

KNN

课前准备

- 下载Anaconda软件，请点击[这里](#)进行下载。

本节要点

- KNN算法的基本原理。
- KNN算法实现分类与回归任务。
- 超参数调整。
- KD树的构建与邻居选择。【扩展】

KNN算法

引言

在讲解KNN之前，我们来看如下的数据集：

语文	数学	学生
95	93	好
90	92	好
91	96	好
85	82	中
83	87	中
80	84	中
61	69	差
66	63	差
72	65	差
83	77	?

我们将以上数据映射到空间中，进行绘制。因为数据具有两个特征，因此，每条数据对应二维空间中的一个点。

```

1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 plt.rcParams["font.family"] = "SimHei"
5 plt.rcParams["axes.unicode_minus"] = False
6 plt.rcParams["font.size"] = 12
7

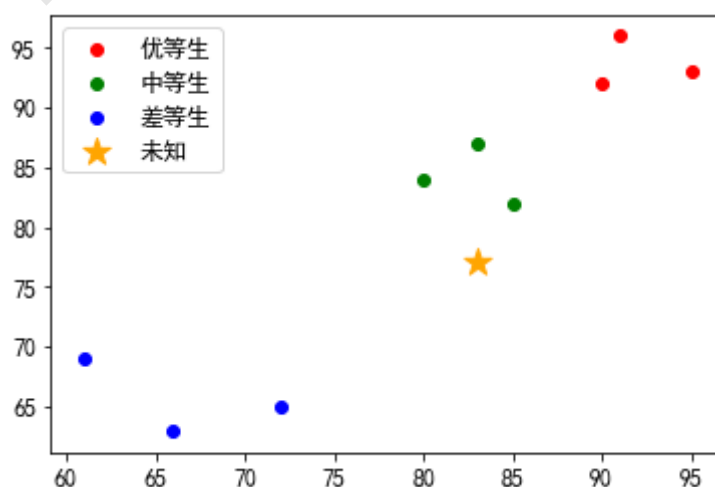
```

```

8 good = np.array([[95, 93], [90, 92], [91, 96]])
9 medium = np.array([[85, 82], [83, 87], [80, 84]])
10 bad = np.array([[61, 69], [66, 63], [72, 65]])
11 unknown = np.array([[83, 77]])
12
13 plt.scatter(good[:, 0], good[:, 1], color="r", label="优等生")
14 plt.scatter(medium[:, 0], medium[:, 1], color="g", label="中等生")
15 plt.scatter(bad[:, 0], bad[:, 1], color="b", label="差等生")
16 plt.scatter(unknown[:, 0], unknown[:, 1], color="orange", marker="*",
17             s=200, label="未知")
18 plt.legend()

```

1 | <matplotlib.legend.Legend at 0x1554b1bec48>



课堂练习

上图中橙色的点（样本），最可能的类别是（ ）。

- A 优等生
- B 中等生
- C 差等生
- D 三种可能概率均等

算法原理

从上例的运行结果中，我们发现，相似度较高的样本，映射到 n 维空间后，其距离会比相似度较低的样本在距离上更加接近，这正是KNN算法的核心思维。

KNN (K-Nearest Neighbor)，即K近邻算法。K近邻就是K个最近的邻居，当需要预测一个未知样本的时候，就由与该样本最接近的K个邻居来决定。KNN既可以用于分类问题，也可以用于回归问题。当进行分类预测时，使用K个邻居中，类别数量最多（或加权最多）者，作为预测结果。当进行回归预测时，使用K个邻居的均值（或加权均值），作为预测结果。

KNN算法的原理在于，样本映射到多维空间时，相似度较高的样本，其距离也会比较接近，反之，相似度较低的样本，其距离也会比较疏远。我们可以将该算法理解为“近朱者赤，近墨者黑”。



