

OPEN DATA : NASA, twitter, Datos Argentina, etc (via Open-API)

Data Mining : Knowledge Discovery from Data ~~Minería~~

↳ Extracción de Patrones Interesantes y conocimiento a partir de grandes cantidades de datos.

- No triviales
- Implícitos
- Previamente desconocidos
- Potencialmente útiles

Knowledge Discovery in Databases (KDD)

• Proceso interactivo e iterativo

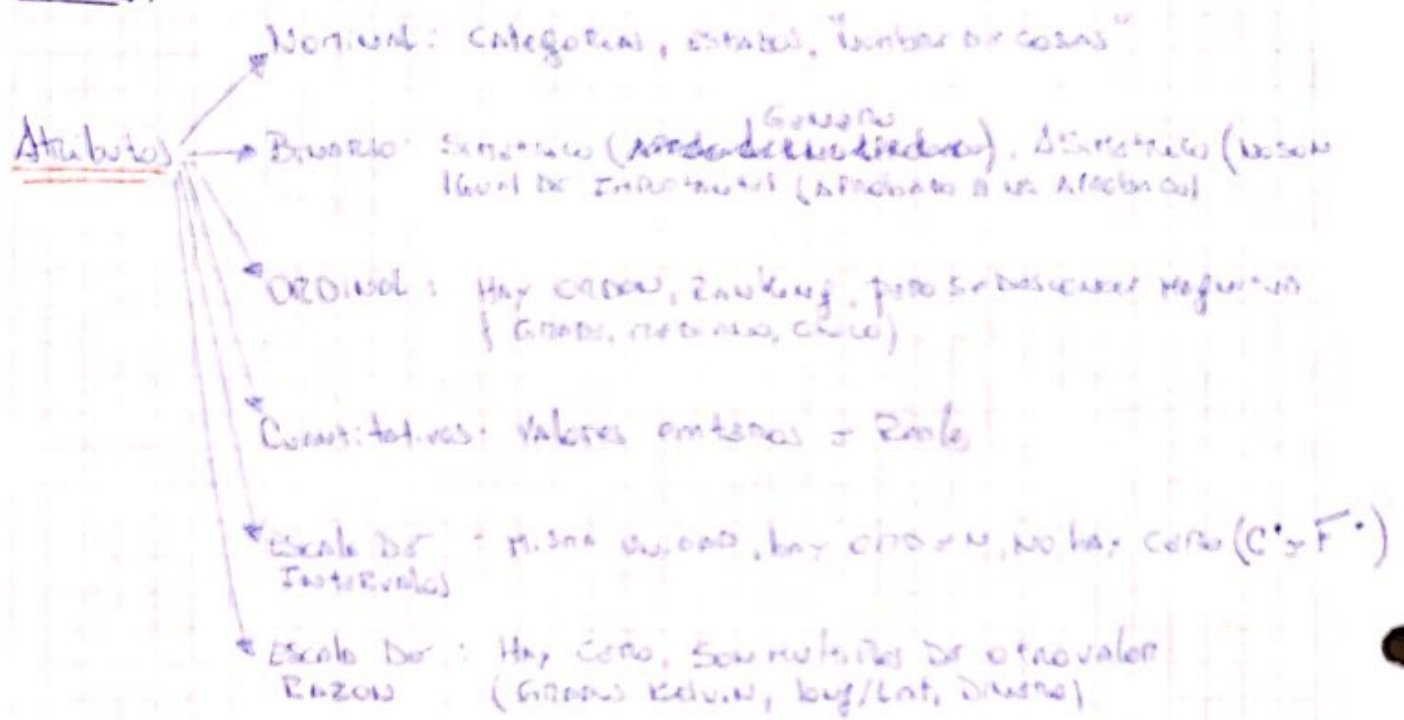
- ① Entender el Dominio

Donde quiero aplicarlo, ¿Cuál es el problema, cuáles los objetivos?
- ② Seleccionar el Conjunto de Datos

Crear el dataset (Averiguar de los datos, obtener adicionales, integrar con otras fuentes)
- ③ Preprocesamiento

Mejorar la fiabilidad de los datos (datos faltantes, Eliminar ruido, outliers (Ojo: no todo es ruido))
- ④ TRANSFORMACIÓN
  - Reducir Dimensionalidad (PCA, X<sup>2</sup>, Correlación)
  - Suavizado (Discretizar, Binning, medias)
  - Agregación, Generalización, Normalización (Z-Score)
  - Construir nuevos atributos
- ⑤ TAREA  
Selección del Método
  - Se extraen patrones y tendencias
  - Predicción, clasificación, Clustering, Asociación
- ⑥ Algoritmo  
Selección del Algoritmo
  - Selección, parámetros, técnica de aprendizaje.
  - Precision vs Capacidad de explicar
- ⑦ Utilización
  - Ejecutar el algoritmo varias veces hasta obtener resultados satisfactorios (Ajustar parámetros)
- ⑧ Evaluación
  - Evaluar e Interpretar Resultados
  - Medidas Cuantitativas
- ⑨ Utilizar el conocimiento
  - Pasar de datos de laboratorio al mundo real

## DM-2



Discretos: Un solo conjunto de valores. Punto o número contable (id. cliente, cp, profesión, etc)

