

机器学习实践记录：从KNN算法到KD树

Author: Limzh

Section1. 序言

在模式识别中，k近邻算法是一个非参数统计学的分类算法。KNN算法秉持着“物以类聚，人以群分”的自然产生的思想，不寻找数据的分布模式（这也是非参数的原因，无模型），而是通过数据集和预测点的距离关系输出预测点的标签。

官方文档的前言如下

最近邻方法背后的原理是从训练样本中找到与新点在距离上最近的预定数量的几个点，然后从这些点中预测标签。这些点的数量可以是用户自定义的常量（K-最近邻学习），也可以根据不同的点的局部密度（基于半径的最近邻学习）确定。距离通常可以通过任何度量来衡量。Neighbors-based（基于邻居的）方法被称为 *非泛化* 机器学习方法，因为它们只是简单地“记住”了其所有的训练数据（可能转换为一个快速索引结构，如 [Ball Tree](#) 或 [KD Tree](#)）。

尽管它简单，但最近邻算法已经成功地适用于很多的分类和回归问题，例如手写数字或卫星图像的场景。作为一个 non-parametric（非参数化）方法，它经常成功地应用于决策边界非常不规则的分类情景下。

通过调用 `sklearn` 中的相关函数，我们可以直接窥得KNN算法的实现结果。

Section2. sklearn 中 KNN 算法

2.1 导入库介绍

```
import numpy as np
import pandas as pandas
import scipy
import matplotlib.pyplot as plt
from matplotlib.colors import ListedColormap
from scipy import stats
from sklearn import neighbors, datasets
from sklearn.datasets.samples_generator import make_classification
```

numpy, pandas, scipy 和 matplotlib 是数据处理和可视化的标准库。sklearn是python的机器学习算法库之一，全称scikit-learn。它集成了四大类机器学习算法，包括分类，回归，降维和聚类。同时它还支持生成算法需要的标数据集，数据的预处理以及数据的引入等功能。`make_classification` 是数据生成器。`neighbors` 包装了大量最近邻有关算法。

2.2 生成数据

```
x, y = make_classification(n_samples=200, n_features=2, n_redundant=0,
                           n_clusters_per_class=1, n_classes=3)
```

上面一行代码生成200个数据样例，每个样本有两个特征。标签值的取值范围为0, 1, 2，意味着所有点被分成三类。

2.3 创建待预测输入

```
# 2. generate to-be-classified data (predicted data)
# -- get the range of to-be-classified point
x_min, x_max = X[:,0].min() - 1, X[:,0].max() + 1
y_min, y_max = X[:,1].min() - 1, X[:,1].max() + 1
# -- mesh
h = .01 # step size for mesh
xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
# -- flatten and pair up
P = np.c_[xx.ravel(), yy.ravel()] # P is an vector of points
```

由于每个样本只具有两个特征，因此整个样本空间为二维平面，由于离样本集合太远的点没有太大的预测价值，因此我们仅关注在样本集合附近的点即可。故而，我们用 `x_min, x_max, y_min, y_max` 将样本集合所在的矩形圈定起来，在该矩形内的所有数据点构成的集合就是我们要预测的样本数据。在这里，我们使用 `np.mesh()` 函数生成网格点。（`xx`和`yy`的具体数值请自己尝试输出）最后，得到`P`作为所有预测点的集合。`P`是一个一维数组，每一个元素是平面一点的坐标。至此，待预测数据已然完成。

2.4 导入分类器并拟合

```
# 3. generate classifier
clf = neighbors.KNeighborsClassifier(n_neighbors=15, weights='distance')
# 4. fit clf with trained-data
clf.fit(X, y)
```

2.5 预测并将图绘制

```
# 5. predict P with fitted classifier and get the result Z full of predicted value
Z = clf.predict(P) # note that it requires P is a vector of 1 dimension

# 6. plot the result.
# -- for scattered points (trained data)
color_bold = ListedColormap(['#156589', '#199934', '#F9AB3B'])

# -- for the predicted plane (to-be-classified data)
color_light = ListedColormap(['#B2EBF2', '#DCEDC8', '#FFE0B2'])

# note that the background(color_light should be printed before to avoid overlapping)
plt.pcolormesh(xx, yy, Z.reshape(xx.shape), cmap = color_light)
plt.scatter(X[:,0], X[:,1], c=y, cmap=color_bold)

# add title
plt.title('sklearn-15NN')
plt.show()
```

