11:    end if
12:  end for
13:  call  $\text{ap-genrules}(f_k, H_{m+1})$ 
14: end if
```

---

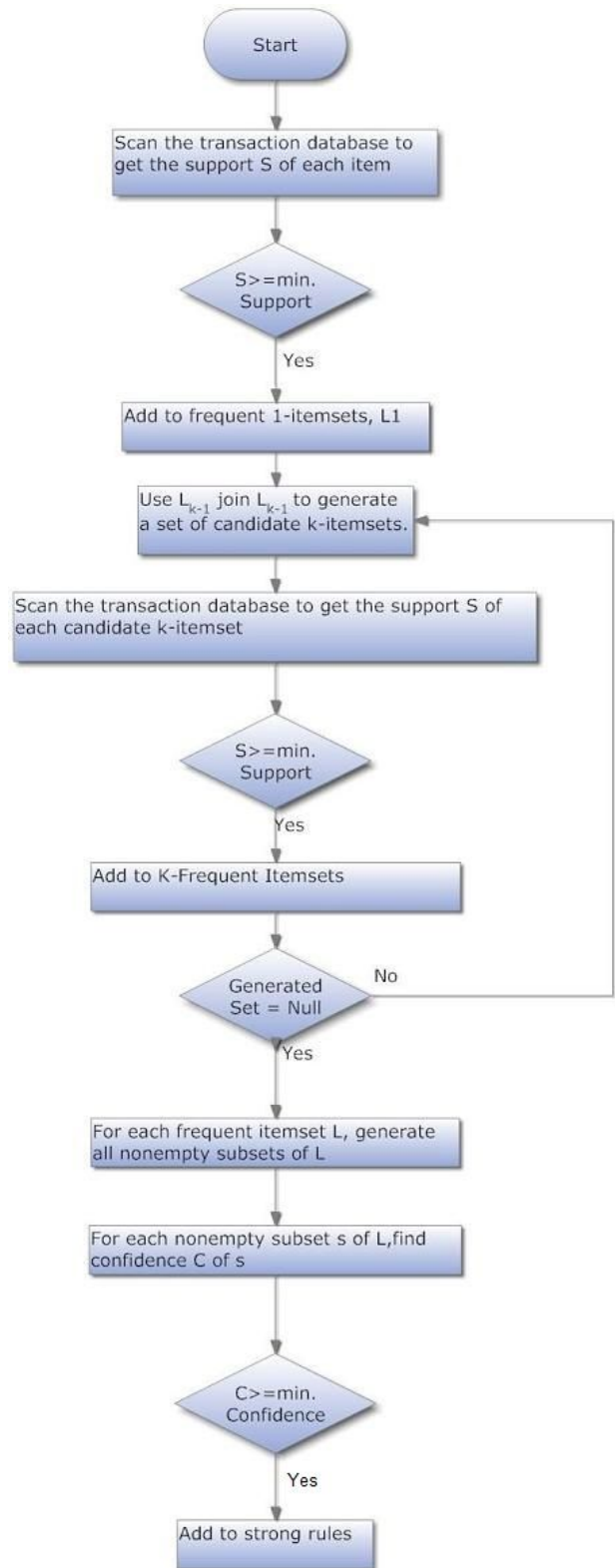
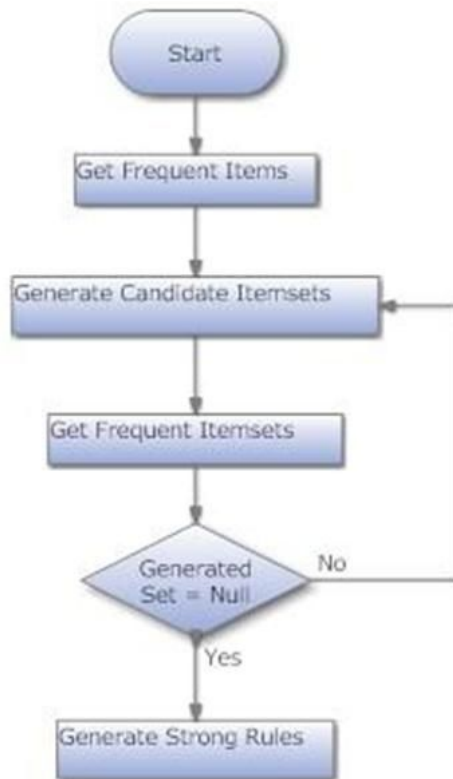
---

