

# MACHINE LEARNING



## Conceptos de Machine Learning

### Modelos de aprendizaje automático

#### Algoritmos de agrupamiento

Existen diversos algoritmos de agrupamiento se utilizan para identificar patrones y estructuras dentro de conjuntos de datos sin etiquetar. Aquí encontraras un informe detallado de las técnicas más utilizadas:

#### K-Means:

El algoritmo de K-Means es un algoritmo popular de agrupamiento en el aprendizaje automático que se utiliza para clasificar instancias en grupos, también conocidos como clústeres, en función de su similitud.



#### Funcionamiento del Algoritmo K-Means:

- 1 Inicialización de Centroides:** El algoritmo comienza con la elección de un número predeterminado de clústeres, denotado como "k". Se seleccionan aleatoriamente "k" instancias del conjunto de datos como centroides iniciales, uno para cada clúster.
- 2 Asignación de Instancias a Clústeres:** Cada instancia del conjunto de datos se asigna al clúster cuyo centroide es el más cercano en términos de distancia euclidiana. La distancia euclidiana entre una instancia y un centroide se calcula en el espacio de características.
- 3 Actualización de Centroides:** Se recalculan los centroides de cada clúster tomando el promedio de todas las instancias asignadas a ese clúster. Este paso se repite iterativamente hasta que los centroides de los clústeres convergen y ya no cambian significativamente.
- 4 Convergencia del Algoritmo:** convergencia se alcanza cuando los centroides de los clústeres ya no experimentan cambios significativos entre iteraciones o cuando se alcanza un número máximo de iteraciones.

**5 Resultados Finales** Al final del proceso, cada instancia del conjunto de datos pertenece a un clúster determinado, y los centroides representan el centro de gravedad de las instancias en ese clúster.

**6 Optimización de la Varianza Intracluster** La optimización del algoritmo de K-Means se basa en minimizar la varianza intracluster, que es la suma de las distancias euclidianas al cuadrado entre cada instancia y el centroide de su clúster asignado. El objetivo es que las instancias dentro de un clúster sean lo más similares posible entre sí, lo que se traduce en una menor varianza intracluster.

#### Desafíos y Consideraciones:



**Elección del Número de Clústeres (k):** Seleccionar el número correcto de clústeres es un desafío y puede requerir métodos como el codo o la silueta para encontrar un valor óptimo.

**Sensibilidad a la Inicialización:** K-Means puede converger a óptimos locales dependiendo de la inicialización de los centroides, por lo que a menudo se ejecuta varias veces con diferentes inicializaciones.

**Manejo de Datos No Lineales:** K-Means es eficaz en datos lineales y globulares, pero puede no funcionar bien en conjuntos de datos no lineales o con formas más complejas.

#### Aplicaciones Prácticas:

**Segmentación de Clientes:** Identificación de grupos de clientes con comportamientos de compra similares para estrategias de marketing personalizadas.

**Análisis de Imágenes:** Agrupamiento de píxeles para segmentar imágenes en regiones homogéneas.

**Agrupamiento de Documentos:** Organización de documentos en categorías o temas similares.

#### Ventajas y Limitaciones:

**Ventajas:** Es computacionalmente eficiente y fácil de entender. Funciona bien en conjuntos de datos grandes y lineales.

**Limitaciones:** Sensible a la elección de "k", no siempre es efectivo en datos no lineales y puede converger a óptimos locales.

**El algoritmo de K-Means es una herramienta valiosa en el aprendizaje automático para la agrupación de datos, siendo especialmente útil en situaciones donde se busca clasificar instancias en grupos homogéneos con respecto a la distancia euclidiana en el espacio de características.**



#### **Agrupamiento Jerárquico:**

El Agrupamiento Jerárquico es un algoritmo de agrupamiento en el aprendizaje automático que se caracteriza por construir una jerarquía de grupos de manera recursiva.

Este enfoque puede ser tanto aglomerativo, donde los datos se agrupan partiendo de instancias individuales hasta formar grupos más grandes, como divisivo, donde se inicia con un grupo que se divide en subgrupos más pequeños.

#### **Funcionamiento del Agrupamiento Jerárquico:**

**Definición de la Distancia entre Grupos o Instancias:** Se define una medida de distancia entre grupos o instancias en función de la similitud entre ellas. Puede ser la distancia euclidiana, la distancia de Manhattan o alguna otra medida de similitud.

