

# IRAT QUE NADIE QUIERE ESTUDIAR

## Capítulo 7 - Clustering:

Clustering es una técnica de análisis de datos que se utiliza en el aprendizaje automático y la minería de datos para agrupar un conjunto de datos en grupos o clústeres, donde los elementos dentro de un mismo clúster son más similares entre sí que con los elementos en otros clústeres. El objetivo principal del clustering es encontrar patrones naturales en los datos y segmentarlos en grupos significativos. Es una forma de aprendizaje no supervisado, lo que significa que no se requieren etiquetas o categorías previas para los datos.

### Capítulo 7.1 - Types of Clustering Techniques (Tipos de Técnicas de Clustering):

.Estos son algunos de los tipos comunes de técnicas de clustering:

1. **Clustering Jerárquico (Hierarchical Clustering):** Esta técnica organiza los datos en una jerarquía de clústeres. Puede ser aglomerativo (combinando clústeres más pequeños en clústeres más grandes) o divisivo (dividiendo clústeres grandes en clústeres más pequeños). Proporciona una vista en forma de árbol de la estructura de clústeres, lo que permite explorar clústeres a diferentes niveles de granularidad.
2. **K-Means Clustering:** K-means es un método popular que agrupa los datos en k clústeres, donde "k" es un valor que se debe especificar previamente. Funciona minimizando la varianza intraclúster, es decir, trata de que los elementos dentro de un clúster sean lo más similares posible en términos de distancia euclidiana.
3. **DBSCAN (Density-Based Spatial Clustering of Applications with Noise):** DBSCAN es una técnica de clustering que se basa en la densidad de los puntos de datos. En lugar de especificar el número de clústeres de antemano, DBSCAN identifica clústeres basados en la densidad de puntos cercanos. Puede identificar clústeres de formas irregulares y detectar puntos de datos anómalos como ruido.
4. **Clustering Basado en Modelos (Model-Based Clustering):** Este enfoque utiliza modelos probabilísticos para representar clústeres. El algoritmo EM (Expectation-

Maximization) es comúnmente utilizado en esta categoría.

5. **Clustering Basado en Densidad (Density-Based Clustering):** Además de DBSCAN, hay otras técnicas basadas en densidad, como OPTICS (Ordering Points to Identify the Clustering Structure), que también se centran en encontrar clústeres densos en el espacio de datos.

Cada uno de estos enfoques tiene sus propias ventajas y desventajas, y la elección de la técnica de clustering adecuada depende de la naturaleza de los datos y los objetivos del análisis.

### **K-Means Clustering:**

El algoritmo de K-Means es una de las técnicas de clustering más utilizadas en el aprendizaje automático y la minería de datos. Su objetivo principal es dividir un conjunto de datos en  $k$  clústeres, donde " $k$ " es un valor predefinido que el usuario debe proporcionar. A continuación, se explica cómo funciona el algoritmo de K-Means:

1. **Inicialización de centroides:** El proceso comienza seleccionando aleatoriamente  $k$  puntos como centroides iniciales. Estos centroides representan el centro de cada uno de los  $k$  clústeres. La elección de centroides iniciales puede afectar el resultado final, por lo que se pueden utilizar diferentes estrategias de inicialización.
2. **Asignación de puntos a clústeres:** Para cada punto de datos en el conjunto de datos, se calcula la distancia entre ese punto y todos los centroides. El punto se asigna al clúster cuyo centroide está más cerca en términos de distancia euclidiana.
3. **Actualización de centroides:** Una vez que todos los puntos han sido asignados a clústeres, se recalculan los centroides de cada clúster tomando la media de todas las coordenadas de los puntos en ese clúster. Los centroides se mueven hacia el centro de gravedad de los puntos asignados a su clúster.
4. **Repetición:** Los pasos 2 y 3 se repiten iterativamente hasta que no haya cambios significativos en la asignación de puntos a clústeres o se alcance un número máximo de iteraciones. En cada iteración, los centroides se actualizan y los puntos se reasignan a los clústeres más cercanos.
5. **Convergencia:** Una vez que el algoritmo converge, los datos están agrupados en  $k$  clústeres, y cada punto de datos pertenece a un clúster específico.