**Algorithm 6.4** Support counting using closed frequent itemsets.

---

```
1: Let  $C$  denote the set of closed frequent itemsets
2: Let  $k_{\max}$  denote the maximum size of closed frequent itemsets
3:  $F_{k_{\max}} = \{f | f \in C, |f| = k_{\max}\}$     {Find all frequent itemsets of size  $k_{\max}$ .}
4: for  $k = k_{\max} - 1$  downto 1 do
5:    $F_k = \{f | f \subset F_{k+1}, |f| = k\}$     {Find all frequent itemsets of size  $k$ .}
6:   for each  $f \in F_k$  do
7:     if  $f \notin C$  then
8:        $f.support = \max\{f'.support | f' \in F_{k+1}, f \subset f'\}$ 
9:     end if
10:  end for
11: end for
```

---



## Técnica de Asociación utilizando el Algoritmo Partition.

### Pseudocódigo

```
Initial Partitions ( examples, final_attributes)  
m ← Number (final_attributes)  
for i ← 1 ... m do  
  e [ ] ← { sorted examples of the attribute i }  
  partitions [ i ] ← Average (e [ ])  
  pointer ← Position (e [ ], partitions [ i ])  
  k ← pointer  
  while ek.class ≠ 1 // seeking next positive  
    if OR [ i ] > 1 then  
      k ← k + 1 // positive association  
    else  
      k ← k - 1 // negative association  
    end-if  
  end-while  
  if pointer ≠ k then  
    if OR [ i ] > 1 then  
      partitions [ i ] ← (ek + ek-1) / 2  
      // positive association  
    else  
      partitions [ i ] ← (ek + ek+1) / 2  
      // negative association  
    end-if  
  end-if  
end-for
```

### Características de Implementación

- Este algoritmo recorre la base de datos sólo dos veces. La primera vez, cada partición es minada independientemente para encontrar todos los conjuntos de ítems frecuentes en la partición y luego se mezclan éstos para generar el conjunto de los conjuntos de ítems candidatos. Muchos de éstos pueden ser falsos positivos, pero ninguno falso negativo.
- En la segunda pasada se cuenta la ocurrencia de cada candidato, aquellos cuyo soporte es mayor que el mínimo soporte especificado se retienen como conjuntos frecuentes.
- Emplea el mecanismo de intersección entre conjuntos para determinar el soporte de dichos conjuntos, en este caso cada ítem en una partición mantiene la lista de los identificadores de las transacciones que contienen a dicho ítem.

## Técnica de Predicción utilizando el Algoritmo de Bayes Naive.

Naive Bayes es una técnica de clasificación y predicción que construye modelos que predicen la probabilidad de posibles resultados. Naive Bayes utiliza datos históricos para encontrar asociaciones y relaciones y hacer predicciones.

---

## Pseudocódigo

---

---

### Algorithm 3.1: Fitting a naive Bayes classifier to binary features

---

```
1  $N_c = 0, N_{jc} = 0;$ 
2 for  $i = 1 : N$  do
3    $c = y_i$  // Class label of  $i$ 'th example;
4    $N_c := N_c + 1;$ 
5   for  $j = 1 : D$  do
6     if  $x_{ij} = 1$  then
7        $N_{jc} := N_{jc} + 1$ 
8  $\hat{\pi}_c = \frac{N_c}{N}, \hat{\theta}_{jc} = \frac{N_{jc}}{N_c}$ 
```

---

---

### Algorithm 3.2: Predicting with a naive bayes classifier for binary features

---

```
1 for  $i = 1 : N$  do
2   for  $c = 1 : C$  do
3      $L_{ic} = \log \hat{\pi}_c;$ 
4     for  $j = 1 : D$  do
5       if  $x_{ij} = 1$  then  $L_{ic} := L_{ic} + \log \hat{\theta}_{jc}$  else  $L_{ic} := L_{ic} + \log(1 - \hat{\theta}_{jc})$ 
6    $p_{ic} = \exp(L_{ic} - \text{logsumexp}(L_{i,:}));$ 
7    $\hat{y}_i = \text{argmax}_c p_{ic};$ 
```

---

---

## Ventajas

---

- Muy simple, fácil de implementar y rápido.