2.6 完整代码

```
#!/usr/bin/env python
# -*- coding:UTF-8 -*-
# AUTHOR: Minzhang Li
# FILE: F:\MyGithubs\Machine-Learning\Supervised-Learning\Nearest-Neighbors\code\KNN_sklearn.py
# DATE: 2020/08/01 Sat
# TIME: 11:01:58

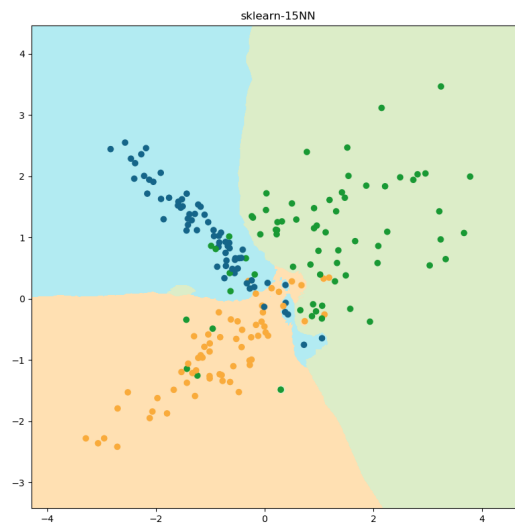
# DESCRIPTION: This file offers an example demonstrating how sklearn lib works.

import numpy as np
import scipy as sp
import matplotlib.pyplot as plt
from matplotlib.colors import ListedColormap
from sklearn.datasets import make_classification
from sklearn import neighbors, datasets
from tqdm import tqdm

def main():
    # 1. generate samples. 200 samples for which are of two dims, 3-classes
    # labels. note that dont mess up 'class' with 'cluster'
    X, y = make_classification(n_samples=200, n_features=2, n_redundant=0,
                              n_clusters_per_class=1, n_classes=3)
    # 2. generate to-be-classified data (predicted data)
    # -- get the range of to-be-classified point
    x_min, x_max = X[:,0].min() - 1, X[:,0].max() + 1
    y_min, y_max = X[:,1].min() - 1, X[:,1].max() + 1
    # -- mesh
    h = .01 # step size for mesh
    xx, yy = np.meshgrid(np.arange(x_min, x_max, h), np.arange(y_min, y_max, h))
    # -- flatten and pair up
    P = np.c_[xx.ravel(), yy.ravel()] # P is a vector of points
    # 3. generate classifier
    clf = neighbors.KNeighborsClassifier(n_neighbors=15, weights='distance')
    # 4. fit clf with trained-data
    clf.fit(X, y)
    # 5. predict P with fitted classifier and get the result Z full of
    # predicted value
    Z = clf.predict(P) # note that it requires P is a vector of 1 dimension
    # 6. plot the result.
    # -- for scattered points (trained data)
    color_bold = ListedColormap(['#156589', '#199934', '#F9AB3B'])
    # -- for the predicted plane (to-be-classified data)
    color_light = ListedColormap(['#B2EBF2', '#DCEDC8', '#FFE0B2'])
    # note that the background(color_light should be printed before to avoid
    # overlapping)
    plt.pcolormesh(xx, yy, Z.reshape(xx.shape), cmap = color_light)
    plt.scatter(X[:,0], X[:,1], c=y, cmap=color_bold)
    # add title
    plt.title('sklearn-15NN')
    plt.show()

if __name__ == '__main__':
    main()
```

2.7 结果



Section3. 性能评估与超参数优化

在以上的实践中，我们做到了使用数据D在分类任务T中获得性能结果，但我们需要对该性能进行评估，并依照评估结果优化超参数K。因为当K过小，对噪音敏感，过拟合；K过大，则欠拟合。选取合适的K值对于分类器性能的好坏至关重要。