Algunos aspectos importantes a considerar sobre el algoritmo de K-Means son:

- Es un algoritmo de aprendizaje no supervisado, lo que significa que no necesita etiquetas previas en los datos.
- La elección de la cantidad de clústeres ( $k$ ) es un parámetro crítico y debe seleccionarse de manera adecuada según el conocimiento del dominio o utilizando métodos como el codo (elbow method) o el índice de silueta.
- K-Means es sensible a la inicialización de los centroides, lo que puede dar lugar a diferentes resultados. Por lo tanto, a menudo se realizan múltiples ejecuciones con diferentes inicializaciones y se selecciona el mejor resultado.
- Funciona bien en datos numéricos y esféricos, pero puede tener dificultades con clústeres de formas irregulares o de tamaños desiguales.

### **DBSCAN (Density-Based Spatial Clustering of Applications with Noise):**

DBSCAN es un algoritmo de clustering basado en la densidad que se utiliza para agrupar datos en clústeres en función de la densidad de los puntos de datos en el espacio. Fue propuesto por Martin Ester, Hans-Peter Kriegel, Jörg Sander y Xiaowei Xu en 1996. La principal ventaja de DBSCAN es su capacidad para identificar clústeres de formas irregulares y manejar la presencia de ruido en los datos. A continuación, se explican los conceptos clave y el funcionamiento de DBSCAN:

#### **Conceptos clave:**

1. **Puntos centrales (Core Points):** Un punto se considera un punto central si dentro de un radio  $\epsilon$  (epsilon) alrededor de él, hay al menos un número mínimo de puntos (MinPts) incluyendo a sí mismo. En otras palabras, los puntos centrales están rodeados por suficientes puntos vecinos.
2. **Puntos de frontera (Border Points):** Un punto de frontera es aquel que no es un punto central, pero está dentro del radio  $\epsilon$  de un punto central. Los puntos de frontera son aquellos que están en el límite de un clúster.
3. **Puntos de ruido (Noise Points):** Los puntos que no son ni centrales ni de frontera se consideran puntos de ruido. Estos son puntos aislados que no pertenecen a ningún clúster.

#### **Funcionamiento de DBSCAN:**

El algoritmo DBSCAN funciona de la siguiente manera:

1. Selecciona un punto de datos aleatorio no visitado en el conjunto de datos.
2. Comprueba si el punto es un punto central, es decir, si tiene al menos MinPts puntos dentro de un radio  $\epsilon$  a su alrededor. Si es un punto central, crea un nuevo clúster y agrega el punto actual y sus vecinos a ese clúster.
3. Si el punto no es un punto central pero es un punto de frontera, se asigna al clúster del punto central más cercano.
4. Se continúa explorando los vecinos de los puntos recién agregados al clúster y se repite el proceso hasta que no se puedan agregar más puntos al clúster.
5. Una vez que se ha formado un clúster, se selecciona un nuevo punto no visitado en el conjunto de datos y se repite el proceso anterior. Esto se repite hasta que todos los puntos de datos hayan sido visitados.
6. Los puntos que no han sido asignados a ningún clúster se consideran puntos de ruido.

#### **Características de DBSCAN:**

- Puede identificar clústeres de formas arbitrarias y no requiere especificar el número de clústeres de antemano, lo que lo hace adecuado para datos con estructuras complejas.
- Es robusto ante la presencia de ruido en los datos, ya que los puntos de ruido se identifican como puntos no asignados a ningún clúster.
- La elección de los parámetros  $\epsilon$  y MinPts es crítica y debe ajustarse cuidadosamente según el dominio y la naturaleza de los datos.

#### **El Capítulo 12 del libro "Predictive Analytics and Data Mining - Concepts and Practice with RapidMiner" de Vijay Kotu y Bala Deshpande "Introducción"**

El Capítulo 12, titulado "Introducción", sirve como punto de partida para comprender los temas esenciales que se abordan en el libro sobre análisis predictivo y minería de datos con la herramienta RapidMiner. El contenido de la introducción incluye:

1. **Contexto de la minería de datos:** Se destaca la importancia de la minería de datos en la actualidad y cómo se ha convertido en una disciplina esencial para las

organizaciones en la toma de decisiones informadas y la obtención de información valiosa a partir de grandes conjuntos de datos.