  - Si el supuesto de independencia condicional NB se mantiene, entonces convergerá más rápido que los modelos discriminatorios como la regresión logística.
  - Incluso si la suposición de NB no se cumple, funciona bien en la práctica.
  - Necesita menos datos de entrenamiento.
  - Altamente escalable. Se escala linealmente con el número de predictores y puntos de datos.
  - Se puede utilizar para problemas de clasificación binarios y de múltiples clases de vidrio.
  - Puede hacer predicciones probabilísticas.
  - Maneja datos continuos y discretos.
  - No es sensible a las características irrelevantes.
-

### Técnica de Predicción utilizando el Algoritmo de Regresión logística

Ventajas	Desventajas
<ul style="list-style-type: none"><li>● Es una técnica muy utilizada porque es muy eficiente.</li><li>● No requiere demasiados recursos informáticos.</li><li>● Es muy fácil de interpretar.</li><li>● No requiere que las funciones de entrada se amplíen.</li><li>● No requiere ningún ajuste.</li><li>● Es fácil de regularizar.</li><li>● Produce probabilidades pronosticadas bien calibradas.</li><li>● Funciona mejor cuando elimina atributos que no están relacionados con la variable de salida, así como atributos que son muy similares (correlacionados) entre sí.</li><li>● Es muy eficiente de entrenar.</li><li>● También es una buena línea de base que puede usar para medir el rendimiento de otros Algoritmos más complejos.</li></ul>	<ul style="list-style-type: none"><li>● No puede resolver problemas no lineales ya que su superficie de decisión es lineal.</li><li>● Cuenta con una alta dependencia de una correcta presentación de sus datos. Esto significa que la regresión logística no es una herramienta útil a menos que ya haya identificado todas las variables independientes importantes.</li><li>● Dado que su resultado es discreto, sólo puede predecir un resultado categórico.</li></ul>

### Técnica de Predicción utilizando el Algoritmo de Regresión logística

Un autómata finito es un modelo computacional que realiza un proceso de forma automática para producir una salida desde cierta entrada. Para implementar un análisis de datos usando este método, la idea será crear un autómata finito que compute automáticamente un proceso que siempre lleve a analizar unos datos de entrada.

#### Definición

Un autómata finito es una 5-tupla que se compone de la siguiente manera:

$(Q, \Sigma, q_0, \delta, F)$ , donde:

- $Q$  es el conjunto finito de estados por el que está compuesto el autómata.
- $\Sigma$  es el alfabeto con el que actúa (dominio).
- $q_0$  pertenece a  $Q$ .
- $\delta$  es la función que define las transiciones entre estados.
- $F$  es el conjunto que contiene a los estados del autómata en los que termina (estados de aceptación).

#### Tabla de Transición

Salida $q \in Q$	Símbolo $s \in \Sigma$	Llegada $q' \in Q$
S1	0	S2
S1	1	S1
S2	0	S1
S2	1	S2

### Tipos

#### Autómatas Finitos Deterministas

Es un autómata finito que además es un sistema determinista; es decir, para cada estado  $q \in Q$  en que se encuentre el autómata, y con cualquier símbolo  $s \in \Sigma$  del alfabeto (dominio) leído, existe siempre una transición posible de un estado a otro. Por ejemplo, el autómata a mano izquierda es un AFD.

#### Autómatas Finitos No Deterministas

Es un autómata finito que, a diferencia de los autómatas finitos deterministas (AFD), posee al menos un estado  $q \in Q$ , tal que para un símbolo  $s \in \Sigma$  del alfabeto, exista más de una transición posible.

Se concluye que lo más apropiado para nuestra solución es implementar un autómata que automatice un procedimiento que analice los datos que se tienen, sería un Autómata finito determinista. En efecto, este permite tener todos los posibles resultados para cualquier tipo de dato que entre (símbolo), a diferencia del AFN.

### Implementación

Se tiene un conjunto (C) de datos. El estado inicial (podría ser  $q_1$  en el ejemplo) del autómata representaría el subconjunto (S0) de datos con el que comienzo a analizar.

