

Qué es DataMining?

Ing. Gustavo Markel

gmarkel@gmail.com

Lic. Cecilia Ruz

ruz.cecilia@gmail.com

Agenda

- Qué es Data Mining?
- Cómo se integra en el proceso de Descubrimiento del conocimiento?
- Funcionalidades del Data Mining
- Técnicas
 - Supervisadas
 - Redes neuronales
 - Árboles
 - Regresión
 - No supervisadas
 - Clustering
 - Reglas de Asociación



Agenda

- **Qué es Data Mining?**
- **Cómo se integra en el proceso de Descubrimiento del conocimiento?**
- **Funcionalidades del Data Mining**
- **Técnicas**
 - **Supervisadas**
 - Redes neuronales
 - Árboles
 - Regresión
 - **No supervisadas**
 - Clustering
 - Reglas de Asociación



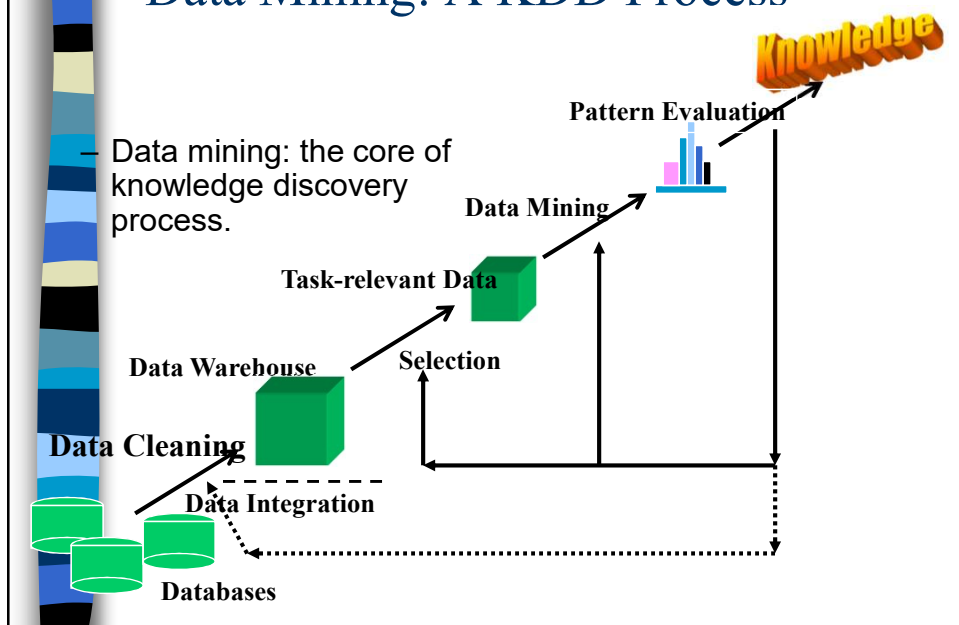
Qué es Data Mining?

- “Es la extracción de patrones o información interesante (no trivial, implícita, previamente desconocida y potencialmente útil) de grandes bases de datos”
- Esta definición tiene numerosas cosas a definir, que quiere decir no-trivial, útiles para quién?

Agenda

- Qué es Data Mining?
- **Cómo se integra en el proceso de Descubrimiento del conocimiento?**
- Funcionalidades del Data Mining
- Técnicas
 - Supervisadas
 - Redes neuronales
 - Árboles
 - Regresión
 - No supervisadas
 - Clustering
 - Reglas de Asociación

Data Mining: A KDD Process



Agenda

- Qué es Data Mining?
- Cómo se integra en el proceso de Descubrimiento del conocimiento?
- **Funcionalidades del Data Mining**
- Técnicas
 - Supervisadas
 - Árboles
 - Redes neuronales
 - Regresión
 - No supervisadas
 - Clustering
 - Reglas de Asociación

Funcionalidades del DM (1)

Descripción de conceptos: Caracterización y discriminación

- Generalizar, Resumir y contrastar las características de la información (por ejemplo las regiones secas vs. Las regiones húmedas)

Asociación (correlación y causalidad)

- Multi-dimensionales vs. unica dimensión
- $\text{age}(X, "20..29") \wedge \text{income}(X, "20..29K") \rightarrow \text{buys}(X, "PC")$ [support = 2%, confidence = 60%]
- $\text{contains}(T, "computer") \rightarrow \text{contains}(x, "software")$ [1%, 75%]

Funcionalidades del DM (2)

Classificación y Predicción

- Encontrar modelos o funciones que describan y distingan clases para futuras predicciones
- Ej, Clasificar países de acuerdo a su clima, clientes de acuerdo a su comportamiento.
- Presentación: árboles de decisión, reglas de clasificación, redes neuronales
- Predicción: Predecir valores numéricos desconocidos o faltantes.

Cluster analysis

- No se sabe a que clase pertenecen los datos : se agrupan datos para formar clases,
- El Clustering se basa en el principio de maximizar la similitud dentro de la clase y minimizar la misma entre clases

Funcionalidades del DM(3)

Análisis de Outliers

- Outlier: un dato (o un objeto) que no respeta el comportamiento general.
- Puede ser ruido o excepciones, pero son muy útiles en la detección de fraudos o eventos raros.