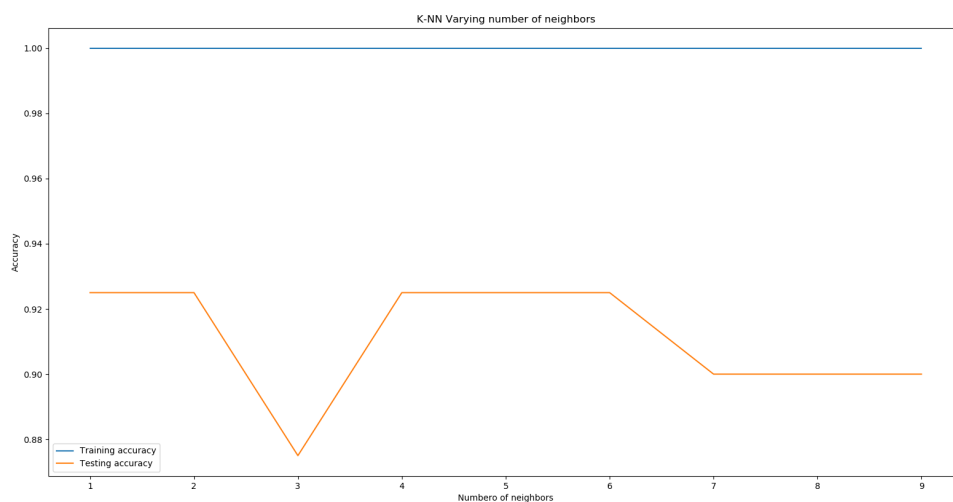
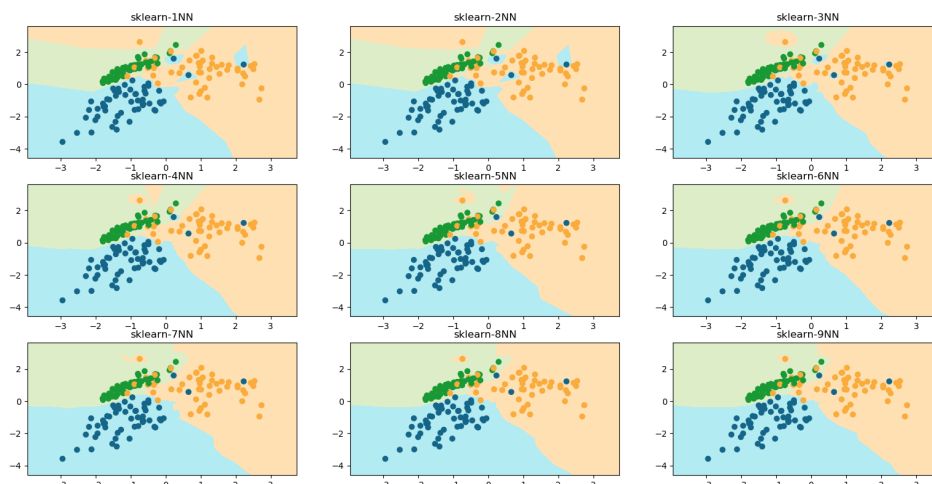
因此我们需要性能评估机制，引入测试集的概念进行超参数优化。

3.1 对不同K的模型的准确度进行评价

```
# To find the best K we should test the result at test set:
K = np.arange(1,10) # [1,10)
test_accuracy = np.zeros(len(K))
train_accuracy = np.zeros(len(K))
res = []
for i, k in enumerate(tqdm(K)):
    # 3. generate classifier
    clf = neighbors.KNeighborsClassifier(n_neighbors=k, weights='distance')
    # 4. fit clf with trained-data
    clf.fit(x_train, y_train)
    # 5. predict P with fitted classifier and get the result Z full of
    # predicted value
    Z = clf.predict(P) # note that it requires P is a vector of 1 dimension
    # push in res list
    res.append(Z)
    # evaluate
    train_accuracy[i] = clf.score(x_train, y_train)
    test_accuracy[i] = clf.score(x_test, y_test)
```

注: `model.score()` 的含义参见附录

3.2 结果



Section4. 动手实践: 蛮力版KNN算法

4.1 仿sklearn接口一览

我们需要实现一个分类器类，并模仿sklearn的风格为其创建接口。这样做是为了增加分类器的拓展性。

```
class KNeighborsClassifier_simple(object):
    def __init__(self, n_neighbors =15, weights = 'distance'):
        pass
    def predict(self, P):
        pass
    def fit(self, X, y):
        pass
    def score(X, y):
        pass
```

初始化的时候，需要传入权重属性和K两个超参数。fit(X, y) 接受训练集的样本数据及其标签集，在fit中进行拟合，学习模型。predict(P) 接受待预测点集合P，输出一个 $1 \times \text{len}(P)$ 的标签集合作为结果。score(X, y) 接受测试集，对相应超参数的情形下模型的精度进行评估。