En seguida entra otro subconjunto de datos (S1) tal que S1 pertenece a C, es decir los subconjuntos son los datos que voy a ir analizando. El AFD, que por el momento se encuentra en el estado inicial, decide que procedimiento realizar dependiendo de la clasificación que tenga S2. (esta clasificación estaría representada por los símbolos posibles que tenga el autómata). En el diagrama a continuación estos serían los símbolos: 0 y 1.

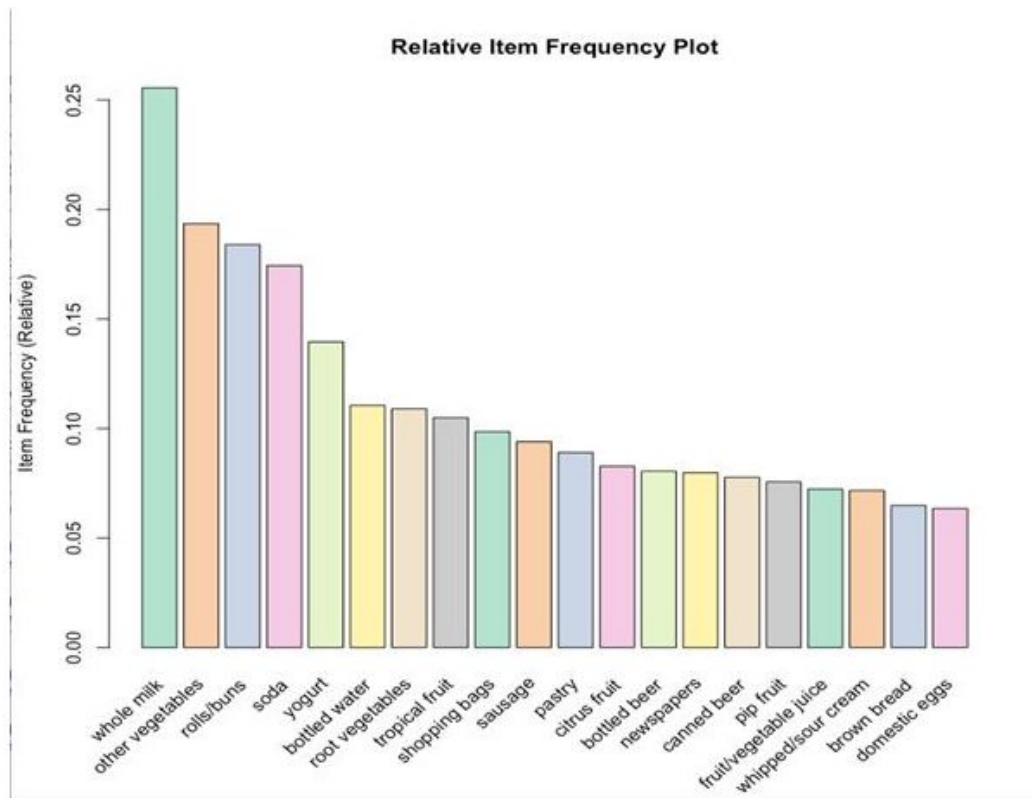
Se continuaría de esta forma hasta que se llegue a un estado final. En el diagrama sería:  $q_0$

Por ejemplo. En el caso hipotético de que se tenga el primer subconjunto S0 ( $q_1$ ) y entre otro subconjunto de datos S1 de tipo 1, entonces el autómata procedería a realizar las operaciones de análisis que lleven directamente al estado final  $q_0$ . En el cual la información que se tiene se puede considerar un análisis relevante.



## Presentación de resultados mediante el uso de histogramas para mostrar la frecuencia de los itemsets.

### Diseño



Se utilizan los componentes de Visual Studio para crear histogramas con la frecuencia de los itemsets de un tamaño dado.

## Presentación de resultados mediante el uso de tablas para mostrar las asociaciones encontradas.

### Diseño

1. Pasillo de comidas - Leche, Huevos, Pan
2. Pasillo de licores - Licor, Red/Blush Vino, Cerveza, Soda
3. Pasillo de desayuno - Cereal Yogurt, Arroz, Avena