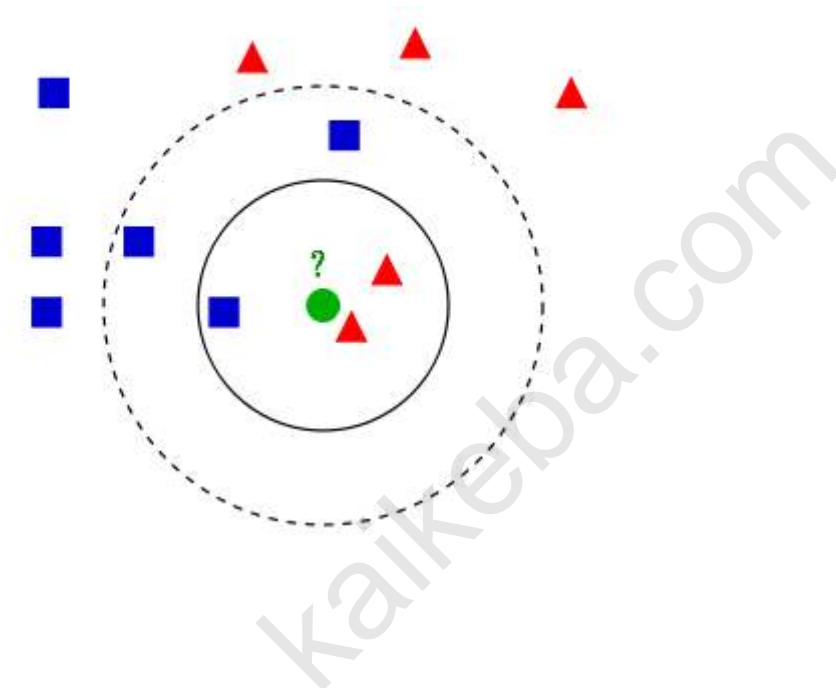
算法超参数

超参数，是指我们在训练模型之前，需要人为指定的参数。该参数不同于模型内部的参数，模型内部的参数是通过训练数据，在训练过程中计算得出的。超参数的不同，可能会对模型的效果产生很大影响。

K值

K 值的选择，会直接影响到预测结果。当 K 值较小时，模型会依赖于附近的邻居样本，具有较好敏感性，但是稳定性会较弱，容易导致过拟合。当 K 值较大时，稳定性增加，但是敏感性会减弱，容易导致欠拟合。

通常情况下，我们可以通过交叉验证的方式，选择最合适的 K 值。



距离度量方式

在scikit-learn中，距离默认使用闵可夫斯基距离（minkowski）， p 的值为2。假设 n 维空间中的两个点为 X 与 Y ：

$$X = (x_1, x_2, \dots, x_n)$$

$$Y = (y_1, y_2, \dots, y_n)$$

则闵可夫斯基距离为：

$$D(X, Y) = (\sum_{i=1}^n |x_i - y_i|^p)^{1/p}$$

可知：当 p 为1时，距离就是曼哈顿距离，当 p 为2时，距离就是欧几里得距离（欧氏距离）。

权重计算方式

权重可以分为两种：

- 统一权重：所有样本的权重相同。
- 距离加权重：样本的权重与待预测样本的距离成反比。



课堂练习



使用KNN预测未知样本 A ，假设 K 的值为3，三个最近的邻居是 B ， C 与 D 。 A 距离 B ， C ， D 的距离分别是2，3，5。如果按照距离加权的方式计算权重，则 B ， C ， D 的权重分别是（ ）。

A $\frac{1}{2}$ ， $\frac{1}{3}$ ， $\frac{1}{5}$

B $\frac{1}{3}$ ， $\frac{1}{3}$ ， $\frac{1}{3}$

C $\frac{1}{2}$ ， $\frac{3}{10}$ ， $\frac{1}{5}$

D $\frac{15}{31}$ ， $\frac{10}{31}$ ， $\frac{6}{31}$



算法步骤

KNN算法的执行过程如下：

1. 确定算法超参数。
 - 确定近邻的数量 K 。
 - 确定距离度量方式。
 - 确定权重计算方式。
 - 其他超参数。
2. 从训练集中选择离待预测样本 A 最近的 K 个样本。

