

ScikitLearn 操作記錄單 3

組別: Team16 學號: 41071102H 姓名: 徐敏皓

Unsupervised Learning

1. 請根據以下教學資源操作: <http://www.cse.msu.edu/~ptan/dmbook/tutorials/tutorial8/tutorial8.html>

請自行查詢了解下列 scikit-learn 模組的功能作用 <https://scikit-learn.org/stable/>

	Module	Function	試寫程式，實驗該函式所提供功能及主要參數設定效果
Clustering	Sklearn.cluster	KMeans()	<p>功能: 根據資料點之間的距離 (通常是歐幾里得距離)，將資料點劃分到最近的群組。共分成 k 個群。</p> <p>KMeans(n_clusters, init, n_init, max_iter, tol, random_state, algorithm)</p> <p>n_clusters: 指定要分的群數 (預設為 8)。</p> <p>Init: 初始化質心的方法，預設為 k-means++。</p> <p>n_init: 初始聚類執行次數，預設為 10。</p>

			<p>max_iter: 最大迭代次數，預設為 300。</p> <p>tol: 收斂的容忍度，預設為 1e-4。</p> <p>random_state: 隨機種子，用於結果的可重現性。</p> <p>algorithm: 演算法選擇，預設為'lloyd'(支持 'auto','full','elkan'等)。</p>
	Sklearn.cluster	AgglomerativeClustering()	<p>功能: 從每個數據點都是單獨的 cluster，按照距離或相似度，逐步合併最相似的 cluster，直到達到指定的群數或距離閾值。會將數據逐步合併成層級結構的聚類樹（dendrogram）。</p> <p>AgglomerativeClustering(n_clusters, metric, linkage, distance_threshold, compute_full_tree, compute_distances)</p> <p>n_clusters: 指定要生成的 cluster 數量（預設為 2）。若設定為 None，則根據距離臨界值</p>

			<p>(<code>distance_threshold</code>) 進行分割。</p> <p><code>metric</code>: 用於計算距離的指標，取代 <code>affinity</code>。</p> <p>預設為 <code>euclidean</code>。</p> <p><code>Linkage</code>: 用於合併 cluster 的方式，選項為：</p> <ul style="list-style-type: none"> • <code>'ward'</code>：最小化總平方差（僅適用於歐幾里得距離）。 • <code>'complete'</code>：最大距離（最遠點）。 • <code>'average'</code>：平均距離。 • <code>'single'</code>：最小距離（最近點）。 <p><code>distance_threshold</code>: 若設置，定義停止合併的距離臨界值（忽略 <code>n_clusters</code>）。</p> <p><code>compute_full_tree</code>: 是否計算整個階層式聚類樹。當 <code>distance_threshold</code> 設定時，會自動設</p>
--	--	--	--

			<p>為 True。</p> <p>compute_distances: 是否在模型中計算兩個 cluster 之間的距離 (False 為預設)。</p>
	Sklearn.cluster	DBSCAN()	<p>功能: 基於密度的分群算法, 適合處理具有噪聲和複雜分佈的資料集。</p> <p>DBSCAN(eps, min_samples, metric, metric_params, algorithm, leaf_size, p)</p> <p>eps: 兩個樣本被認為是鄰居的最大距離 (即鄰域的半徑)。</p> <p>min_samples: 形成一個核心點所需的最少鄰居數量 (包括核心點本身)。</p> <p>metric: 距離計算方式 (預設為 'euclidean')。</p> <p>可以設為 'manhattan'、'cosine' 等, 或自定義函數。</p>

			<p>metric_params: 距離計算的其他參數。</p> <p>algorithm: 計算最近鄰的算法（可選 'auto'、'ball_tree'、'kd_tree' 或 'brute'）。</p> <p>leaf_size: 用於 BallTree 或 KDTree 的葉節點大小，影響速度。</p> <p>p: 用於 minkowski 距離的冪參數（例如 p=2 對應歐幾里得距離）。</p>
Clustering Evaluation	Sklearn.metrics.cluster	normalized_mutual_info_score()	<p>功能: 評估分群結果與真實標籤的一致性的指標。</p> <p>normalized_mutual_info_score(labels_true, labels_pred, average_method)</p> <p>labels_true: 類別標籤的真實值 (ground truth labels)，即分群的真實分配。</p> <p>labels_pred: 類別標籤的預測值，即分群演算</p>

		<p>法的分配結果。</p> <p>average_method: 預設值為 'arithmetic'。指定正規化的計算方式，選項包括：</p> <ul style="list-style-type: none"> • 'min'：使用最小值正規化。 • 'geometric'：使用幾何平均正規化。 • 'arithmetic'：使用算術平均正規化。 • 'max'：使用最大值正規化。
	Sklearn.metrics.cluster	<p>silhouette_score()</p> <p>功能: 評估聚類結果的質量，其功能是衡量每個數據點與其所屬群集的相似性（凝聚度）與其與最近的其他群集的相似性（分離度）之間的平衡。</p> <p>silhouette_score(X, labels, metric, sample_size, random_state)</p> <p>X: 數組或類似數組的結構（如 DataFrame），</p>

			<p>表示特徵數據或距離矩陣。</p> <p>labels: 數組或列表，群集標籤。</p> <p>metric: 字串，計算相似度的度量方式（默認為'euclidean'）。包括：</p> <p>'euclidean','cosine','manhattan','precomputed'等。</p> <p>sample_size: 整數，指定用於計算的樣本大小（默認為 None，即使用全部數據）。</p> <p>random_state: 整數或 np.random.RandomState，當使用 sample_size 時用於再現性。</p> <p>kwargs: 傳遞給距離計算函數的其他參數。</p>
--	--	--	---

其他參考資源:

- machine learning 參考書: "[Introduction to Machine Learning with Python](#)" 之 github code

github code of the books “Introduction to Machine Learning with Python”

https://github.com/amueller/introduction_to_ml_with_python/blob/master/03-unsupervised-learning.ipynb

Scikit Learn documentation(<http://scikit-learn.org/stable/index.html>)

- 尋搜尋其他可信網路資源