Binarizados como particular de Discretos

Continuos: Vars. Reales como valores de atributo (temp, altura, peso)  
# Línea de Dptos

Tendencia Central

MEDIA - MODA - MEDIANA - SIMETRÍA / SESGO

Dispersión

Cuantiles, Outliers, Boxplot, Varianza,  $\sigma$

Histogramas, Geo Plot Cuantiles, Q-Plot, Scatter Plot, Plot, Rotas, Carta de Control

### Pre-Processamiento

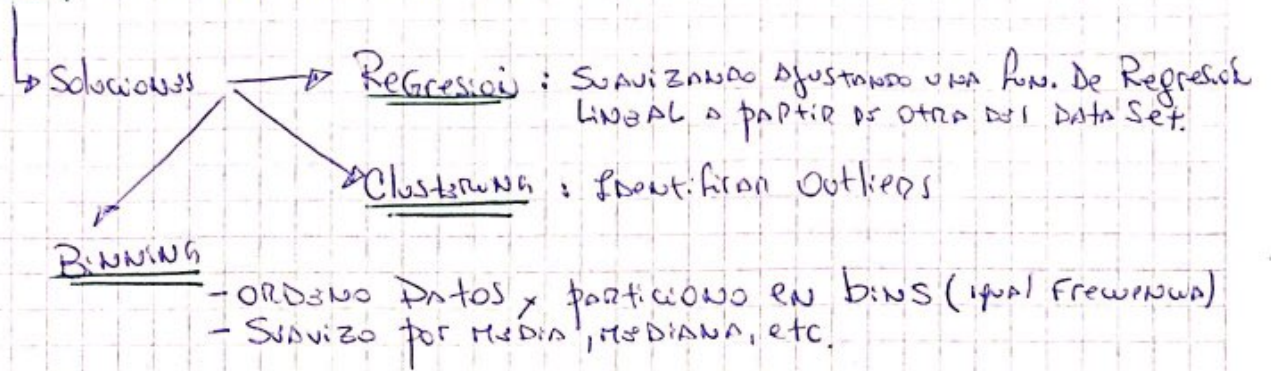
- Campos desolados o faltantes
- Valores Faltantes
- Outliers
- Formato Incorrecto
- Valores no consistentes

### Calidad

- Precisión, Integridad, Consistencia, Rentabilidad, Confabilidad, Interpretabilidad



# Ruido: Error Aleatorio & Varianza en una variable med.



Robustez: Capacidad de construcción modelos insensible a datos corruptos

Integración de Datos: Combinar  
Identificar  
Resolver Conflictos

↓  
ETL

Datos Redundantes los detecto con Analisis de Correlación  
Analisis de Covarianza

## DM-3: OUTLIERS

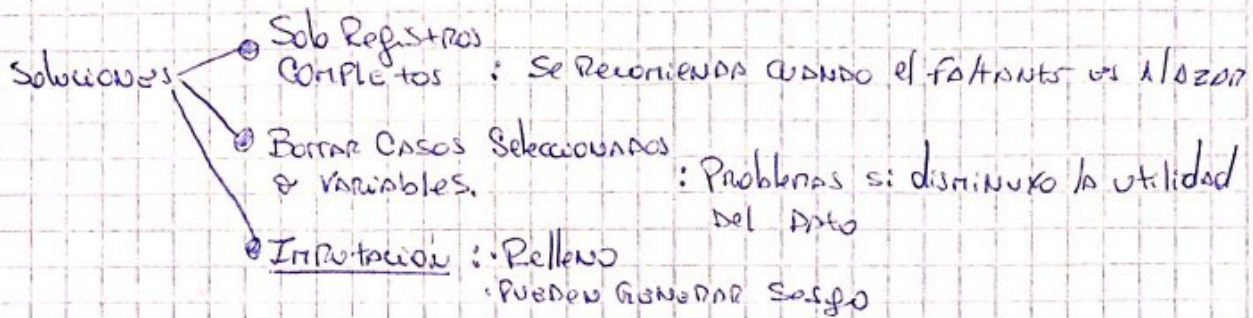
UNIVARIADO	MULTIVARIADO
<ul style="list-style-type: none"> <li>• Valores Atípicos q' se encuentran en una simple variable.</li> <li>• Buenos para detectar extremos pero no en otros casos.</li> </ul>	<ul style="list-style-type: none"> <li>• Se encuentran en un espacio N-DIMENSIONAL</li> <li>• Para detectarlos hay que ajustar el modelo</li> </ul>
<ul style="list-style-type: none"> <li>• BoxPlot, z-score, IQR</li> </ul>	<ul style="list-style-type: none"> <li>• Clustering</li> <li>• LOF</li> <li>• Isolation Forest (Ensamblados)</li> </ul>
<p><u>Boxplots</u>: Problema cuando los Atípicos no son muy grandes o muy chicos.</p> <ul style="list-style-type: none"> <li>• Mover tamaño de Bigotes y tamaño cajas</li> </ul>	<p><u>LOF: Local Outlier Factor</u></p> <ul style="list-style-type: none"> <li>• Basado en densidad y utiliza vecindades locales.</li> </ul>
<p><u>Z-SCORES</u>: # desviaciones estándar tiene una observación</p> $Z_i = \frac{x_i - \mu}{\sigma}$ <p>Se fija un umbral</p>	<p> <math>LOF_2(d) = 2.004</math> → Valor Atípico  <math>LOF_2(b) = 5.336</math>  <math>LOF_2(a) = 3.507</math>  <math>LOF_2(c) = 3.501</math> </p> <p>Hay que definir umbral de cutoffs.</p>



## DM-4: MISSING DATA

Razones: Mal diseño, Negativa a Responder, Respuestas Inaplicables.