3. 根据这 K 个样本预测 A 。

- 对于分类，使用 K 个样本的类别（或加权类别）预测 A 。
- 对于回归，使用 K 个样本目标值（ y ）的均值（或加权均值）预测 A 。



课堂练习



给定待预测样本 A ，如果在训练集中，存在两个（或更多）不同的邻居 N_1 与 N_2 （ N_1 与 N_2 的标签不同），二者距离 A 的距离相等，假设 K 的值为1，则此时选择哪个邻居较为合理？

- A 选择 N_1
- B 选择 N_2
- C 随机选择
- D 根据 N_1 与 N_2 在训练集中出现的先后顺序选择。
- E C或D



程序示例

使用KNN实现分类

建模预测

我们以鸢尾花数据集为例，通过KNN算法实现分类任务。同样，为了方便可视化，我们只取其中的两个特征。

```
1 from sklearn.neighbors import KNeighborsClassifier
2 from sklearn.datasets import load_iris
3 from sklearn.model_selection import train_test_split
4 from sklearn.metrics import classification_report
5
6 iris = load_iris()
7 x = iris.data[:, :2]
8 y = iris.target
```

```

9 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.25,
    random_state=0)
10 # n_neighbors: 邻居的数量。
11 # weights: 权重计算方式。可选值为uniform与distance。
12 #     uniform: 所有样本统一权重。
13 #     distance: 样本权重与距离成反比。
14 knn = KNeighborsClassifier(n_neighbors=3, weights="uniform")
15 knn.fit(X_train, y_train)
16 y_hat = knn.predict(X_test)
17 print(classification_report(y_test, y_hat))

```

	precision	recall	f1-score	support
0	1.00	1.00	1.00	13
1	0.78	0.44	0.56	16
2	0.44	0.78	0.56	9
accuracy			0.71	38
macro avg	0.74	0.74	0.71	38
weighted avg	0.77	0.71	0.71	38

超参数对模型的影响

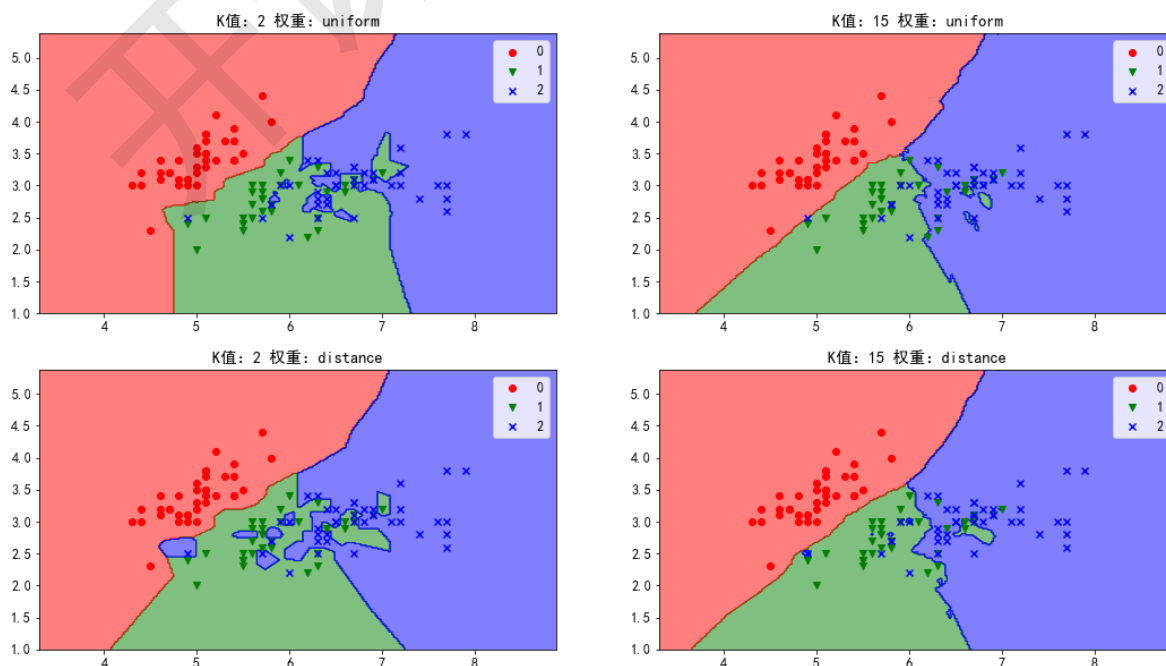
当然，不同的超参数值，会直接影响模型的分类效果。我们可以从决策边界中发现这一点。

```

1 from matplotlib.colors import ListedColormap
2
3 def plot_decision_boundary(model, X, y):
4     color = ["r", "g", "b"]
5     marker = ["o", "v", "x"]
6     class_label = np.unique(y)
7     cmap = ListedColormap(color[: len(class_label)])
8     x1_min, x2_min = np.min(X, axis=0)
9     x1_max, x2_max = np.max(X, axis=0)
10    x1 = np.arange(x1_min - 1, x1_max + 1, 0.02)
11    x2 = np.arange(x2_min - 1, x2_max + 1, 0.02)
12    X1, X2 = np.meshgrid(x1, x2)
13    Z = model.predict(np.c_[X1.ravel(), X2.ravel()])
14    Z = Z.reshape(X1.shape)
15    plt.contourf(X1, X2, Z, cmap=cmap, alpha=0.5)
16    for i, class_ in enumerate(class_label):
17        plt.scatter(x=X[X[y == class_, 0], y=X[X[y == class_, 1],
18                c=cmap.colors[i], label=class_, marker=marker[i])
19    plt.legend()