Análisis de tendencias y evolución

- Tendencia y Desvíos: análisis de regresión
- Análisis de patrones secuenciales
- Análisis de similitudes

Otros análisis estadísticos o de patrones



Agenda

- Qué es Data Mining?
- Cómo se integra en el proceso de Descubrimiento del conocimiento?
- Funcionalidades del Data Mining
- Técnicas
 - Supervisadas
 - **Redes neuronales**
 - Árboles
 - Regresión
 - No supervisadas
 - Clustering
 - Reglas de Asociación



Redes Neuronales (1)

Son sistemas :

- Capaces de aprender
- Adaptarse a a condiciones variantes
- Adaptarse al ruido
- Predecir el estado futuro
- Enfrentar problemas que eran resueltos sólo por el cerebro humano



Redes Neuronales (2)

No son algorítmicas

- No se programan haciéndoles seguir una secuencia predefinida de instrucciones.
- Las RNA generan ellas mismas sus propias "reglas", para asociar la respuesta a su entrada;
- Aprenden por ejemplos y de sus propios errores.
- Utilizan un procesamiento paralelo mediante un gran número de elementos altamente interconectados.



Redes Neuronales (3)

Para mejorar su performance las RNA pueden ser combinadas con otras herramientas

- Lógica Difusa (Fuzzy Logic)
- Algoritmos Genéticos
- Sistemas expertos
- Estadísticas
- Transformadas de Fourier
- Wavelets.



Redes Neuronales – Aplicaciones

La clase de problemas que mejor se resuelven con las redes neuronales son los mismos que el ser humano resuelve mejor pero a gran escala.

- Asociación,
- Evaluación
- Reconocimiento de Patrones.

Las redes neuronales son ideales para problemas que son muy difíciles de calcular

- No requieren de respuestas perfectas,
- Sólo respuestas rápidas y buenas.

Ejemplos

- Escenario bursátil: ¿Compro? ¿Vendo? ¿Mantengo?
- Reconocimiento: ¿se parece? ¿es lo mismo con una modificación?

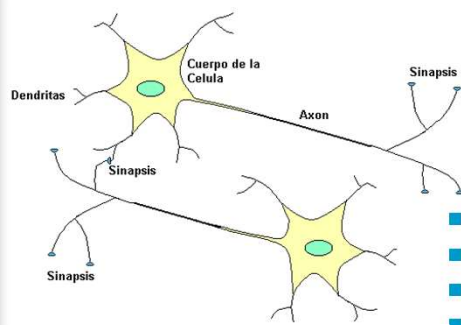


Redes Neuronales - Fallas

Las RNA no son buenas para:

- Cálculos precisos,
- Procesamiento serie,
- Reconocer nada que no tenga inherentemente algún tipo de patrón.

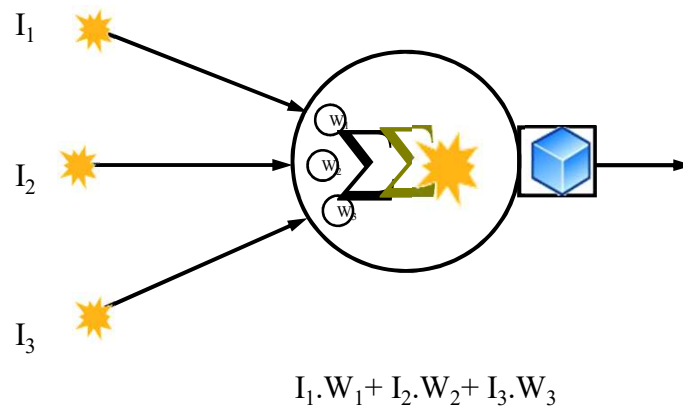
Redes Neuronales - Biología



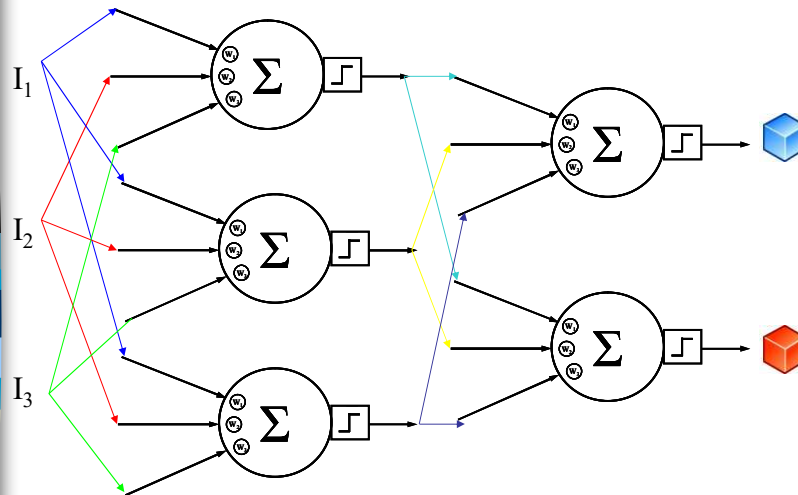
Las Redes Neuronales Artificiales (RNA) se basan en modelos simplificados de las neuronas las reales. Las partes mas comúnmente modelados son:

- Axón
- Dendrita
- Sinapsis
- Cuerpo de la célula

Redes Neuronales - Neurona Modelo



Redes Neuronales - Red Modelo



Redes Neuronales - Aprendizaje

Regla Delta Generalizada o Back Propagation

- Para que una RNA aprenda o se Entrene se deben hacer pasar a todos los valores de entrenamiento por el siguiente proceso, según la topología de la red este ciclo puede repetirse varias veces y con los datos en diferente orden.
- Calcular la diferencia de la salida con la esperada
- Corregir los valores de los W que intervienen en esa salida de modo que se achique esa diferencia
- Se utiliza una constante muy pequeña (Delta)
- No se busca que la diferencia tienda a cero sino que se minimice de a poco
- Si la constante es muy grande o se minimiza la diferencia muy de golpe se corre el riesgo de que cada vez que se aprende algo nuevo se modifique demasiado lo que aprendió anteriormente



Redes Neuronales – Tipos

Tipos de Redes mas utilizados

- Perceptrón Multicapa
- Red de Hopfield (Mapas Asociativos)
- Red de Kohonen (Mapas Autoorganizativos)



Redes de Kohonen o SOM

(Self-Organizing Map).