2. **Naturaleza de los datos:** Se discute la diversidad de los datos que las organizaciones recopilan y almacenan, incluidos datos estructurados y no estructurados, y cómo estos datos pueden ser una fuente valiosa de información.
3. **Desafíos y oportunidades:** Se abordan los desafíos que enfrentan las organizaciones al tratar con grandes volúmenes de datos y la importancia de la minería de datos para descubrir patrones, tendencias y conocimientos que pueden impulsar la toma de decisiones estratégicas.
4. **RapidMiner como herramienta de minería de datos:** Se introduce RapidMiner como una plataforma de código abierto ampliamente utilizada para la minería de datos y el análisis predictivo. Se destaca su versatilidad y facilidad de uso en la creación de flujos de trabajo de análisis de datos.
5. **Estructura del libro:** Se presenta una descripción general de la estructura del libro, que incluye los temas clave a lo largo de los capítulos, como la preparación de datos, la exploración de datos, la modelización predictiva y el despliegue de modelos, todo ello utilizando RapidMiner como plataforma central.

### **Resumen de la sección 12.1 - Clasifying Feature Selection Methods:**

- La sección comienza destacando la relevancia de la selección de características en el proceso de análisis de datos. La selección de características implica identificar las variables más significativas o relevantes en un conjunto de datos y eliminar aquellas que no contribuyen de manera significativa a la tarea en cuestión.
- Se subraya que la selección de características es un paso crucial en el preprocesamiento de datos, ya que puede tener un impacto significativo en la precisión y eficiencia de los modelos de aprendizaje automático. Reducir la dimensionalidad de los datos al seleccionar solo las características más informativas puede llevar a modelos más simples y eficaces.
- La sección menciona que existen varios enfoques y métodos para la selección de características, y que estos métodos se pueden clasificar en diferentes categorías. Se discuten algunas de las categorías comunes de métodos de selección de características, que incluyen:

- **Filtros (Filter Methods):** Estos métodos aplican pruebas estadísticas u otras técnicas para evaluar la relación entre cada característica y la variable objetivo. Luego, se seleccionan las características que superan un umbral predefinido.
- **Envoltura (Wrapper Methods):** Los métodos de envoltura evalúan la calidad de las características al entrenar y validar modelos de aprendizaje automático utilizando diferentes conjuntos de características. Esto se hace de manera iterativa para encontrar el conjunto óptimo de características que mejora el rendimiento del modelo.
- **Incrustación (Embedded Methods):** Estos métodos incorporan la selección de características directamente en el proceso de entrenamiento del modelo. Algunos algoritmos de aprendizaje automático, como la regresión Lasso, automáticamente seleccionan características durante el proceso de modelización.
- La sección resalta que la elección del método de selección de características dependerá del problema específico y de las características de los datos. No hay un método universalmente mejor, y la elección adecuada debe basarse en el conocimiento del dominio y la exploración de los datos.
- Se enfatiza la importancia de realizar evaluaciones de validación cruzada o pruebas en datos independientes para garantizar que la selección de características no conduzca a un sobreajuste (overfitting) y que mejore el rendimiento general del modelo.

### **Resumen de la sección 12.2 - Principal Component Analysis (PCA):**

- La sección comienza introduciendo el concepto de PCA como una técnica que se utiliza para reducir la dimensionalidad de los datos. La reducción de la dimensionalidad implica tomar un conjunto de datos con muchas variables (características) y proyectarlo en un nuevo conjunto de dimensiones más bajo, manteniendo la mayor cantidad de información posible.
- Se explica que PCA busca encontrar nuevas dimensiones, llamadas componentes principales, que son combinaciones lineales de las características originales y que

maximizan la varianza en los datos. Esto significa que los componentes principales capturan la mayor cantidad de información o variabilidad en los datos.

