

# 机器学习补充练习

## 示例1：线性回归

读取single.txt文件中的样本，定义线性回归模型，并训练，绘制训练的模型，打印模型的R2分数

```
1  # 线性回归示例
2  import numpy as np
3  # 线性模型
4  import sklearn.linear_model as lm
5  # 模型性能评价模块
6  import sklearn.metrics as sm
7  import matplotlib.pyplot as mp
8
9  x, y = [], []    # 输入、输出样本
10 with open("single.txt", "rt") as f:
11     for line in f.readlines():
12         data = [float(substr) for substr in
13 line.split(",")]
14         x.append(data[:-1])
15         y.append(data[-1])
16
17 x = np.array(x)    # 二维数据形式的输入矩阵，一行一样本，一列一特征
18 y = np.array(y)    # 一维数组形式的输出序列，每个元素对应一个输入样
19 本
20
21 print(x)
22 print(y)
23
24 # 创建线性回归器
25 model = lm.LinearRegression()
26 # 用已知输入、输出数据集训练回归器
27 model.fit(x, y)
28
29 # 根据训练模型预测输出
30 pred_y = model.predict(x)
```

```

28 # 评估指标
29 err = sm.mean_absolute_error(y, pred_y) # 评价绝对值误差
30 print(err)
31 err2 = sm.mean_squared_error(y, pred_y) # 平均平方误差
32 print(err2)
33 err3 = sm.median_absolute_error(y, pred_y) # 中位绝对值误差
34 print(err3)
35 err4 = sm.r2_score(y, pred_y) # R2得分，范围[0, 1]，分值越大越好
36 print(err4)
37
38 # 可视化回归曲线
39 mp.figure('Linear Regression', facecolor='lightgray')
40 mp.title('Linear Regression', fontsize=20)
41 mp.xlabel('x', fontsize=14)
42 mp.ylabel('y', fontsize=14)
43 mp.tick_params(labelsize=10)
44 mp.grid(linestyle=':')
45 # 绘制样本点
46 mp.scatter(x, y, c='dodgerblue', alpha=0.8, s=60,
47            label='Sample')
48 # 绘制拟合直线
49 sorted_indices = x.T[0].argsort()
50 mp.plot(x[sorted_indices], pred_y[sorted_indices],
51         c='orangered', label='Regression')
52 mp.show()

```

## 示例2：利用随机森林实现共享单车投放量预测

- 数据集：一段时期内共享单车使用量，特征：日期、季节、年、月、小时、是否是假期、星期几、是否为工作日、天气、温度、体感温度、湿度、风速；标签：游客使用量、注册用户使用量、总使用量
- 实现代码：

```

1 # -*- coding: utf-8 -*-

```

```

2 # 使用随机森林实现共享单车使用量预测
3
4 import csv
5 import numpy as np
6 import sklearn.utils as su
7 import sklearn.ensemble as se
8 import sklearn.metrics as sm
9 import matplotlib.pyplot as mp
10
11 # 读取共享单车使用率文件中的数据
12 ##### 基于天的数据训练与预测 #####
13 with open("bike_day.csv", "r") as f:
14     reader = csv.reader(f)
15     x, y = [], []
16     for row in reader:
17         x.append(row[2:13]) # 第1列序号掐掉，挑出其中的输入
18         y.append(row[-1]) # 最后一列是输出
19
20 fn_dy = np.array(x[0]) # 保存特征名称
21 x = np.array(x[1:], dtype=float) # 去掉第1行标题部分
22 y = np.array(y[1:], dtype=float) # 去掉第1行标题部分
23
24 # 将矩阵打乱
25 x = su.shuffle(x, random_state=7)
26 y = su.shuffle(y, random_state=7)
27
28 # 计算训练数据的笔数，创建训练集、测试集
29 train_size = int(len(x) * 0.9) # 用90%的数据来训练模型
30
31 train_x = x[:train_size] # 训练输入
32 train_y = y[:train_size] # 训练输出
33
34 test_x = x[train_size:] # 测试输入
35 test_y = y[train_size:] # 测试输出
36
37 # 创建随机森林回归器，并进行训练
38 model = se.RandomForestRegressor(max_depth=10, #最大深度