**Inicialización de Grupos Individuales:** Cada instancia individual del conjunto de datos se considera como un grupo inicial.

**Iteraciones Recursivas:** Durante cada iteración, se fusionan o dividen grupos basándose en la distancia o similitud entre ellos. En el enfoque aglomerativo, se fusionan los grupos más cercanos iterativamente hasta que todos los datos forman un solo grupo. En el enfoque divisivo, se divide el grupo más grande en subgrupos más pequeños.

**Construcción de la Jerarquía:** Cada iteración crea un nivel en la jerarquía que representa la estructura de agrupamiento en ese momento. El resultado final es un dendrograma, una representación visual de la jerarquía que muestra cómo los grupos se fusionan o dividen a lo largo de las iteraciones.

**Elección del Número de Clústeres:** La elección del número de clústeres se realiza examinando el dendrograma. Pueden seleccionarse clústeres en función de una altura específica del dendrograma. La elección de cómo fusionar o dividir grupos se determina mediante métodos de enlace como el enlace simple, el enlace completo, el enlace promedio, entre otros. Estos métodos definen la distancia entre grupos en función de las distancias entre sus instancias más cercanas, más lejanas o promedio.

#### Aplicaciones Prácticas:



**Biología y Genómica:** Agrupamiento de perfiles genéticos para identificar patrones en expresiones genéticas y clasificación de especies.

**Análisis de Imágenes:** Segmentación de imágenes en regiones homogéneas basada en la similitud de intensidad de píxeles.



**Análisis de Redes Sociales:** Identificación de comunidades o grupos de usuarios con intereses similares en redes sociales.

#### Ventajas y Limitaciones:

**Ventajas:** Permite explorar la estructura jerárquica de los datos. No es necesario especificar el número de clústeres de antemano.

**Limitaciones:** Puede ser computacionalmente costoso, especialmente en grandes conjuntos de datos. La elección del método de enlace y la interpretación del dendrograma pueden ser subjetivas.

**Criterios de Corte:** Para determinar el número de clústeres, se utilizan criterios de corte, como cortar el dendrograma en un nivel que minimice la pérdida de información o maximice la homogeneidad intracluster.



**El Agrupamiento Jerárquico es una técnica versátil en el aprendizaje automático que permite explorar la estructura jerárquica de los datos. Su capacidad para construir dendrogramas facilita la interpretación de la jerarquía de agrupamientos y su aplicabilidad a una variedad de problemas hace que sea una opción valiosa en el análisis de conjuntos de datos no etiquetados.**

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

Es una técnica de agrupamiento en el aprendizaje automático que se destaca por su capacidad para identificar clusters de formas arbitrarias y manejar la presencia de instancias consideradas como ruido. DBSCAN se basa en la densidad de las instancias en el espacio de características y define clusters como áreas de alta densidad, separadas por regiones de baja densidad.

#### Funcionamiento de DBSCAN:

**Se definen los parámetros Clave** Radio de Búsqueda (epsilon,  $\epsilon$ ): Especifica la distancia máxima entre dos instancias para que una de ellas sea considerada como vecina de la otra. Número Mínimo de Vecinos (MinPts): Establece el número mínimo de instancias dentro del radio de búsqueda para que un punto sea considerado como núcleo.

**Identificación de Núcleos y Vecinos:** Cada instancia se clasifica como núcleo, borde o ruido. Un punto es un

núcleo si tiene al menos MinPts puntos dentro de su radio de búsqueda. Un punto es un borde si está dentro del radio de búsqueda de un núcleo pero no es un núcleo en sí mismo.

**Formación de Clústeres** Los núcleos se agrupan en clusters, y los puntos borde se asocian con el clúster de su núcleo más cercano. Los puntos que no son núcleos ni borde se consideran ruido y no se asignan a ningún cluster.

**Detección de Densidad Variable** DBSCAN es capaz de detectar clusters de formas arbitrarias y es menos sensible a la geometría global del conjunto de datos. Puede manejar clusters de diferentes formas y tamaños, así como la presencia de áreas de baja densidad.

#### Ventajas de DBSCAN:



**Manejo de Ruido:** La capacidad de identificar y manejar puntos de datos considerados como ruido es una ventaja crucial en conjuntos de datos del mundo real.