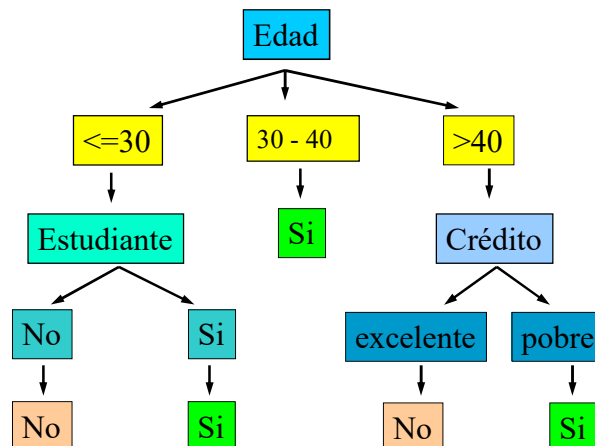
- Hay evidencia que en el cerebro existen neuronas que se organizan de forma que la información se representa internamente en forma de capas bidimensionales.
- En el sistema visual se han detectado mapas del espacio visual en zonas de córtex (capa externa del cerebro).
- En el sistema auditivo se detecta organización según la frecuencia a la que cada neurona alcanza la mayor respuesta (organización tonotópica).

Basado en estas evidencias el Dr. Teuvo Kohonen desarrollo las redes que el prefiere llamar SOM. En las cuales la actualización Delta se realiza solo en la neurona cuyos pesos tengan la mínima distancia con el valor a entrenar y en menor medida se actualizan los pesos de las neuronas vecinas.

Agenda

- Qué es Data Mining?
- Cómo se integra en el proceso de Descubrimiento del conocimiento?
- Funcionalidades del Data Mining
- Técnicas
 - Supervisadas
 - Redes neuronales
 - **Árboles**
 - Regresión
 - No supervisadas
 - Clustering
 - Reglas de Asociación

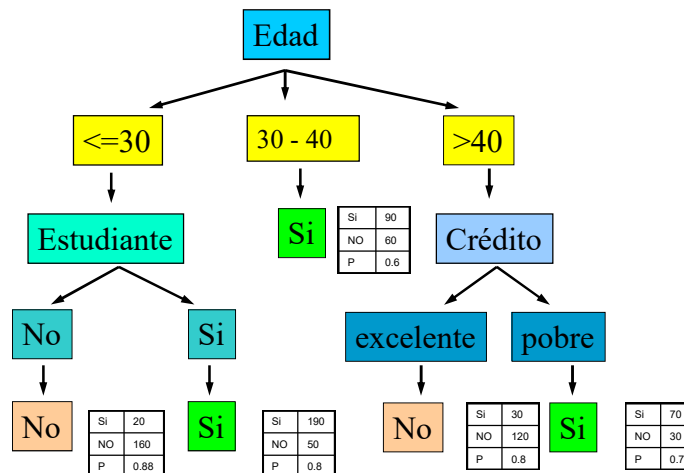
Árbol de Decisión para ver quien compra una computadora



Clasificación por medio de Árboles de Decisión

- Árboles de Decisión
 - Los nodos internos son preguntas sobre los atributos
 - Las hojas representan las etiquetas o clases resultantes
- La generación del árbol tiene fundamentalmente dos pasos
 - Construcción
 - Al comienzo todos los ejemplos están en la raíz del árbol
 - Se dividen los ejemplos en forma recursiva basado en atributos elegidos
 - Pruning
 - Identificar y remover ramas que representan outliers o ruido
- Uso de los árboles de decisión: clasificación de un ejemplo desconocido
 - Se controlan los valores de los atributos del ejemplo para asignarle la clase

Árbol de Decisión con Probabilidad



Extracción de reglas de clasificación a partir de los árboles

- Representa el conocimiento en la forma de reglas de IF-THEN
- Se genera una regla para cada camino desde la raíz hasta las hojas.
- Cada par atributo – valor forma una conjunción
- La hoja tiene la clase a predecir
- Las reglas son fácilmente entendibles por los seres humanos
- Ejemplos

```
IF edad = "<=30" AND estudiante = "no" THEN compra_PC = "no"
IF edad = "<=30" AND estudiante = "yes" THEN compra_PC = "si"
IF edad = "31 - 40" THEN compra_PC = "si"
IF edad = ">40" AND credito = "excelente" THEN compra_PC = "si"
IF edad = ">40" AND credito = "pobre" THEN compra_PC = "no"
```

Evitar el Overfitting en la clasificación

- El árbol obtenido puede hacer overfitting sobre el conjunto de entrenamiento
 - Si hay demasiadas ramas algunas pueden reflejar anomalías
 - Como consecuencia de esto se tiene una performance muy mala sobre ejemplos nuevos
- Dos aproximaciones para evitar el overfitting
 - Prepruning: Interrumpir la construcción del árbol en forma anticipada. No partir un nodo si la mejora que esto produce está por debajo de un cierto umbral.
 - Es difícil encontrar el umbral adecuado
 - Postpruning: quitar ramas de un árbol ya contruido
 - Se puede usar un conjunto diferente del de entrenamiento para hacer esto.

Detección de Valores Extremos, Outliers

Los conjuntos de datos que analizamos generalmente proporcionan un subconjunto de datos en el que existe una variabilidad y/o una serie de errores. Estos datos siguen un comportamiento diferente al resto del conjunto ya sea en una o varias variables. Muchas veces es útil estudiarlos para detectar anomalías, mientras que otras veces es mejor descartarlos de los análisis porque ensucian o influyen en los resultados (por ejemplo en los promedios).