Problemas: Dificultan el análisis, introducir sesgos, Pérdida de eficiencia



### Imputación

- Sustitución de Casos: Se reemplaza con valores no observados (EXPERTO DE DATOS)
- Sustitución de Media: Reemplazo con promedio (SI NO SON NORMALES USAR MEDIA MEDIANA)
- Cold-Deck:
  - Selecciona valores y usar relaciones obtenidas de fuentes distintas de la base actual.
  - Se sustituye un valor constante derivado de otras fuentes
  - misma desventaja que la Media (VAR ATENUADA, DISTORSION DISTRIBUCION, REPETICION DE VALOR)
- Hot-Deck:
  - Reemplaza los faltantes con los más similares
  - ⊕ Simple, funciona con conjunto completo
  - ⊖ Difícil determinar que es similar
- Regresiones:
  - Predecir valores faltantes
  - Simples o múltiples
  - Se identifican INDEP y DEP
- Imputaciones Múltiples (MICO):
  - Supone que el origen del faltante es RANEO (MAR) (la proba de q' falta depende sob de los valores observados y no de los no observados)
  - Aplicar sólo a DATOS NO MAR pueden generar Sesgo!



Mye Steps: Se Repite el ciclo varias veces (10 aprox)

- 70 times
- ↳ ① Imputacion Simple (Reemplaza por Media o Moda)
  - ↳ ② UNA de las Faltantes originales se las vuelve a poner como Faltantes
  - ↳ ③ Predice esa Faltante con el Resto
  - ↳ ④ Se Reemplaza por las predicciones

④ No produce sesgo  
/ Se puede usar en cualquier Analisis  
/ Fácil de usar

⑤ Costoso Computacionalmente  
/ Genera un Dataset completo por iteracion  
/ Hay que pensar el modelo de Imputacion Ademas del de Analisis

## DM-5: Reduccion De la Dimensionalidad

- Reducir el volumen pero aun produciendo casi los mismos Resultados Analiticos.

- Maldicion: - Mas Dimensionalidad los datos se vuelven mas Ralos.  
- El Analisis de Atidicos se vuelve mas significativo  
- Las combinaciones de subespacios crece exponencialmente.

① Eliminar Columnas C/ Datos Faltantes :  
/ Atributos con menos del x% de valores se eliminan  
/ Se busca valor de corte  
/ Aplica Para Numericas y Categoricas

② Low Variance Filter : Mido var de una columna  
/ Valores constantes  $\Rightarrow$  var = 0  $\Rightarrow$  No ayudan a Discriminar.  
/ Remueve Atributos que estan por debajo de un umbral de Var  
(Hay q' normalizar rangos de datos)

③  $\chi^2$  :  
/ Mide Dependencia.  
/ Sacamos variables con mayor probabilidad de ser independientes de la clase (irrelevantes para clasificar)  
/ 1 mas Relacionado a 0 menos.

④ Componentes Principales (PCA) :  
/ Encuentra una proyeccion que capture la mayor cantidad de variacion entre los Datos  
/ Los Datos originales se proyectan en espacios mas pequeños  
/ Busco Autovalores de la Mat Cov.



- Reducing Highly Correlated Columns
  - Atributos Correlacionados introducen Redundancia.
  - Se puede eliminar sin disminuir drásticamente la cant. de info.
  - Eliminar Pares correlacionados a partir de la matriz de correlaciones.
  - Variables continuas & discretas (Pearson,  $\chi^2$  de Pearson)
- Variables Importantes (RF)
  - Derivadas de la salida de ensemble Random y Forest.
  - Arbols  $\Rightarrow$  Medidas internas de importancia.
  - Calcular la importancia de cada variable y con eso el promedio de importancia de cada una.
- Backward Feature Elimination
  - Usa Algoritmo de AA para medir como disminuye el error al quitar un atributo.
  - Alto número de iteraciones (Ⓢ computo)
- Forward Feature Construction
  - Conjunto vacío a va agregando el mejor atributo determinada

## DM-6: FEATURE ENGINEERING

- Transformación de Datos, mejoran el Rendimiento. transformando el Feature Space
- Construcción de variables por
  - Discretización
  - Normalización
  - Binning
- Evaluación de las Transformaciones

### NORMALIZACIÓN

- Mover los Features a un rango mas pequeño
- Ayuda a evitar que atributos con valores grandes tengan mas peso

MIN-MAX: Cuanto mas grande es el valor actual del valor mínimo y escala x rango  
 $\in [0, 1]$

$$X'_{min} = \frac{X - \min(x)}{\text{Rango}(x)} = \frac{X - \min(x)}{\max(x) - \min(x)}$$



Z-SCORES: Se normaliza  $X$  en base a la media y el desvío estándar  $X$ .

$$Z\text{-Score} = \frac{X - \text{MEDIA}(X)}{Sd(X)}$$

- Util cuando no conozco el min ni el max
- Util cuando outliers dominan el min || max

Decimal Scaling:

Asegura que cada valor normalizado  $\in [-1, 1]$

$$X_{\text{decimal}} = \frac{X}{10^d}$$

$d$  número de dígitos en los valores de la variable con valor absoluto mas grandes.