```

```
1 from itertools import product
2 weights = ['uniform', 'distance']
3 ks = [2, 15]
4
5 plt.figure(figsize=(18, 10))
6 # 计算weights与ks的笛卡尔积组合。这样就可以使用单层循环取代嵌套循环，
7 # 增加代码可读性与可理解性。
8 for i, (w, k) in enumerate(product(weights, ks), start=1):
9     plt.subplot(2, 2, i)
10    plt.title(f"K值: {k} 权重: {w}")
11    knn = KNeighborsClassifier(n_neighbors=k, weights=w)
12    knn.fit(X, y)
13    plot_decision_boundary(knn, X_train, y_train)
```



通过决策边界，我们可以得出如下结论：

- K 的值越小，模型敏感度越强（稳定性越弱），模型也就越复杂，容易过拟合。
- K 的值越大，模型敏感度越弱（稳定性越强），模型也就越简单，容易欠拟合。

超参数调整

在实际应用中，我们很难单凭直觉，就能够找出合适的超参数，通常，我们可以通过网格交叉验证的方式，找出效果最好的超参数。

```

1 from sklearn.model_selection import GridSearchCV
2
3 knn = KNeighborsClassifier()
4 # 定义需要尝试的超参数组合。
5 grid = {"n_neighbors": range(1, 11, 1), "weights": ['uniform', 'distance']}
6 # estimator: 评估器, 即对哪个模型调整超参数。
7 # param_grid: 需要检验的超参数组合。从这些组合中, 寻找效果最好的超参数组合。
8 # scoring: 模型评估标准。
9 # n_jobs: 并发数量。
10 # cv: 交叉验证折数。
11 # verbose: 输出冗余信息, 值越大, 输出的信息越多。
12 gs = GridSearchCV(estimator=knn, param_grid=grid, scoring="accuracy",
13                   n_jobs=-1, cv=3, verbose=10)
14 gs.fit(X_train, y_train)