Orígenes de la Variación

Variabilidad de la fuente. Es la que se manifiesta en las observaciones y que se puede considerar como un comportamiento natural de la población en relación a la variable que se estudia.

Errores del medio. Son los que se originan cuando no se dispone de la técnica adecuada para valorar la variable sobre la población, o cuando no existe un método para realizar dicha valoración de forma exacta. En este tipo de errores se incluyen los redondeos forzados que se han de realizar cuando se trabaja con variables de tipo continuo.

Errores del experimentador. Son los atribuibles al experimentador, y que fundamentalmente se pueden clasificar de la siguiente forma:

- **Error de Planificación.** Se origina cuando el experimentador no delimita correctamente la población, y realiza observaciones que pueden pertenecer a una población distinta.
- **Error de Realización.** Se comete al llevar a cabo una valoración errónea de los elementos. Aquí se incluyen, entre otros, transcripciones erróneas de los datos, falsas lecturas realizadas sobre los instrumentos de medida, etc.

Definiciones



A la vista de lo anterior, podemos clasificar las observaciones atípicas o anómalas como:

- Observación **atípica**: Es aquel valor que presenta una gran variabilidad de tipo inherente.
- Observación **errónea**: Es aquel valor que se encuentra afectado de algún tipo de error, sea del medio, del experimentador, o de ambos.

Se llamará "outlier" a aquella observación que siendo atípica y/o errónea, tiene un comportamiento muy diferente respecto al resto de los datos, en relación al análisis que se desea realizar sobre las observaciones. Análogamente, se llamará "inlier" a toda observación no considerada como outlier.

Outliers Peligrosos



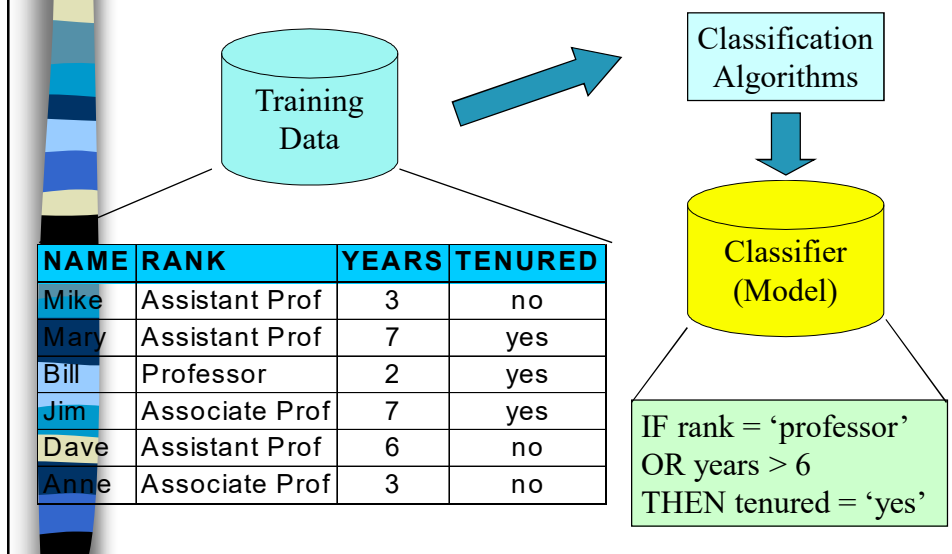
El "agujero de ozono" sobre la antártida es un ejemplo de uno de los outliers más infames de la historia reciente. Es también un buen ejemplo para decir a los que eliminan sistemáticamente los outliers de un dataset simplemente porque son outliers. En 1985 tres investigadores (Farman, Gardinar y Shanklin) fueron desconcertados por un ciertos datos recopilados por el "examen antártico británico" que demostraba que los niveles del ozono para la antártida habían caído el 10% debajo de los niveles normales de enero. El problema era, porqué el satélite Nimbo 7, que tenía instrumentos a bordo para medir con precisión los niveles del ozono, no había registrado concentraciones de ozono semejantemente bajas. Cuando examinaron los datos del satélite no les tomó mucho darse cuenta de que el satélite de hecho registraba estos niveles de concentraciones bajos y lo había estado haciendo por años. ¡Pero como las concentraciones de ozono registradas por el satélite fueron tan bajas eran tratadas como outliers por un programa de computadora y desechadas! El satélite Nimbo 7 de hecho había estado recolectando la evidencia de los niveles bajos de ozono desde 1976. El daño a la atmósfera causada por los clorofluocarburos pasó desapercibido y no fue tratado por nueve años porque los outliers fueron desechados sin ser examinados.

Moraleja: No tirar los outliers sin examinarlos, porque pueden ser los datos más valiosos de un dataset.