Escalados Robustos: muchos valores atípicos, usar media y var no funciona bien

USO ESTIMACIONES MAS SOLIDAS PARA EL CENTRO Y RANGO

MEDIANA (o algún percentil) e IQR

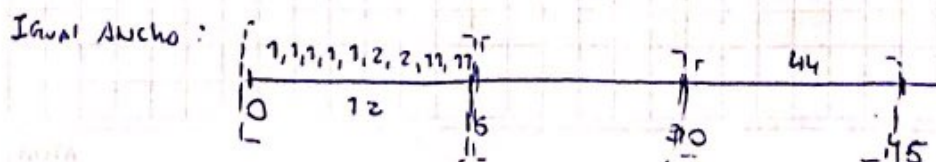
Discretización

- Dividir Rango continuo en intervalos
- Por Class  $\Rightarrow$  SUPERVISADA (ENTROPÍA)
- Sin Class  $\Rightarrow$  NO-SUPERVISADA (BINNING, RF, Quantiles, Math)

Binning

- TOP-DOWN
- Cantidad específica de bins
- Agrupamiento
  - $\rightarrow$  Por Frecuencia: misma # de observaciones por bin
  - $\rightarrow$  Igual Ancho: Definir Rango e intervalos por bin.
- NO-SUPERVISADA
- Para cada Agrupamiento puede Reemplazarse Por Media, Mediana, etc.

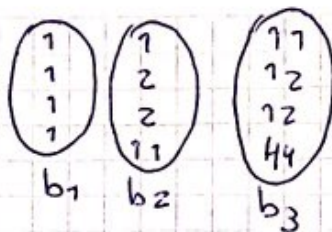
$X = \{1, 1, 1, 1, 1, 2, 2, 11, 11, 12, 44\}$  3 categorías



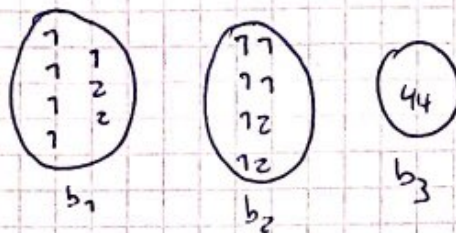


Igual  
Frecuencia

$m = 12$   
 $\text{bins} = 3$   
 $\frac{m}{\text{bins}} = 4$



K-MEANS: Identifican lo que  
parece ser intuitivamente  
Correcto



Rank

El Ranking de un número es su tamaño relativo a otros valores de una variable numérica.

- Mismos valores reciben el mismo Rank, los duplicados afectan las posiciones de los posteriores.
- Solo si para el mismo valor puede tener dos ranks en diferentes listas.

Quantiles

- Útiles para como Rank puede un  $X$  tener distintos cuantiles si la lista cambia.

Math  
Function

$\text{Floor}(\log(x))$ : efectivo binning para las distribuciones sesgadas

Entropía

- SUPERVISADA y SEPARACION TOP-DOWN
- Explora y Determina Split-Point

$$\text{Entropía Class } E(S) = \sum p_i \ln p_i$$

$$\text{Entropía Class n Split-Point: } E(S, A) = \sum \frac{|S_v|}{|S|} E(S_v)$$

$$\text{Information-Gain} = E(S) - E(S, A)$$

Variable Flag

- Regresión Requiere numéricas.
- Si tengo categóricas lo puedo recodificar en una o mas Dummy y Flag.

$$\begin{cases} \text{if sex = female then sex.Flag = 0} \\ \text{if sex = male then sex.Flag = 1} \end{cases}$$



## DM-7: No-SQL

- Heterogeneidad de los datos y la cantidad llevaron al límite los RDBMS
- Conceptos que permitan el procesamiento poniendo foco en la performance, confiabilidad y agilidad.

RDBMS: controla las tx a través del uso de la Atomicidad, consistencia, independencia y persistencia.