**Flexibilidad Geométrica:** Puede identificar clusters de formas arbitrarias sin verse afectado por la geometría global del conjunto de datos.

**Número Variable de Clústeres:** No es necesario especificar el número de clústeres de antemano, lo que lo hace adecuado para conjuntos de datos con una cantidad desconocida de clusters.

#### Desafíos y Consideraciones:

**Sensibilidad a Parámetros:** La elección de los parámetros epsilon y MinPts puede afectar la calidad de los resultados y puede requerir ajustes.

**Densidades Variables:** DBSCAN puede tener dificultades en conjuntos de datos con clusters de densidades variables.

#### **Aplicaciones Prácticas:**

**Detección de Anomalías:** Identificación de instancias atípicas o anómalas en un conjunto de datos.

**Segmentación de Imágenes:** Agrupamiento de regiones de una imagen basado en la densidad de píxeles.

**Monitoreo de Redes:** Identificación de patrones de tráfico anómalos en redes.

Comparando DBSCAN con K-Means podemos ver que Mientras que K-Means asume la existencia de clusters de forma esférica y requiere que se especifique el número de clusters, DBSCAN puede identificar clusters de formas más complejas y no requiere la especificación del número de clusters.

DBSCAN es una herramienta valiosa en el aprendizaje automático para el agrupamiento basado en densidad, especialmente cuando se enfrenta a conjuntos de datos con clusters de formas arbitrarias y densidades variables. Su capacidad para manejar ruido y su flexibilidad geométrica lo convierten en una opción poderosa en diversas aplicaciones.

#### **Mean Shift:**



El algoritmo Mean Shift es un método de agrupamiento en el aprendizaje automático que se destaca por su capacidad para encontrar modas locales en el espacio de características y agrupar instancias en función de la convergencia hacia estas modas. A diferencia de otros algoritmos de agrupamiento, Mean Shift no requiere la especificación previa del número de clústeres, y puede adaptarse de manera efectiva a la forma y densidad de los datos.

#### Funcionamiento de Mean Shift:

**Asignación de Ventanas** Se inicia asignando ventanas a cada instancia en el conjunto de datos. La forma y tamaño de la ventana se eligen según la distribución de datos.

### **Cálculo del Vector de Desplazamiento (Mean Shift):**

Para cada instancia, se calcula un vector de desplazamiento que apunta hacia la dirección en la que la densidad de instancias es mayor. Este vector se calcula como la media ponderada de las instancias dentro de la ventana, donde las instancias más cercanas tienen mayor peso.

**Actualización de Posiciones:** Las instancias se desplazan en la dirección del vector calculado. Este proceso se repite iterativamente hasta que las instancias convergen hacia las modas locales, donde no hay más cambios significativos.

**Agrupamiento de Instancias:** Las instancias que convergen hacia la misma moda local se asignan al mismo clúster.

### **Ventajas de Mean Shift:**

**Adaptabilidad a la Forma de los Datos:** Mean Shift puede adaptarse eficazmente a la forma y densidad variable de los datos, identificando modas locales incluso en conjuntos de datos complejos.

**Eliminación de la Necesidad de Especificar el Número de Clústeres:** A diferencia de otros algoritmos, Mean Shift no requiere la especificación previa del número de clústeres, ya que identifica naturalmente las modas en los datos.

**Manejo de Densidades Variables:** Puede manejar clusters de diferentes tamaños y densidades, ya que ajusta dinámicamente el tamaño de las ventanas según la distribución de los datos.

### **Desafíos y Consideraciones:**



**Sensibilidad al Parámetro de Banda (Bandwidth):** La elección del tamaño de la ventana (ancho de banda) puede afectar significativamente los resultados y puede requerir ajustes empíricos.

**Eficiencia Computacional:** En conjuntos de datos grandes, el rendimiento computacional puede ser un desafío debido a la necesidad de calcular distancias y actualizar posiciones en cada iteración.

### **Aplicaciones Prácticas:**

**Segmentación de Imágenes:** Identificación de regiones de interés basadas en la convergencia hacia modas locales en el espacio de características.

**Detección de Objetos:** Agrupamiento de píxeles para la detección de objetos en imágenes.

**Análisis de Datos Geoespaciales:** Identificación de patrones de densidad en conjuntos de datos geoespaciales.

Comparando Mean Shift con K-Means podemos ver que Mientras que K-Means asume la existencia de clusters de forma esférica y requiere la especificación del número de clusters, Mean Shift se adapta de manera más flexible a la forma y densidad de los datos, sin requerir la predefinición del número de clusters.