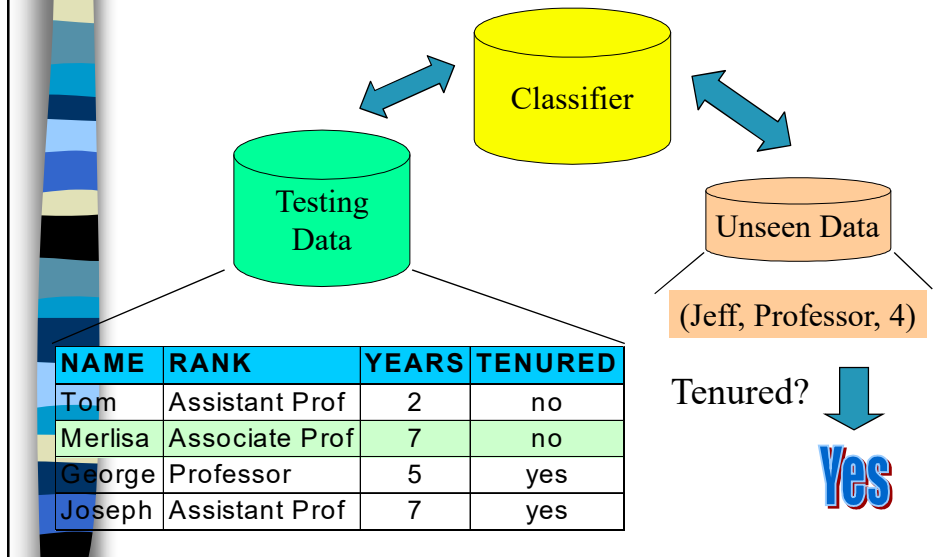
Clasificación—Un proceso de dos pasos

- Construcción del modelo: descripción de las clases existentes
 - Cada ejemplo pertenece a una clase determinada
 - El training set es el conjunto de ejemplos que se usa para entrenar el modelo
 - El modelo se representa por medio de reglas de clasificación, árboles o fórmulas matemáticas
- Uso del modelo: para clasificar ejemplos futuros o desconocidos
 - Estimar la precisión del modelo
 - Para esto se aplica el modelo sobre un conjunto de test y se compara el resultado del algoritmo con el real.
 - Precisión es el porcentaje de casos de prueba que son correctamente clasificados por el modelo
 - El conjunto de entrenamiento debe ser independiente del de test para evitar “overfitting”

Proceso de Clasificación (1): Construcción del Modelo



Proceso de Clasificación (2): Uso del modelo para predecir



Agenda

- Qué es Data Mining?
- Cómo se integra en el proceso de Descubrimiento del conocimiento?
- Funcionalidades del Data Mining
- Técnicas
 - Supervisadas
 - Redes neuronales
 - Árboles
 - **Regresión**
 - No supervisadas
 - Clustering
 - Reglas de Asociación



Regresión Lineal

Para poder crear un modelo de regresión lineal, es necesario que se cumpla con los siguientes supuestos:

- La relación entre las variables es lineal.
- Los errores son independientes.
- Los errores tienen varianza constante.
- Los errores tienen una esperanza matemática igual a cero.
- El error total es la suma de todos los errores.



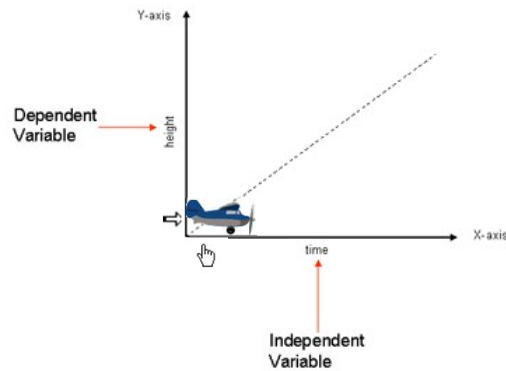
Tipos de Regresión Lineal

- Regresión lineal simple. Sólo se maneja una variable independiente
- Regresión lineal múltiple. Maneja varias variables independientes.

Represión Lineal Ejemplo

Variables dependientes:

Son las variables de respuesta que se observan en el estudio y que podrían estar influidas por los valores de las variables independientes.



Variables independientes: Son las que se toman para establecer agrupaciones en el estudio, clasificando intrínsecamente a los casos del mismo

Regresión Logística

- La regresión logística Se aplica cuando la variable dependiente es dicotómica o politómica y no numérica
- Para poder aplicar una regresión se asocia la variable dependiente a su probabilidad de ocurrencia.
- Por lo tanto el resultado de una regresión logística es la probabilidad de ocurrencia del suceso

Por qué usar clasificación Bayesiana?

- Aprendizaje probabilístico: calcular explícitamente las probabilidades de la hipótesis está entre las aproximaciones más prácticas para cierto tipo de problemas.
- Incremental: Cada ejemplo de entrenamiento puede incrementar / disminuir la probabilidad de que una hipótesis sea correcta.
- Predicción probabilística: se pueden efectuar múltiples predicciones, cada una pesada por su probabilidad

Teorema de Bayes

- Dado un set D de datos de entrenamiento, la probabilidad a posteriori de una hipótesis h $P(h|D)$ sigue el teorema de Bayes

$$P(h|D) = \frac{P(D|h)P(h)}{P(D)}$$

Estimando las probabilidades a posteriori

Teorema de Bayes :

$$P(C|X) = P(X|C) \cdot P(C) / P(X)$$

- $P(X)$ es constante para todas las clases
- $P(C)$ = frecuencia relativa de los ejemplos de la clase C
- C tal que $P(C|X)$ sea maxima =
- C tal que $P(X|C) \cdot P(C)$ sea maxima
- Problema: calcular $P(X|C)$ is imposible!

Naïve Bayes

- Asume que los atributos son independientes:

$$P(C_j|V) \propto P(C_j) \prod_{i=1}^n P(v_i|C_j)$$

- Reduce fuertemente el costo de cálculo, solo cuenta la distribución de la clase.

Naïve Bayes

- Asunción Naïve : los atributos son independientes

$$P(x_1, \dots, x_k | C) = P(x_1 | C) \cdot \dots \cdot P(x_k | C)$$

- Si el atributo i es **categorico** :
 $P(x_i | C)$ se calcula como la frecuencia relativa de los ejemplos que tienen el valor x_i en el atributo i y corresponde a la clase C
- Si el atributo i no es **continuo**:
 $P(x_i | C)$ se estima por medio de una función de densidad
- Se calcula fácilmente.