**A**tonicity: La tx no se puede separar, de hacer o no se hace.

**C**onsistency: La tx pasa la BD de un estado consistente a otro consistente.

**I**solation: La tx se ejecuta independientemente.

**D**urability: Los cambios son permanentes, se persisten.

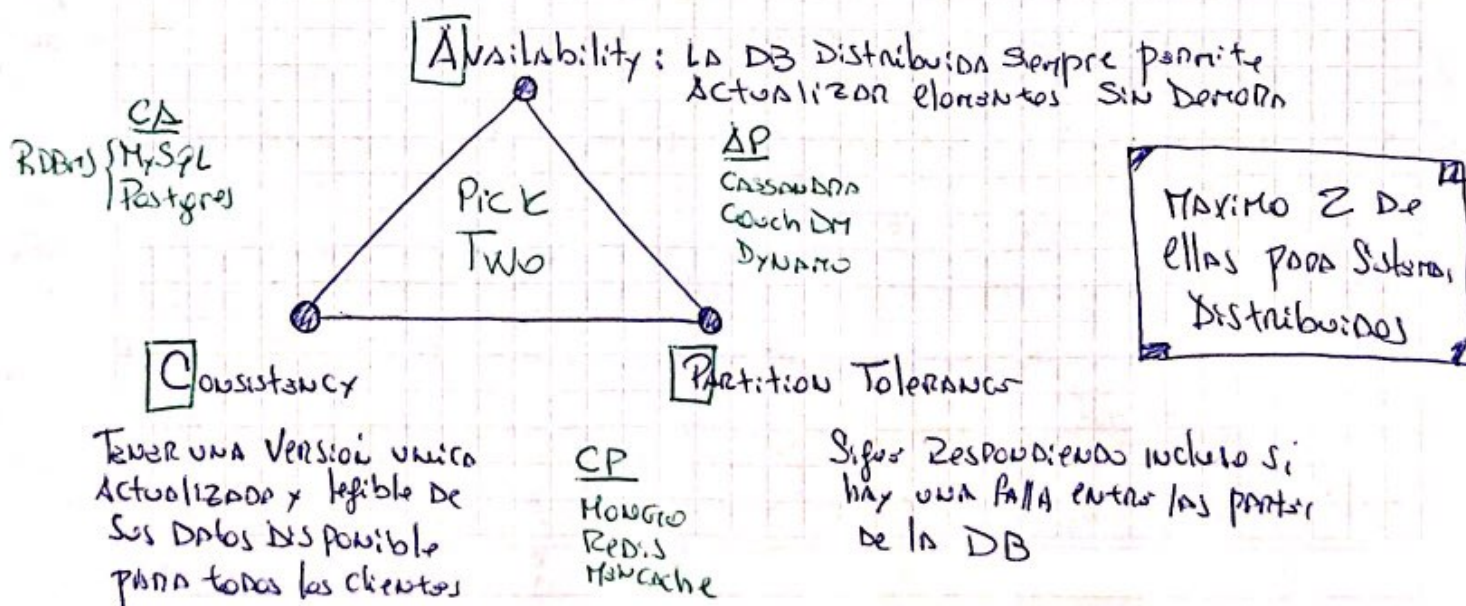
ACID hace foco en Consistencia e Integridad. Bloquea Recursos

**B**asic **A**vailability: Permite que los sistemas temporalmente sean inconsistentes para que las tx sean manejables.

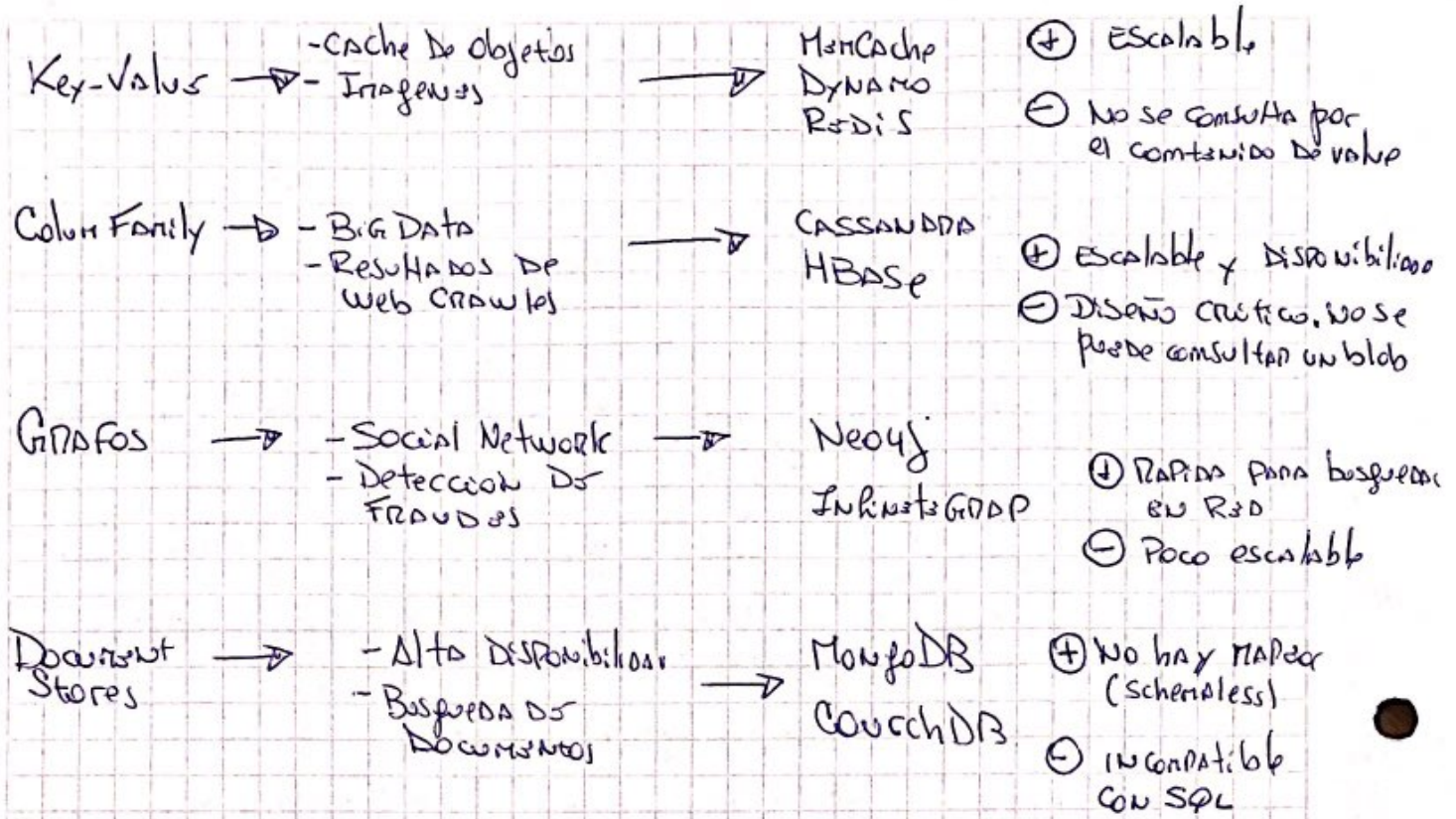
**S**oft-state: Permitir inconsistencias temporales, los datos pueden cambiar mientras se usan.