```

```

39         n_estimators=1000, #树数
    量
40         min_samples_split=2) #最
    小样本数量，小于该数就不再划分子节点
41 model.fit(train_x, train_y) # 训练
42
43 # 基于天统计数据特征的重要性
44 fi_dy = model.feature_importances_
45 # print(fi_dy)
46 pre_test_y = model.predict(test_x)
47 print(sm.r2_score(test_y, pre_test_y)) #打印r2得分
48
49 # 可视化
50 mp.figure('Bike', facecolor='lightgray')
51 mp.subplot(211)
52 mp.title('Day', fontsize=16)
53 mp.ylabel('Importance', fontsize=12)
54 mp.tick_params(labelsize=10)
55 mp.grid(axis='y', linestyle=':')
56 sorted_idx = fi_dy.argsort()[::-1]
57 pos = np.arange(sorted_idx.size)
58 mp.bar(pos, fi_dy[sorted_idx], facecolor='deepskyblue',
    edgecolor='steelblue')
59 mp.xticks(pos, fn_dy[sorted_idx], rotation=30)
60
61 ##### 基于小时的数据训练与预测 #####
62 with open("bike_hour.csv", "r") as f_hr:
63     reader = csv.reader(f_hr)
64     x, y = [], []
65     for row in reader:
66         x.append(row[2:13]) # 第1列序号掐掉，挑出其中的输入
67         y.append(row[-1]) # 输出
68
69 fn_hr = np.array(x[0])
70
71 x = np.array(x[1:], dtype=float)
72 y = np.array(y[1:], dtype=float)
73

```

```

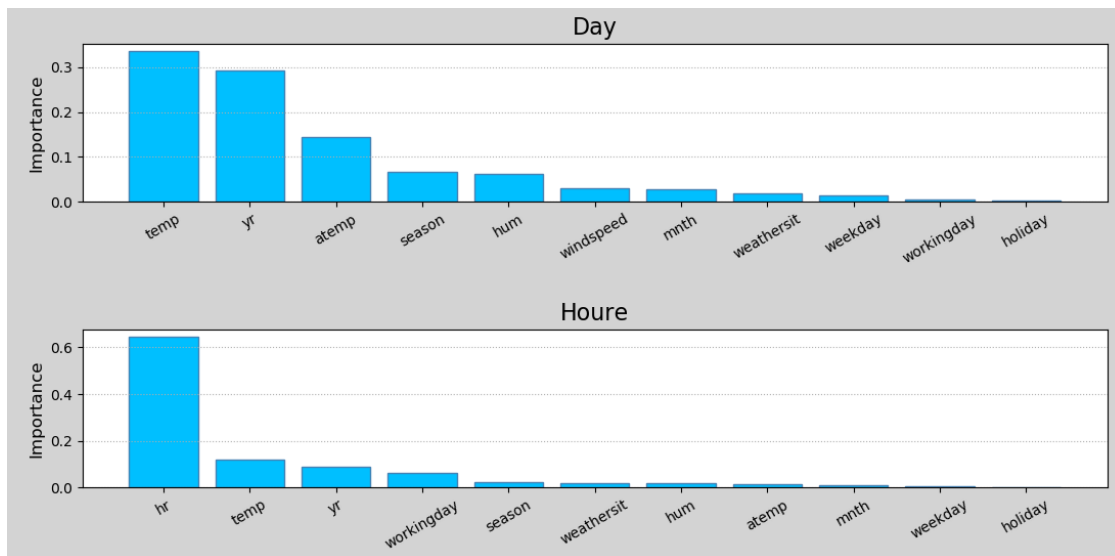
74 x = su.shuffle(x, random_state=7)
75 y = su.shuffle(y, random_state=7)
76
77 # 计算训练数据的笔数，创建训练集、测试集
78 train_size = int(len(x) * 0.9)
79 train_x = x[:train_size] # 训练输入
80 train_y = y[:train_size] # 训练输出
81 test_x = x[train_size:] # 测试输入
82 test_y = y[train_size:] # 测试输出
83
84 # 创建随机森林回归器，并进行训练
85 model = se.RandomForestRegressor(max_depth=10,
86                                   n_estimators=1000, min_samples_split=2)
87
88 fi_hr = model.feature_importances_ # 基于小时数据的特征重要性
89
90 pre_test_y = model.predict(test_x)
91 print(sm.r2_score(test_y, pre_test_y)) #打印r2得分
92
93 #可视化
94 mp.subplot(212)
95 mp.title('Houre', fontsize=16)
96 mp.ylabel('Importance', fontsize=12)
97 mp.tick_params(labelsize=10)
98 mp.grid(axis='y', linestyle=':')
99 sorted_idx = fi_hr.argsort()[::-1]
100 pos = np.arange(sorted_idx.size)
101 mp.bar(pos, fi_hr[sorted_idx], facecolor='deepskyblue',
102        edgecolor='steelblue')
103 mp.xticks(pos, fi_hr[sorted_idx], rotation=30)
104 mp.tight_layout()
105 mp.show()