Mean Shift es una técnica de agrupamiento robusta y versátil en el aprendizaje automático, especialmente adecuada para conjuntos de datos con formas y densidades variables. Su capacidad para adaptarse naturalmente a la distribución de los datos y su flexibilidad en la identificación de modas locales lo convierten en una opción valiosa en diversas aplicaciones de agrupamiento.

### **Gaussian Mixture Models (GMM):**



Los Modelos de Mezcla de Gaussianas son algoritmos de agrupamiento en el aprendizaje automático que asumen que el conjunto de datos es una combinación de múltiples distribuciones gaussianas. Este enfoque permite modelar estructuras más complejas y encontrar clusters de forma elíptica en lugar de asumir la forma esférica de los clusters.



## Funcionamiento de Gaussian Mixture Models (GMM)

### **Asunción de una Mezcla de Distribuciones Gaussianas:**

Se parte de la premisa de que el conjunto de datos es una combinación (mezcla) de varias distribuciones gaussianas. Cada componente gaussiano representa un cluster potencial en el conjunto de datos.

**Estimación de Parámetros:** Se utilizan algoritmos de maximización de la expectativa (Expectation-Maximization o EM) para estimar los parámetros del modelo GMM. Los parámetros incluyen las medias, covarianzas y pesos de cada componente gaussiano.

**Asignación de Instancias a Componentes:** Para cada instancia en el conjunto de datos, se calcula la probabilidad de pertenecer a cada componente gaussiano. La instancia se asigna al componente con la probabilidad más alta.

**Flexibilidad para Modelar Estructuras Complejas:** Debido a la naturaleza de la mezcla de distribuciones gaussianas, GMM es capaz de modelar estructuras más complejas en los datos. Puede identificar clusters con formas elípticas y tamaños variables.

## Ventajas de Gaussian Mixture Models (GMM):

**Flexibilidad en Formas de Cluster:** A diferencia de K-Means, GMM no asume la forma esférica de los clusters y puede modelar estructuras más complejas y no lineales.

**Manejo de Clusters con Diferentes Varianzas:** Puede adaptarse a clusters con varianzas diferentes, lo que lo hace adecuado para conjuntos de datos donde los clusters tienen formas y tamaños variados.

**Modelado de Distribuciones Mixtas:** GMM es especialmente útil cuando los datos provienen de una mezcla de diferentes distribuciones.

## Desafíos y Consideraciones:

**Sensibilidad al Número de Componentes:** La elección del número correcto de componentes gaussianos puede ser un desafío y puede requerir métodos como el criterio de información bayesiano (BIC) o la validación cruzada.

**Eficiencia en Conjuntos de Datos Grandes:** En conjuntos de datos grandes, la convergencia del algoritmo puede ser computacionalmente costosa.

## Aplicaciones Prácticas:



**Reconocimiento de Patrones en Imágenes:** Identificación de objetos y patrones en imágenes.

**Análisis de Series Temporales:** Agrupamiento de series temporales para descubrir patrones en datos secuenciales.

**Segmentación de Clientes:** Identificación de grupos de clientes con comportamientos de compra similares.

Comparando GMM con K-Means tenemos que Mientras que K-Means asume clusters con formas esféricas y requiere especificar el número de clusters de antemano, GMM es más flexible al no hacer estas suposiciones y al permitir una modelización más rica de la distribución de los datos.

Los Modelos de Mezcla de Gaussianas (GMM) son herramientas potentes en el aprendizaje automático para el agrupamiento, especialmente cuando se trata de conjuntos de datos que contienen clusters con formas y tamaños variados. Su capacidad para modelar distribuciones mixtas y adaptarse a estructuras complejas los hace relevantes en diversas aplicaciones.

## Spectral Clustering:



Es una técnica de agrupamiento en el aprendizaje automático que utiliza la información espectral de una matriz de afinidad para realizar una reducción de dimensionalidad antes de aplicar técnicas de agrupamiento. A diferencia de algunos algoritmos tradicionales que asumen estructuras lineales en los datos, Spectral Clustering es especialmente útil para identificar clusters en conjuntos de datos con estructuras no lineales o de geometría más compleja.