**E**ventual Consistency: Cuando toda la lógica se ejecuta se llega a un estado consistente.

BASE promueve el no-bloqueo ante la consistencia. Disponibilidad







## DM-8: Datos NO ESTRUCTURADOS

### TEX MINING

- Lexicon: Conjunto completo de las distintas palabras para definir el CORPUS
- CORPUS: Conjunto de datos que corresponde a una colección de documentos.

• Bolsa D/palabras →  $\frac{\text{Terminos}}{\text{Totales}} = \text{Frecuencia}$

↓  
Cual es la proba de sacar = frecuencia una palabra  
(se tratan como dimensiones o features)

- STOP-WORDS: Palabras frecuentes que no ayudan y se fuerzan sacar.

### • MATRIZ TERMINO DOCUMENTO

↳ Filas son los documentos

↳ Columnas terminos para analisis

↳ Matriz llena de ceros. Alta dimensión, dispersa x no negativa.

	Hola	Mundo	Manolita
Hola Mundo	1	1	
Hola Manolita	1		1

- Stemming: extracción de raíz morfológica

Stu	dies	electu	electu
Stem	Subj	electu	electu
		electu	electu
		electu	electu



- Lemma: Usa la parte específica del Habla para Determinar la Raíz.
- Dependencia del IDIOMA.

FORMA	INFO	LEMMA
NINAS	Femenino Plural de NIÑO	NIÑO

## • NER: Named Entity Recognition

- Identificar y clasificar expresiones de un texto que hacen Referencia a personas, lugares, marcas, fechas, horas, "ENTIDADES"

## • Part-Of-Speech

Tagging: Averiguar que son Sustantivos, Verbos, etc.

## DM-9: Reglas De Asociación

- DADO UN CONJUNTO DE TX ENCONTRAR REGLAS QUE PUEDAN PREDICIR LA OCURRENCIA DE UN ITEM BASADO EN LA PRESENCIA DE OTROS
- LA IMPLICACION INDICA CO-OCURRENCIA (NO CASUALIDAD).

Item set: Colección de 1 o MAS Items

K-Itemset: un itemset con K elementos

Support Count ( $\kappa$ ): CANTIDAD DE OCURRENCIAS DE UN itemset.

Support ( $s$ ): Fracción de tx que contienen a un itemset

$$s = \frac{\kappa(x)}{|T|}$$

Itemset Frequent: un itemset cuyo Support ( $s$ ) es MAYOR o igual al umbral de MinSup.



Tid	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

$$\sigma(\{\text{Milk, Diaper}\}) = 3$$

$$S(\{\text{Milk, Diaper}\}) = \frac{\sigma(\{\text{Milk, Diaper}\})}{|T_r|} = \frac{3}{5}$$

$$\sigma(\{\text{Beer}\}) = 2 \Rightarrow S(\{\text{Beer}\}) = \frac{2}{5}$$

Supongo  $\text{minSup} = 0.6$

$$S(\{\text{Beer}\}) = 0.4 \quad (\text{X})$$

$$S(\{\text{Milk, Diaper}\}) = 0.6$$

$$S(\{\text{Bread}\}) = \frac{4}{5} = 0.8$$

$$S(\{\text{Coke}\}) = \frac{2}{5} = 0.4 \quad (\text{X})$$

$$S(\{\text{Milk}\}) = \frac{4}{5} = 0.8$$

Itemset  
Frecuente

$$S(\{\text{Milk, Diaper, Beer}\}) = \frac{2}{5} = 0.4 \quad (\text{X})$$

Regla De Asociación:

Expresión  $X \rightarrow Y$  Donde  $X$  e  $Y$  son Itemset frecuentes.

Métricas

Support:  $S$  = Fracción de transacciones que contienen a  $X$  e  $Y$

Confidence:  $C$  : Mide con qué frecuencia  $X$  aparece en las transacciones en las que también aparece  $Y$

$$C(X \rightarrow Y) = \frac{\sigma(X \cup Y)}{\sigma(X)}$$

$$\{\text{Milk, Diaper}\} \rightarrow \{\text{Beer}\}$$

$$S = \frac{\sigma(\{\text{Milk, Diaper, Beer}\})}{|T_r|}$$

$$S = \frac{2}{5} = 0.4$$



$$C = \frac{\sigma(\{milk, Diaper, Beer\})}{\sigma(milk, Diaper)} = \frac{\frac{2}{5}}{\frac{3}{5}} = \frac{2}{3} = 0.67 = \frac{0.4}{0.6}$$

Reglas Creadas a Partir del mismo Itemset tienen mismo Soporte pero pueden tener Distinta Confianza.