除此之外，KNN算法的三大元素：距离算法，权重算法（分类决策规则），K的选择中，距离算法和权重算法需要从 `predict(self, P)` 中剥离出来，以方便替换。而 K 的选择则通过简单交叉验证在主函数中得到即可。

故而，最后我们得到的完整接口列表是：

```
class KNeighborsClassifier_simple(object):
    def __init__(self, n_neighbors =15, weights = 'distance'): pass
    def __output(self, neighbors): pass
    def __distance(self, x, p): pass
    def __predict_point(self, p): pass
    def predict(self, P): pass
    def fit(self, X, y): pass
    def score(X, y): pass
    def score(X, y, Z): pass
```

在蛮力版KNN算法中，距离算法选用欧氏距离，分类决策规则使用多数投票规则加距离权重。而最重要的，在 `__predict_point(self, p)` 中，k neighbors的寻找采用蛮力排序算法。

4.2 完整代码

```
class KNeighborsClassifier_simple(object):
    def __init__(self, n_neighbors = 15, weights = 'distance'):
        self.n_neighbors = n_neighbors
        self.weights = weights
        # original data
        self.X = None
        self.y = None
        # combination of samples
        self.samples = None
        self.dim = None

    def __distance(self, x, p):
        '''distance algorithm'''
        res = 0
        for i in range(self.dim):
            res += (x[i] - p[i]) ** 2
        # n root
        return np.power(res, 1 / self.dim)

    def __output(self, neighbors):
        '''when k neighbors array already, output labels'''
        # given neighbors array, output predicted label
        n_labels = len(set(self.y))
        labels = np.zeros(n_labels)
        # count num based on distance
        for i in neighbors:
            if i[-1] != 0: labels[int(i[-2])] += 1/ i[-1]
            else: labels[int(i[-2])] += 1e9
        return np.argmax(labels)

    def __predict_point(self, p):
        '''predict one point'''
        # brute-forcelly calculate distance
        for i in range(len(self.samples)):
            self.samples[i][-1] = self.__distance(self.samples[i], p)
```

```

        # sort and find k neighbors
        ord = list(self.samples.copy())
        ord.sort(key = lambda x: x[-1])
        ord = np.array(ord)
        neighbors = ord[:self.n_neighbors]
        # predict the res
        return self.__output(neighbors)

def fit(self, X, y):
    assert len(X) == len(y) and len(X) != 0
    # X, y
    self.X = X
    self.y = y
    # get dim
    self.dim = len(X[0])
    self.samples = np.c_[X, y]
    distances = np.zeros(len(self.samples), dtype = float)
    # add one more dim in samples for distance
    self.samples = np.c_[self.samples, distances]
def predict(self, P):
    res = np.zeros(len(P))
    for i, p in enumerate(P):
        res[i] = self.__predict_point(p)
    return res
def score(self, X, y):
    '''score without predicted data'''
    y_predicted = self.predict(X)
    res = np.sum(y_predicted == y) / len(y)
    return res
def score(self, X, y, Z):
    '''score with predicted data'''
    y_predicted = Z
    res = np.sum(y_predicted == y) / len(y)
    return res

