

Folha de Dicas: Construindo Modelos de Aprendizado Não Supervisionado

Modelos de aprendizado não supervisionado

Nome do Modelo	Descrição Breve	Sintaxe do Código
UMAP	<p>UMAP (Uniform Manifold Approximation and Projection) é usado para redução de dimensionalidade.</p> <p>Prós: Alto desempenho, preserva a estrutura global.</p> <p>Contras: Sensível a parâmetros.</p> <p>Aplicações: Visualização de dados, extração de características.</p> <p>Hiperparâmetros chave:</p> <ul style="list-style-type: none">n_neighbors: Controla o tamanho do bairro local (padrão = 15).min_dist: Controla a distância mínima entre pontos no espaço embutido (padrão = 0.1).n_components: A dimensionalidade da incorporação (padrão = 2).	<pre>from umap.umap_ import UMAP umap = UMAP(n_neighbors=15, min_dist=0.1, n_components=2)</pre>
t-SNE	<p>t-SNE (t-Distributed Stochastic Neighbor Embedding) é uma técnica de redução de dimensionalidade não linear.</p> <p>Prós: Bom para visualizar dados de alta dimensionalidade.</p> <p>Contras: Computacionalmente caro, propenso a overfitting.</p> <p>Aplicações: Visualização de dados, detecção de anomalias.</p> <p>Hiperparâmetros chave:</p> <ul style="list-style-type: none">n_components: O número de dimensões para a saída (padrão = 2).perplexity: Equilibra a atenção entre aspectos locais e globais dos dados (padrão = 30).learning_rate: Controla o tamanho do passo durante a otimização (padrão = 200).	<pre>from sklearn.manifold import TSNE tsne = TSNE(n_components=2, perplexity=30, learning_rate=200)</pre>
PCA	<p>PCA (análise de componentes principais) é usado para redução de dimensionalidade linear.</p> <p>Prós: Fácil de interpretar, reduz ruído.</p> <p>Contras: Linear, pode perder informações em dados não lineares.</p> <p>Aplicações: Extração de características, compressão.</p> <p>Hiperparâmetros chave:</p> <ul style="list-style-type: none">n_components: Número de componentes principais a serem retidos (padrão = 2).whiten: Se deve escalar os componentes (padrão = False).svd_solver: O algoritmo para calcular os componentes (padrão = 'auto').	<pre>from sklearn.decomposition import PCA pca = PCA(n_components=2)</pre>
DBSCAN	<p>DBSCAN (Density-Based Spatial Clustering of Applications with Noise) é um algoritmo de agrupamento baseado em densidade.</p> <p>Prós: Identifica outliers, não requer o número de clusters.</p> <p>Contras: Difícil com clusters de densidade variável.</p> <p>Aplicações: Detecção de anomalias, agrupamento de dados espaciais.</p> <p>Hiperparâmetros chave:</p> <ul style="list-style-type: none">eps: A distância máxima entre dois pontos para serem considerados vizinhos (padrão = 0.5).min_samples: Número mínimo de amostras em um bairro para formar um cluster (padrão = 5).	<pre>from sklearn.cluster import DBSCAN dbscan = DBSCAN(eps=0.5, min_samples=5)</pre>
HDBSCAN	<p>HDBSCAN (Hierarchical DBSCAN) melhora o DBSCAN ao lidar com clusters de densidade variável.</p> <p>Prós: Melhor manejo de densidades variáveis.</p> <p>Contras: Pode ser mais lento que o DBSCAN.</p> <p>Aplicações: Grandes conjuntos de dados, problemas de agrupamento complexos.</p> <p>Hiperparâmetros chave:</p> <ul style="list-style-type: none">min_cluster_size: O tamanho mínimo dos clusters (padrão = 5).min_samples: Número mínimo de amostras para formar um cluster (padrão = 10).	<pre>import hdbscan clusterer = hdbscan.HDBSCAN(min_cluster_size=5)</pre>

Nome do Modelo	Descrição Breve	Sintaxe do Código
Agrupamento K-Means	<p>K-Means é um algoritmo de agrupamento baseado em centróides que agrupa dados em k clusters.</p> <p>Prós: Eficiente, simples de implementar.</p> <p>Contras: Sensível aos centróides iniciais dos clusters.</p> <p>Aplicações: Segmentação de clientes, reconhecimento de padrões.</p> <p>Hiperparâmetros chave:</p> <ul style="list-style-type: none">• n_clusters: Número de clusters (padrão = 8).• init: Método para inicializar os centróides ('k-means++' ou 'random', padrão = 'k-means++').• n_init: Número de vezes que o algoritmo será executado com diferentes sementes de centróides (padrão = 10).	<pre>from sklearn.cluster import KMeans kmeans = KMeans(n_clusters=3)</pre>

Funções associadas usadas

Método	Descrição Breve	Sintaxe do Código
make_blobs	Gera blobs gaussianos isotrópicos para clustering.	<pre>from sklearn.datasets import make_blobs X, y = make_blobs(n_samples=100, centers=2, random_state=42)</pre>
multivariate_normal	Gera amostras de uma distribuição normal multivariada.	<pre>from numpy.random import multivariate_normal samples = multivariate_normal(mean=[0, 0], cov=[[1, 0], [0, 1]], size=100)</pre>
plotly.express.scatter_3d	Cria um gráfico de dispersão 3D usando Plotly Express.	<pre>import plotly.express as px fig = px.scatter_3d(df, x='x', y='y', z='z') fig.show()</pre>
geopandas.GeoDataFrame	Cria um GeoDataFrame a partir de um DataFrame Pandas.	<pre>import geopandas as gpd gdf = gpd.GeoDataFrame(df, geometry='geometry')</pre>
geopandas.to_crs	Transforma o sistema de referência de coordenadas de um GeoDataFrame.	<pre>gdf = gdf.to_crs(epsg=3857)</pre>
contextily.add_basemap	Adiciona um mapa base a um gráfico de GeoDataFrame para contexto.	<pre>import contextily as ctx ax = gdf.plot(figsize=(10, 10)) ctx.add_basemap(ax)</pre>
pca.explained_variance_ratio_	Retorna a proporção da variância explicada por cada componente principal.	<pre>from sklearn.decomposition import PCA pca = PCA(n_components=2) pca.fit(X) variance_ratio = pca.explained_variance_ratio_</pre>

Autor

[Jeff Grossman](#)
[Abhishek Gagneja](#)



Skills Network