ENCONTRAR TODAS LAS REGLAS /  $\text{Soporte} \geq \text{MinSup}$   
 $\text{Confianza} \geq \text{MinConf}$  } ENCONTRAR TODAS (FUERZA bruta)  
 ES MUY COSTOSO!

- ① Generar todos los Itemsets /  $\text{Support} \geq \text{MinSup}$
  - ② Generar a partir de los Itemsets Subconjuntos con las Reglas /  $\text{Confianza} \geq \text{MinConf}$
- Dados  $d$  items hay  $2^d$  Itemsets candidatos

Fuerza bruta  $\sim O(NMw)$ ,  $M = 2^d$ ,  $N = tx$

- Reducir candidatos ( $M$ ) mediante reglas
- Reducir comparaciones ( $MN$ ), utilizando estructuras eficientes
- Reducir # de  $T_x$ : en DHP

Principio Apriori: Reduce candidatos

"Si un itemset es frecuente  $\Rightarrow$  todos sus subsets deben serlo"

Propiedad  
Antimonotonia

$$\forall x, y: (x \subseteq y) \Rightarrow \text{Support}(x) \geq \text{Support}(y)$$

El soporte de  $y$  nunca supera el de sus subsets

• Si  $x$  itemset no satisface  $\text{MinSupport} \Rightarrow x$  no es frecuente

Si Agrego  $x_2$  a  $x$ :  $(x \cup x_2) \Rightarrow$  el resultado no puede ser mas frecuente

luego

$$\text{Support}(x \cup x_2) < \text{MinSupp}$$



APRORI: Factores que afectan la complejidad

↳ Elegir Min-Support

↳ Bajarlo incrementa la cantidad de itemsets frecuentes.

(puede incrementar candidatos y la longitud máxima de itemsets frecuentes)

↳ Dimensionalidad del Dataset (# Items)

↳ MAS espacio para almacenar el count Supp de cada item.

↳ Computacionalmente MAS costoso

↳ Tamaño de la BD (# transacciones)

↳ APRORI hace varias pasadas, el tiempo de ejecución aumenta con el número de Tx.

Tips Especiales:

es útil seleccionar un conjunto representativo de itemsets y a partir de ellos derivar conjuntos de itemsets frecuentes.

Maximal Frequent Itemset (MFI)

- Un itemset es Maximal si ninguno de sus superset es frecuente.
- Proporcionan una representación compacta del conjunto de elementos frecuentes.

Closed Frequent Itemset (CFI)

- Un itemset es closed si ninguno de sus inmediatos superset tienen el mismo support que el itemset.
- Representación mínima de los itemset sin perder info de supports.
- CFI permite Remover Reglas Redundantes.
- Se puede usar para determinar el support de un item que no es closed.

Reglas Redundantes

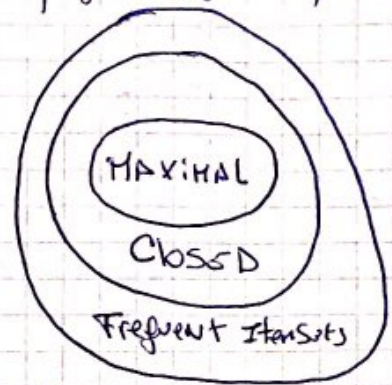
$X \rightarrow Y$  es Redundante si  $\exists X' \rightarrow Y', X \subseteq X' \wedge Y \subseteq Y' /$  Support Confianza son iguales.

CFI las evita



Todos Los Max Frequent Itemsets son Closed, porque ninguno Superset tiene el mismo Support.

MAXIMAL es Closed



La confianza es limitada, Reglas con alta confianza pueden ocurrir de CASUALIDAD.

LIFT

- Asumiendo INDEPENDENCIA
- $X \rightarrow Y$ , el Lift es la confianza de la Regla Dividida la confianza esperada.
- mide que tan lejos de la independencia estan  $X$  e  $Y$

$$\text{Lift}(X \rightarrow Y) = \frac{\text{Conf}(X \rightarrow Y)}{\text{Supp}(Y)}$$

CERCANOS A 1 implican que  $X$  e  $Y$  son indep x por lo tanto la Regla es Poco interesante.



## DM-10: Frequent Pattern Growth (FP-Growth) - Algorithm

Apriori	FP-Growth
<ul style="list-style-type: none"><li>• Utiliza un enfoque de Generar y Probar: Genera itemsets candidatos y prueba si son frecuentes.</li><li>• Generar itemsets candidatos es costoso</li><li>• El control de soporte es costoso</li></ul>	<ul style="list-style-type: none"><li>• Descubrir itemsets sin generar itemsets candidatos</li><li>• Paso-1: Crea estructura compacta: FP-Tree<ul style="list-style-type: none"><li>↳ Haciendo 2 pasadas sobre el conjunto de datos</li></ul></li><li>• Paso-2: Extrae itemsets frecuentes directamente desde el FP-Tree.<ul style="list-style-type: none"><li>↳ A través del recorrido por el FP-Tree</li></ul></li></ul>



# DM-11: SISTEMAS DE RECOMENDACION

## MATRIZ

UTILIDAD: Los valores representan grados de preferencia de los usuarios sobre los items.