Clasificación bayesiana

- El problema de clasificación puede formalizarse usando **probabilidades a-posteriori** :

$$P(C|X) = \text{prob. que la tupla } X = \langle x_1, \dots, x_k \rangle \text{ sea de la clase } C.$$

- Ej. $P(\text{class} = N \mid \text{outlook} = \text{sunny}, \text{windy} = \text{true}, \dots)$

- Idea: asignar el ejemplo X a la clase C que haga que $P(C|X)$ sea máxima

Play-tennis example: estimating $P(x_i|C)$

Outlook	Temperature	Humidity	Windy	Class
sunny	hot	high	false	N
sunny	hot	high	true	N
overcast	hot	high	false	P
rain	mild	high	false	P
rain	cool	normal	false	P
rain	cool	normal	true	N
overcast	cool	normal	true	P
sunny	mild	high	false	N
sunny	cool	normal	false	P
rain	mild	normal	false	P
sunny	mild	normal	true	P
overcast	mild	high	true	P
overcast	hot	normal	false	P
rain	mild	high	true	N

$$P(p) = 9/14$$

$$P(n) = 5/14$$

outlook

$$P(\text{sunny}|p) = 2/9$$

$$P(\text{sunny}|n) = 3/5$$

$$P(\text{overcast}|p) = 4/9$$

$$P(\text{overcast}|n) = 0$$

$$P(\text{rain}|p) = 3/9$$

$$P(\text{rain}|n) = 2/5$$

temperature

$$P(\text{hot}|p) = 2/9$$

$$P(\text{hot}|n) = 2/5$$

$$P(\text{mild}|p) = 4/9$$

$$P(\text{mild}|n) = 2/5$$

$$P(\text{cool}|p) = 3/9$$

$$P(\text{cool}|n) = 1/5$$

humidity

$$P(\text{high}|p) = 3/9$$

$$P(\text{high}|n) = 4/5$$

$$P(\text{normal}|p) = 6/9$$

$$P(\text{normal}|n) = 2/5$$

windy

Naive Bayes

Dado un conjunto de entrenamiento se pueden calcular las probabilidades

Outlook	P	N	Humidity	P	N
sunny	2/9	3/5	high	3/9	4/5
overcast	4/9	0	normal	6/9	1/5
rain	3/9	2/5			
Temperature			Windy		
hot	2/9	2/5	true	3/9	3/5
mild	4/9	2/5	false	6/9	2/5
cool	3/9	1/5			

Ejemplo de jugar al tenis: clasificar a X

■ $X = \langle \text{rain, hot, high, false} \rangle$

■ $P(X|p) \cdot P(p) =$

$P(\text{rain}|p) \cdot P(\text{hot}|p) \cdot P(\text{high}|p) \cdot P(\text{false}|p) \cdot P(p)$
 $= 3/9 \cdot 2/9 \cdot 3/9 \cdot 6/9 \cdot 9/14 = 0.010582$

■ $P(X|n) \cdot P(n) =$

$P(\text{rain}|n) \cdot P(\text{hot}|n) \cdot P(\text{high}|n) \cdot P(\text{false}|n) \cdot P(n)$
 $= 2/5 \cdot 2/5 \cdot 4/5 \cdot 2/5 \cdot 5/14 = 0.018286$

■ El ejemplo X es clasificado en la clase n
 (no juega)

La hipótesis de independencia...

■ ... hace que se pueda calcular

■ ... raramente se satisface en la práctica, dado
 que las variables frecuentemente se encuentran
 correlacionadas.

■ Algunas formas de superar esta limitación:

- Redes Bayesianas combinan el razonamiento bayesiano con relaciones entre los atributos
- Árboles de Decisión, toman de un atributo por vez empezando por los más importantes

Agenda

- Qué es Data Mining?
- Cómo se integra en el proceso de Descubrimiento del conocimiento?
- Funcionalidades del Data Mining
- Técnicas
 - Supervisadas
 - Redes neuronales
 - Árboles
 - Regresión
 - No supervisadas
 - **Clustering**
 - Reglas de Asociación

Qué es el análisis de clusters?

- Cluster: una colección de objetos
 - Similares dentro del cluster
 - Diferentes de los objetos en los otros clusters
- Cluster analysis
 - Agrupar un conjunto de datos en un cluster
- Clustering es **clasificación no supervisada** : no hay clases predefinidas
- Aplicaciones típicas
 - Como una herramienta independiente para tener una idea sobre la distribución de los datos
 - Como un proceso previo a usar otros algoritmos

Qué es un buen Clustering?

Un buen método de clustering produce clusters de alta calidad con

- Alta similitud en la clase
- Baja similitud entre clases

La calidad de un clustering depende de la medida de “similitud” usada por el método y de la forma en que está implementado.

Medición de la calidad de un cluster

- Medida de similitud: La similitud está expresada en base a una función de distancia
- Hay una función separada que mide la bondad del clustering
- Las funciones de distancia a utilizar son muy diferentes de acuerdo al tipo de dato.
- Algunas veces es necesario asignarle “peso” a las variables dependiendo del significado que tienen para el problema

Distancias

$$d_{ij} = \sum_{k=1}^p w_k |x_{ik} - x_{jk}|$$

City-Block (Manhattan)

$$d_{ij} = \sqrt{\sum_{k=1}^p w_k (x_{ik} - x_{jk})^2}$$

Euclídea

$$d_{ij} = \lambda \sqrt{\sum_{k=1}^p w_k (x_{ik} - x_{jk})^{\lambda}} \quad \lambda > 0$$