```

1 Fitting 3 folds for each of 20 candidates, totalling 60 fits

```

1 [Parallel(n_jobs=-1)]: Using backend LokyBackend with 12 concurrent
   workers.
2 [Parallel(n_jobs=-1)]: Done   1 tasks      | elapsed:    0.8s
3 [Parallel(n_jobs=-1)]: Done   8 tasks      | elapsed:    0.8s
4 [Parallel(n_jobs=-1)]: Done  17 tasks      | elapsed:    0.9s
5 [Parallel(n_jobs=-1)]: Done  26 tasks      | elapsed:    0.9s
6 [Parallel(n_jobs=-1)]: Batch computation too fast (0.1930s.) Setting
   batch_size=2.
7 [Parallel(n_jobs=-1)]: Done  37 tasks      | elapsed:    0.9s
8 [Parallel(n_jobs=-1)]: Done  44 out of  60 | elapsed:    0.9s remaining:
   0.3s
9 [Parallel(n_jobs=-1)]: Done  51 out of  60 | elapsed:    0.9s remaining:
   0.1s
10 [Parallel(n_jobs=-1)]: Done  58 out of  60 | elapsed:    0.9s remaining:
   0.0s
11 [Parallel(n_jobs=-1)]: Done  60 out of  60 | elapsed:    0.9s finished

```

```

1 GridSearchCV(cv=3, error_score=nan,
2              estimator=KNeighborsClassifier(algorithm='auto', leaf_size=30,
3                                              metric='minkowski',
4                                              metric_params=None,
5                                              n_jobs=None,
6                                              n_neighbors=5, p=2,
7                                              weights='uniform'),
8              iid='deprecated', n_jobs=-1,
9              param_grid={'n_neighbors': range(1, 11),
10                          'weights': ['uniform', 'distance']},
11              pre_dispatch='2*n_jobs', refit=True, return_train_score=False,
12              scoring='accuracy', verbose=10)

```

当网格交叉验证结束后, 我们就可以通过GridSearchCV对象的相关属性, 来获取最好的参数, 分值与模型。


```
1 # 最好的分值。
2 print(gs.best_score_)
3 # 最好的超参数组合。
4 print(gs.best_params_)
5 # 使用最好的超参数训练好的模型。
6 print(gs.best_estimator_)
```

```
1 0.7773826458036984
2 {'n_neighbors': 9, 'weights': 'uniform'}
3 KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
4                       metric_params=None, n_jobs=None, n_neighbors=9, p=2,
5                       weights='uniform')
```

最后，我们使用最好的模型在测试集上进行测试，实现最后的检验。

```
1 estimator = gs.best_estimator_
2 y_hat = estimator.predict(X_test)
3 print(classification_report(y_test, y_hat))
```

		precision	recall	f1-score	support
0		1.00	1.00	1.00	13
1		0.75	0.38	0.50	16
2		0.41	0.78	0.54	9
accuracy				0.68	38
macro avg		0.72	0.72	0.68	38
weighted avg		0.76	0.68	0.68	38