- Se destaca que PCA es especialmente útil cuando se trabaja con datos de alta dimensionalidad, ya que permite reducir la complejidad de los datos y simplificar el análisis sin perder información crítica.
- El proceso de PCA se describe en los siguientes pasos:
  1. Estandarización de datos: Se asegura que todas las características tengan media cero y desviación estándar uno para que las unidades de medida no afecten el análisis.
  2. Cálculo de la matriz de covarianza: Se calcula la matriz de covarianza de las características estandarizadas.
  3. Cálculo de los autovalores y autovectores: Se encuentran los autovalores y autovectores de la matriz de covarianza. Los autovectores son los componentes principales.
  4. Selección de componentes principales: Los componentes principales se eligen en función de los autovalores, donde los componentes con autovalores más grandes capturan la mayor cantidad de variabilidad en los datos.
  5. Proyección de datos: Los datos se proyectan en el espacio de los componentes principales, lo que resulta en una nueva representación de los datos con menos dimensiones.
- Se enfatiza que PCA es una técnica no supervisada, lo que significa que no requiere etiquetas o categorías previas para los datos, y se utiliza principalmente para explorar y visualizar la estructura de los datos.
- Se discuten las aplicaciones de PCA, que incluyen la reducción de dimensionalidad, la eliminación de la multicolinealidad (correlación entre características), la visualización de datos en un espacio de dimensiones reducidas y la compresión de datos.
- La sección concluye resumiendo los beneficios de PCA en la simplificación de datos y la identificación de patrones subyacentes en conjuntos de datos complejos.

### **Sección 10.1 - The Challenge of Unsupervised Learning (El Desafío del Aprendizaje No Supervisado):**

- Esta sección aborda la naturaleza y los desafíos del aprendizaje no supervisado. A diferencia del aprendizaje supervisado, en el que se dispone de etiquetas o respuestas para guiar el modelado, en el aprendizaje no supervisado, no hay respuestas predefinidas. El objetivo es descubrir patrones, estructuras y relaciones en los datos sin una guía explícita.
- Se resalta que el aprendizaje no supervisado se utiliza para explorar la estructura oculta en los datos y puede ser útil en la segmentación de datos en grupos o en la reducción de dimensionalidad.
- También se mencionan algunas de las técnicas de aprendizaje no supervisado comunes, como el análisis de componentes principales (PCA) y el clustering (agrupamiento), que se tratan en las siguientes secciones del capítulo.

## **Sección 10.2 - Principal Component Analysis (Análisis de Componentes Principales o PCA):**

- Esta sección introduce el concepto de PCA, una técnica de reducción de dimensionalidad. PCA se utiliza para transformar un conjunto de datos de alta dimensionalidad en un conjunto de datos de menor dimensionalidad, donde las nuevas dimensiones (llamadas componentes principales) son combinaciones lineales de las características originales.
- Se explica que PCA busca retener la máxima varianza en los datos en las primeras componentes principales, lo que permite resumir la información en un espacio de menor dimensión. Estas componentes principales también son ortogonales entre sí.
- Se ilustra cómo PCA se utiliza para reducir la dimensionalidad de los datos y cómo es útil en la visualización y la identificación de patrones subyacentes en conjuntos de datos complejos.

### **Sección 10.3.2 - Hierarchical Clustering (Clustering Jerárquico):**

- Esta sección se centra en el clustering jerárquico, una técnica de aprendizaje no supervisado utilizada para agrupar datos en una jerarquía de clústeres. El enfoque puede ser aglomerativo o divisivo, dependiendo de si se comienza con los datos individuales y se agrupan en clústeres más grandes o viceversa.
- Se explica cómo el clustering jerárquico se basa en la similitud o distancia entre observaciones, y cómo esta técnica puede revelar estructuras jerárquicas en los



datos, lo que permite la exploración de diferentes niveles de granularidad en la segmentación de datos.

- Se mencionan enfoques específicos utilizados en el clustering jerárquico, como el método de enlace (linkage method), que determina cómo se mide la distancia entre clústeres y cómo se fusionan.