Minkowski

Otras

$$d_{ij} = \frac{\sum_{k=1}^p x_{ik} \cdot x_{jk}}{\sqrt{\sum_{k=1}^p x_{ik}^2} \cdot \sqrt{\sum_{k=1}^p x_{jk}^2}}$$

$$d_{ij} = \frac{\sum_{k=1}^p (x_{ik} - \bar{x}_i) \cdot (x_{jk} - \bar{x}_j)}{\sqrt{\sum_{k=1}^p (x_{ik} - \bar{x}_i)^2} \cdot \sqrt{\sum_{k=1}^p (x_{jk} - \bar{x}_j)^2}}$$

Definición de la distancia: La distancia Euclídea

$$D_{ij} = \sqrt{\sum_{k=1}^n (x_{ki} - x_{kj})^2}$$

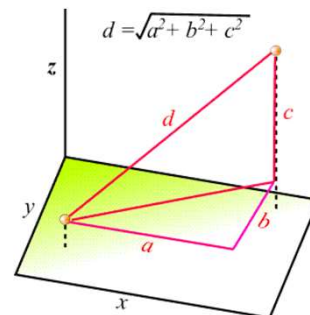
D_{ij} distancia entre los casos i y j
 x_{ki} valor de la variable X_k para el caso j

Problemas:

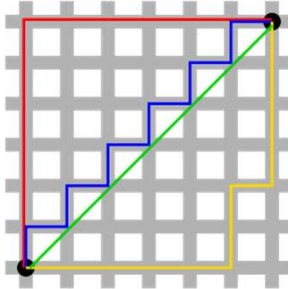
- Diferentes medidas = diferentes ponderaciones
- Correlación entre variables (redundancia)
- Variables faltantes (Missing Values)
- Variables de distinto tipo.
- Incompatibilidad en las Unidades de Medida

Soluciones:

- Análisis de Componentes Principales
- Normalización o Estandarización de las Variables



Manhattan versus Euclidean



El rojo, azul, y amarillo representan la distancia Manhattan, todas tienen el mismo largo(12), mientras que la verde representa la distancia Euclídea con largo de $6 \times \sqrt{2} \approx 8.48$.

Variables numéricas

■ Estandarizar los datos

- Calcular la desviación absoluta de la media

$$s_f = \frac{1}{n}(|x_{1f} - m_f| + |x_{2f} - m_f| + \dots + |x_{nf} - m_f|)$$

donde $m_f = \frac{1}{n}(x_{1f} + x_{2f} + \dots + x_{nf})$

- Normalizar (z-score)

$$z_{if} = \frac{x_{if} - m_f}{s_f}$$

Similitud entre objetos

- Las distancias se usan habitualmente para medir la similitud entre dos objetos

- Algunas de las más conocidas: distancia de *Minkowski*

$$d(i,j) = \sqrt[q]{(|x_{i1} - x_{j1}|^q + |x_{i2} - x_{j2}|^q + \dots + |x_{ip} - x_{jp}|^q)}$$

Donde $i = (x_{i1}, x_{i2}, \dots, x_{ip})$ y $j = (x_{j1}, x_{j2}, \dots, x_{jp})$ son dos objetos de p dimensiones y q es un entero positivo

- Si $q = 1$, d es la distancia de Manhattan

$$d(i,j) = |x_{i1} - x_{j1}| + |x_{i2} - x_{j2}| + \dots + |x_{ip} - x_{jp}|$$

59

Similitud entre objetos (cont)

- Si $q = 2$, d es la distancia euclídeana:

$$d(i,j) = \sqrt{(|x_{i1} - x_{j1}|^2 + |x_{i2} - x_{j2}|^2 + \dots + |x_{ip} - x_{jp}|^2)}$$

– Propiedades de cualquier función de distancia

- $d(i,j) \geq 0$
- $d(i,i) = 0$
- $d(i,j) = d(j,i)$
- $d(i,j) \leq d(i,k) + d(k,j)$

60

Variables binarias

■ Una tabla de contingencia

		Object <i>j</i>		
		1	0	<i>sum</i>
Object <i>i</i>	1	<i>a</i>	<i>b</i>	<i>a+b</i>
	0	<i>c</i>	<i>d</i>	<i>c+d</i>
	<i>sum</i>	<i>a+c</i>	<i>b+d</i>	<i>p</i>

Coeficiente simple

$$d(i, j) = \frac{b + c}{a + b + c + d}$$

Coeficiente de Jaccard :

$$d(i, j) = \frac{b + c}{a + b + c}$$

61

Variables Nominales

■ Pueden tomar más de dos estados : estado civil

■ Método1: Macheo Simple

– *m*: # de coincidencias, *p*: # total de variables

$$d(i, j) = \frac{p - m}{p}$$

■ Método 2: transformación de las variables en dummy

62

Variables ordinales

■ Puede ser discreta o continua, el orden es importante, por ejemplo nivel de educación

■ Pueden ser tratadas como las numéricas comunes

– Reemplazando por su lugar en el ranking

$$r_{if} \in \{1, \dots, M_f\}$$

– normalizar

$$z_{if} = \frac{r_{if} - 1}{M_f - 1}$$

63

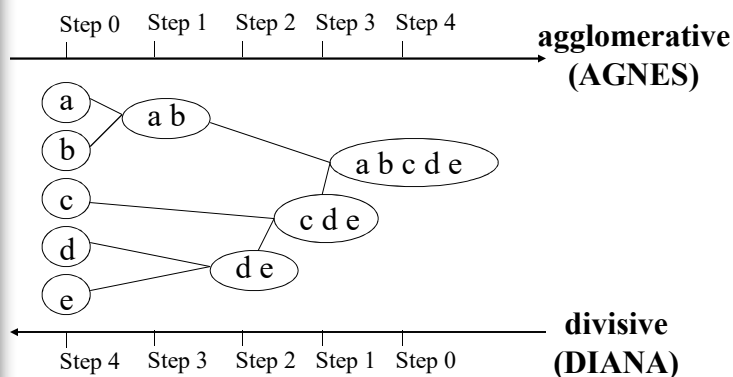
Formas de obtener un cluster

■ Jerárquicas

■ No jerárquicas

Clustering Jerárquico

- Usa la matriz de distancia como criterio. No requiere que el número de cluster sea uno de los parámetros de input