USUARIOS	Items					
	i1	i2	i3	i4	i5	i6
U1	?	1	?	5	?	?
U2	4	?	2	3	?	?
U3	?	4	?	8	5	3
U4	?	?	5	?	3	?

Explícita

USUARIOS	Items					
	i1	i2	i3	i4	i5	i6
U1	0	1	0	1	0	0
U2	1	0	1	1	0	0
U3	0	1	0	0	1	1
U4	0	0	1	0	1	0

Implícita

Los valores representan grados de preferencia de los usuarios sobre los items.

## ETAPAS RECOMENDACION PARA USUARIO "U"

### E1: Predicción

→ El sistema asigna un score a cada item "i" observado por "U".

### E2: Recomendación

→ Se genera una lista de items ordenada por valor de score y se recomiendan los primeros k elementos de esta lista.

## Clasificación De Algoritmos

- Populacional
- Basados en Contenido
- Asociación De Productos
- Filtrado Colaborativo
- Híbrido y Ensamblados

## Asociación De Productos

- Cada User es una Tx
- Se calculan Reglas entre todos los pares de items
- Se utiliza una matriz de Asociación (Support, Lift, etc)
- Se genera matriz cuadrada  $S$  tamaño  $|I| \times |I|$
- Si la matriz es un índice de similitud  $\Rightarrow$  Matriz Similitud y es Mat. de Similitud

Usuario U  
Interactúa con  
i1

Predicción

Score(i1) = 0.1  
Score(i2) = 0.2  
Score(i3) = 0.6  
Score(i5) = 0  
Score(i6) = 0.9

	Items					
	i1	i2	i3	i4	i5	i6
i1	1	2	5	1	2	9
i2	2	1	0.5	2	4	0
i3	5	0.5	1	6	2	0
i4	1	2	6	1	0	9
i5	2	4	2	0	1	7
i6	9	0	0	9	7	1

Recomendación:

[i6, i3, i2, i1, i5]



## Similitud Jaccard

↳ Los items se Representan por conjuntos de usuarios

$$i_1 = \{u_1, u_2, u_5\}$$

$$i_2 = \{u_1, u_6\}$$

$$\text{Sim}(A, B) = \frac{\text{Supp}(A \& B)}{\text{Supp}(A) + \text{Supp}(B) - \text{Supp}(A \& B)}$$

$$\text{Sim}(i_1, i_2) = \frac{|\{u_1\}|}{|\{u_1, u_2, u_5, u_6\}|} = \frac{1}{4} = 0.25$$

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{|A \cap B|}{|A| + |B| - |A \cap B|}$$

## Similitud Coseno

↳ Los items se Representan como vectores con tantas dimensiones como usuarios existan.

$$i_1 = \{u_1, u_2, u_5\} \rightarrow [1, 1, 0, 0, 1, 0]$$

$$i_2 = \{u_1, u_6\} \rightarrow [1, 0, 0, 0, 0, 1]$$

$$\text{Sim}(A, B) = \frac{A \cdot B}{\|A\| \cdot \|B\|}$$

$$\text{Sim}(A, B) = \frac{\text{Supp}(A \& B)}{\sqrt{\text{Supp}(A) * \text{Supp}(B)}}$$

$$\text{Sim}(i_1, i_2) = \frac{1}{\sqrt{2 * 3}} = 0.41$$



## Filtrado Colaborativo

### Predicciones

### USER to USER

- Dado un usuario  $U$
- Se calculan similitudes contra todos los usuarios
- Se eligen a los  $k$  vecinos mas cercanos
- Por cada item  $i$  desconocido por  $U$  se calcula como score a la suma de todas las indices de similitud de los  $k$  vecinos que contengan  $i$ .

### Items to items

- Generalización de las Recomendaciones de Asociación de productos.
- Para las predicciones se suman todas las similitudes a los items conocidos por  $U$ .

$$U_1 = \{i_1, i_3, i_4\}$$

$$\text{Score}(i_2) = \text{Sim}(i_2, i_1) + \text{Sim}(i_2, i_3) + \text{Sim}(i_2, i_4)$$

$$\text{Score}(i_2) = 0.2 + 0.05 + 0.2 = \boxed{0.45}$$

$$\text{Score}(i_5) = 0.2 + 0.2 + 0 = \boxed{0.4}$$

$$\text{Score}(i_6) = 0.9 + 0 + 0.9 = \boxed{1.8}$$

## Factorización de Matrices

### Método De Reducción De Dimensionalidad

- SV
- Gradient Descent.
- Se descomponen la matriz de utilidad en nuevas dimensiones
- Estas Dimensiones Latentes/Ocultas captan  $\neq$  características de los items o usuarios
- Los items y los usuarios quedan representados en este espacio Latente.
- Al estar representados en un mismo espacio pueden calcularse directamente la similitud/distancia entre un usuario  $x$  un item,