通过网格交叉验证，可以帮助我们完成调参工作，因此，即使我们不太了解超参数的含义，也可以顺利完成调参任务。这种说法正确吗？

- A 正确
- B 不正确



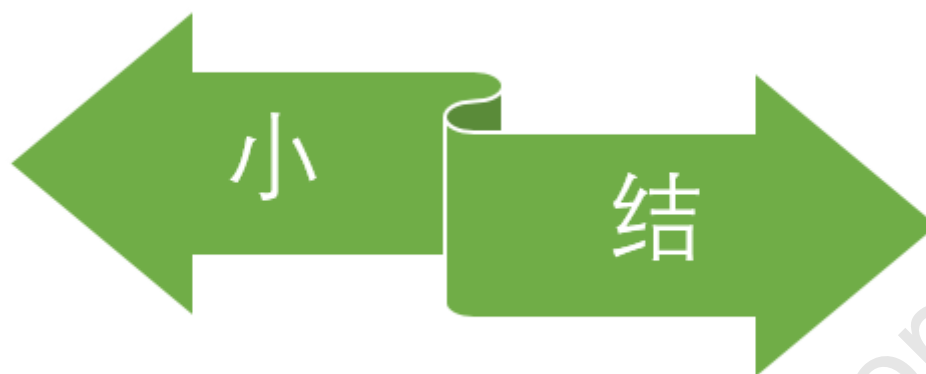
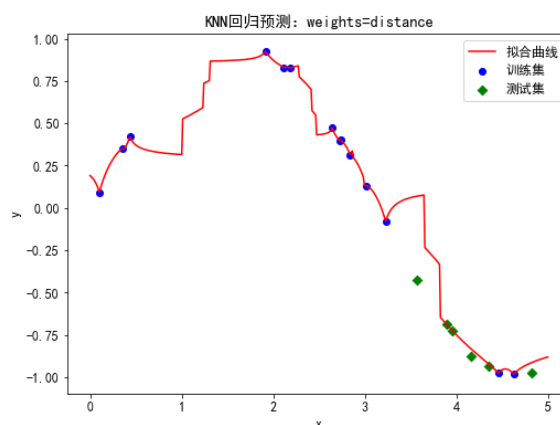
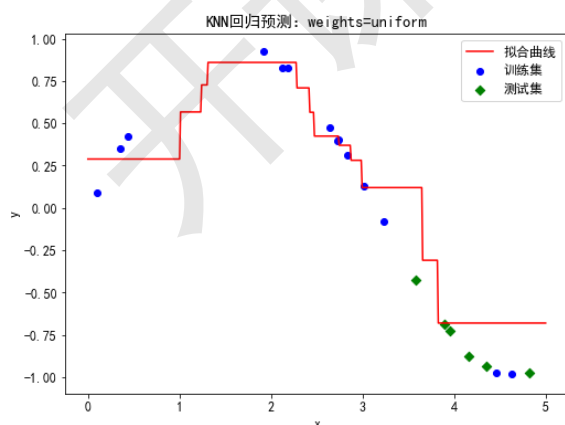
使用KNN回归预测

```
1 from sklearn.neighbors import KNeighborsRegressor
2
3 np.random.seed(0)
4 # 在[0, 5)的范围内，随机生成若干x点。
5 x = 5 * np.random.random(20)
6 x = x[:, np.newaxis]
7 y = np.sin(x) + np.random.normal(0, 0.01, size=len(x))
8 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
9                                                       random_state=0)
```

```

9 plt.figure(figsize=(18, 6))
10 for i, weights in enumerate(['uniform', 'distance']):
11     plt.subplot(1, 2, i + 1)
12     knn = KNeighborsRegressor(n_neighbors=3, weights=weights)
13     knn.fit(X_train, y_train)
14     plt.scatter(X_train, y_train, c="b", marker="o", label="训练集")
15     plt.scatter(X_test, y_test, c="g", marker="D", label="测试集")
16     t = np.linspace(0, 5, 500).reshape(-1, 1)
17     plt.plot(t, knn.predict(t), c="r", label="拟合曲线")
18     plt.legend()
19     plt.title(f"KNN回归预测: weights={weights}")
20     plt.xlabel("x")
21     plt.ylabel("y")

```



KD树

当样本数量较少时，我们可以使用遍历所有样本的方式，找出最近的K的邻居，然而，如果数据集庞大，这种方式会造成大量的时间开销，此时，我们可以使用KD树的方式来选择K个邻居。

KD树算法中，首先是对训练数据进行建模，构建KD树，然后再根据建好的模型来获取邻近样本数据。

拓展点

- KD-Tree的构建。
- 使用KD-Tree寻找最近的邻居。

作业

1. 在数据集中，如果特征之间的量纲差异较大时，使用KNN算法之前，是否需要对其进行标准化操作？
2. 在逻辑回归算法中，样本概率是通过sigmoid函数来计算的。KNN在分类时，是如何预测概率的。
3. 在个人信息数据集上，使用KNN算法建模，并预测一个人的收入是否会超过5万美元。