### Funcionamiento de Spectral Clustering:

**Construcción de la Matriz de Afinidad:** Se calcula una matriz de afinidad que mide la similitud entre todas las instancias del conjunto de datos. La matriz de afinidad puede construirse utilizando funciones de similitud, como la similitud euclidiana o la similitud gaussiana.

**Transformación Espectral:** Se realiza una descomposición espectral de la matriz de afinidad para obtener sus eigenvectores y eigenvalores. Los eigenvectores asociados a los eigenvalores más grandes forman la nueva representación de las instancias en un espacio de menor dimensión.

**Clustering en el Nuevo Espacio:** Las instancias se representan en el nuevo espacio definido por los eigenvectores. Se aplica un algoritmo de agrupamiento, como K-Means, en este espacio de menor dimensión para asignar instancias a clusters.

### Ventajas de Spectral Clustering:

**Manejo de Estructuras No Lineales:** Al realizar una transformación espectral, Spectral Clustering es capaz de capturar estructuras no lineales en los datos, lo que lo hace más efectivo en comparación con algunos algoritmos que asumen linealidad.

**Robustez ante Ruido:** Spectral Clustering puede ser más robusto ante ruido y no depende de supuestos específicos sobre la forma de los clusters.

**Flexibilidad en Número de Clusters:** No requiere especificar el número de clusters a priori y puede adaptarse naturalmente a la complejidad del conjunto de datos.

### Desafíos y Consideraciones:

**Elección de Parámetros:** La elección de parámetros, como la función de similitud y la cantidad de eigenvalores/eigenvectores utilizados, puede afectar la calidad de los resultados y puede requerir ajustes.

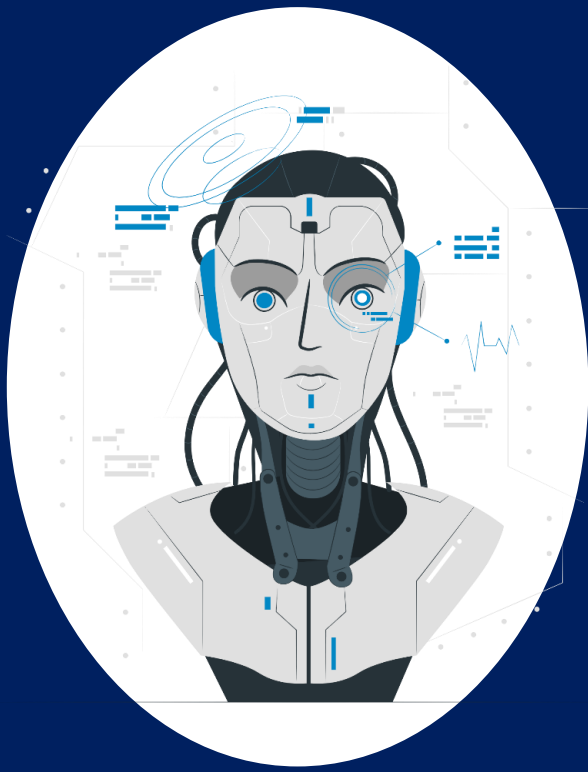
**Eficiencia Computacional:** Dependiendo del tamaño del conjunto de datos, la computación de la matriz de afinidad y la descomposición espectral puede ser computacionalmente costosa.

### Aplicaciones Prácticas:

**Segmentación de Imágenes:** Identificación de regiones y objetos en imágenes.

**Análisis de Redes Sociales:** Agrupamiento de usuarios basado en interacciones y similitudes.

**Biología Computacional:** Agrupamiento de perfiles genéticos para identificar patrones en expresiones genéticas.



Comparando Spectral Clustering con K-Means tenemos que Mientras que K-Means asume que los clusters son esféricos y puede ser menos efectivo en conjuntos de datos no lineales, Spectral Clustering puede capturar de manera más robusta estructuras no lineales y adaptarse a la forma y complejidad de los datos.

Spectral Clustering es una técnica valiosa en el aprendizaje automático para el agrupamiento, especialmente cuando se enfrenta a conjuntos de datos con estructuras no lineales. Su capacidad para realizar una reducción de dimensionalidad basada en la información espectral permite una representación más rica de los datos, lo que puede mejorar la efectividad del agrupamiento en conjuntos de datos complejos.