Agrupamiento aglomerativo

Métodos de enlace

- Enlace simple (distancia mínima)
- Enlace Completo (distancia máxima)
- Enlace promedio

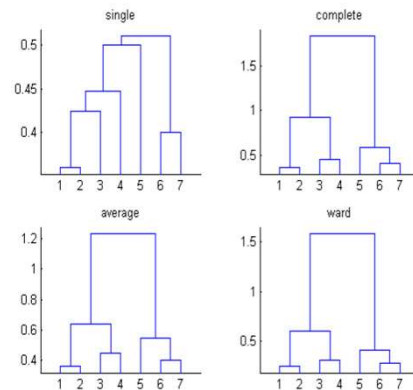
Método de Ward

1. Calcular la suma de las distancias al cuadrado dentro de los clusters
2. Agregar clusters con incremento mínimo en la suma de cuadrados total

Método del centroide

- La distancia entre dos clusters se define como la distancia entre los centroides (medias de los cluster)

Dendrogramas: Otros Métodos



No Jerárquicas: algoritmo básico

- Método de particionamiento: Construir una partición de la base de datos D de n objetos en k clusters
- Dado k encontrar una partición de k clusters que optimice el criterio de partición usado
 - Optimo Global: enumerar todas las particiones posibles
 - Métodos heurísticos:
 - k -means (MacQueen'67): cada cluster esta representado por el centro del cluster
 - k -medoids or PAM (Partition around medoids) (Kaufman & Rousseeuw'87): cada cluster está representado por uno de los objetos del cluster

Métodos jerárquicos vs no jerárquicos

Agrupamiento jerárquico

- No hay decisión acerca del número de clusters
- Existen problemas cuando los datos contienen un alto nivel de error
- Puede ser muy lento

Agrupamiento no jerárquico

- Más rápido, más fiable
- Es necesario especificar el número de clusters (arbitrario)
- Es necesario establecer la semilla inicial (arbitrario)

Agenda

- Qué es Data Mining?
- Cómo se integra en el proceso de Descubrimiento del conocimiento?
- Funcionalidades del Data Mining
- Técnicas
 - Supervisadas
 - Redes neuronales
 - Árboles
 - Regresión
 - No supervisadas
 - Clustering
 - Reglas de Asociación

Propósito de MBA

- Generar reglas del tipo:
 - IF (SI) **condición** ENTONCES (THEN) **resultado**
- Ejemplo:
 - **Si** **producto B** ENTONCES **producto C**
- Association rule mining:
 - “Finding frequent patterns, associations, correlations, or causal structures among sets of items or objects in transaction databases, relational databases, and other information repositories.”

Tipos de reglas según su utilidad

- **Útiles / aplicables** : reglas que contienen buena calidad de información que pueden traducirse en acciones de negocio.
- **Triviales** : reglas ya conocidas en el negocio por su frecuente ocurrencia
- **Inexplicables** : curiosidades arbitrarias sin aplicación práctica



¿Cuán buena es una regla?

■ Medidas que califican a una regla:

- Soporte
- Confianza
- Lift (Improvement)



Ejemplo

T1 = {A, B, C, D}

T2 = {B, C}

T3 = {A, B, C}

T4 = {B, C, D}

T5 = {A, D}

T6 = {A, B}



Soporte

- Es la cantidad (%) de transacciones en donde se encuentra la regla.
 - Ej : “Si B entonces C” está presente en 4 de 6 transacciones.
 - Soporte (B/C) : 66.6%



Confianza

- Cantidad (%) de transacciones que contienen la regla referida a la cantidad de transacciones que contienen la cláusula condicional
 - Ej : Para el caso anterior, si B está presente en 5 transacciones (83.33%)
 - Confianza (B/C) = $66.6/83.3 = 80\%$

Mejora (Improvement)

- Capacidad predictiva de la regla:
 - Mejora = $p(B/C) / p(B) * p(C)$
 - Ej:


$$p(B/C) = 0,67 ; p(B) = 0,833; p(C) = 0,67$$

$$\text{Improv (B/C)} = 0,67(0,833*0,67) = 1.2$$

Mayor a 1 : la regla tiene valor predictivo

Tipos de Reglas

- Booleanas o cuantitativas (de acuerdo a los valores que manejan)
 - $\text{buys}(x, \text{"SQLServer"}) \wedge \text{buys}(x, \text{"DMBook"}) \rightarrow \text{buys}(x, \text{"DBMiner"})$ [0.2%, 60%]
 - $\text{age}(x, \text{"30..39"}) \wedge \text{income}(x, \text{"42..48K"}) \rightarrow \text{buys}(x, \text{"PC"})$ [1%, 75%]
- Una dimensión o varias dimensiones
- Con manejo de jerarquías entre los elementos (taxonomías) o con elementos simples

- 
- Esta presentación fue hecha en base al material que acompaña al libro “Data Mining : Concepts and Techniques” de Han - Kamber



Referencias

- <http://www.kdnuggets.com/>
- <http://www.acm.org/sigkdd/>
- http://www.computer.org/portal/site/transactions/tkde/content/index.jsp?pageID=tkde_home
- <http://domino.research.ibm.com/comm/research.nsf/pages/r.kdd.html>
- <http://www.cs.waikato.ac.nz/~ml/weka/>
- http://www.cs.umd.edu/users/nfa/dm_people_papers.html

Preguntas



Muchas Gracias!!!