```

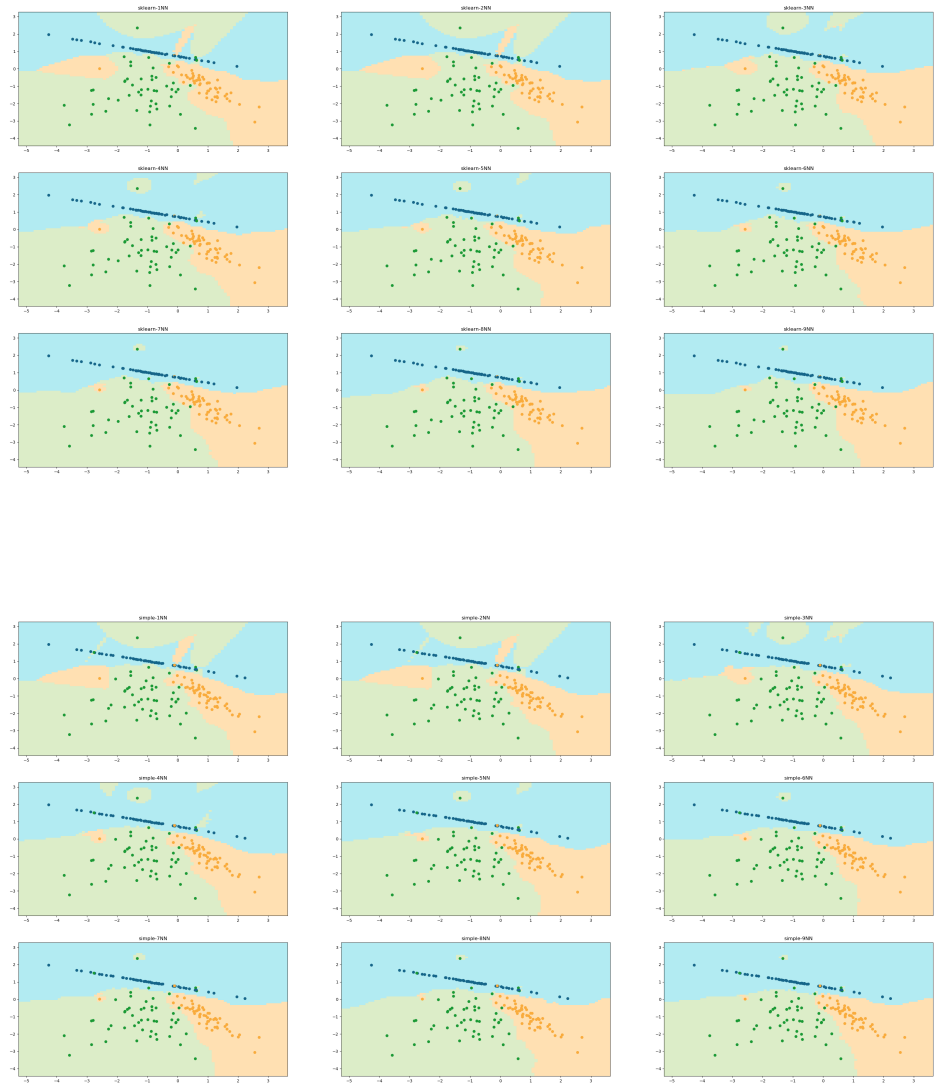
4.3 复杂度分析

假设样本空间的维度是 r , 数据集的大小是 p , 预测输入的大小是 n . 那么对于预测输入的每一个预测点, 都要对数据集的所有点进行一次距离计算, 这意味着 $O(r \cdot n \cdot p)$ 的时间复杂度, 而对于每一个点, 计算完和每一个数据集点的距离之后都要进行一次排序来找到 k 个邻居, 这意味着 $O(n \cdot \log n)$ 的时间复杂度。最后, 决策分类规则要对 k 个邻居的标签进行处理, 需要 $O(n \cdot k)$ 的时间复杂度。所以, 蛮力算法的时间复杂度为 $O(n \cdot (r \cdot p + n \cdot \log n + k))$. 这意味着, 至少需要 $O(n^2 \log n)$ 。

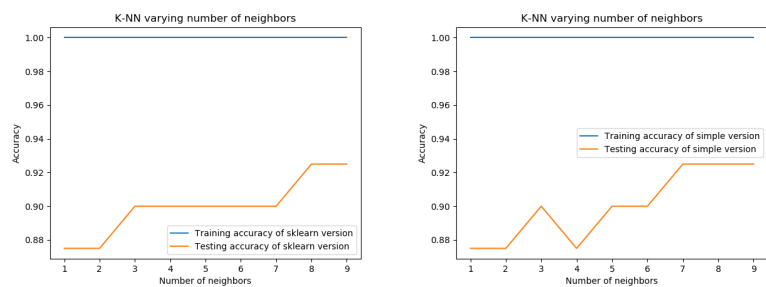
4.4 结果比较

选取数据集 200×2 , 待预测点集合规模10000左右, 有以下结果。

4.4.1 图



4.4.2 准确度



4.4.3 运行时间

版本	数据规模	时间
sklearn	数据集, 200 x 2. 待预测集, 10000. K: 1 - 9	0:00:01.44
simple	数据集, 200 x 2. 待预测集, 10000. K: 1 - 9	0:01:09.12

由此, simple版本的效率远远慢于sklearn.

附录

1. [sklearn库的结构](#)

2. [model.score\(\)的含义](#)