```

- 打印输出

```
1 0.8915180372559434
2 0.9185448658002986
```

- 特征重要性可视化



## 示例3：利用SVM预测交通流量

利用支持向量机预测体育场周边交通流量。样本特征分别为：星期、时间、对手球队、棒球比赛是否正在进行、通行汽车数量。

```
1 # 利用支持向量机实现交通流量预测
2 # 数据集：17568笔样本
3 # 特征分别为星期、时间、对手球队、棒球比赛是否正在进行，标签为通行汽车数量
4 import numpy as np
5 import sklearn.model_selection as ms
6 import sklearn.svm as svm
7 import sklearn.metrics as sm
8 import matplotlib.pyplot as mp
9 import sklearn.preprocessing as sp
10
11 # 自定义编码器
12 class DigitEncoder():
13     def fit_transform(self, x):
14         return x.astype(int)
15
```

```
16     def transform(self, x):
17         return x.astype(int)
18
19     def inverse_transform(self, x):
20         return x.astype(str)
21
22
23 data = []
24 with open("../data/traffic.txt", "r") as f:
25     for line in f.readlines():
26         line = line.replace("\n", "")
27         data.append(line.split(","))
28 data = np.array(data).T
29
30 encoders, x = [], []
31 for row in range(len(data)):
32     if data[row, 0].isdigit(): # 数值，使用自定义编码器
33         encoder = DigitEncoder()
34     else: # 字符串，使用标签编码器
35         encoder = sp.LabelEncoder()
36
37     if row < len(data) - 1: # 不是最后一行：特征
38         x.append(encoder.fit_transform(data[row]))
39     else: # 最后一行：标签
40         y = encoder.fit_transform(data[row])
41
42     encoders.append(encoder) # 记录编码器
43
44 x = np.array(x).T # 转置还原
45
46 # 划分训练集、测试集
47 train_x, test_x, train_y, test_y = ms.train_test_split(
48     x, y, test_size=0.25, random_state=5)
49
50 # 基于径向基核函数的支持向量机回归器
51 model = svm.SVR(kernel="rbf", c=10, epsilon=0.2)
52 model.fit(train_x, train_y)
53 pred_test_y = model.predict(test_x)
```

```

54
55 print("r2_score:", sm.r2_score(test_y, pred_test_y))
56
57 data = [["Tuesday", "13:35", "San Francisco", "yes"]] #
    待预测数据
58 data = np.array(data).T
59 x = []
60 # 对样本进行编码
61 for row in range(len(data)):
62     encoder = encoders[row]
63     x.append(encoder.transform(data[row]))
64
65 x = np.array(x).T
66 pred_y = model.predict(x)
67 print(int(pred_y))

```

执行结果：

```

1 r2_score: 0.6379517119380995
2 27

```

## 示例 4：凝聚层次对中心不明显的数据聚类

下面来看一个中心点不明显的凝聚层次聚类示例。

```

1 # 凝聚层次聚类示例
2 import numpy as np
3 import sklearn.cluster as sc
4 import matplotlib.pyplot as mp
5 import sklearn.neighbors as nb
6
7 n_sample = 500
8 t = 2.5 * np.pi * (1 + 2 * np.random.rand(n_sample, 1)) #
    产生随机角度
9
10 # 产生数据样本(阿基米德螺线)
11 x = 0.05 * t * np.cos(t)

```

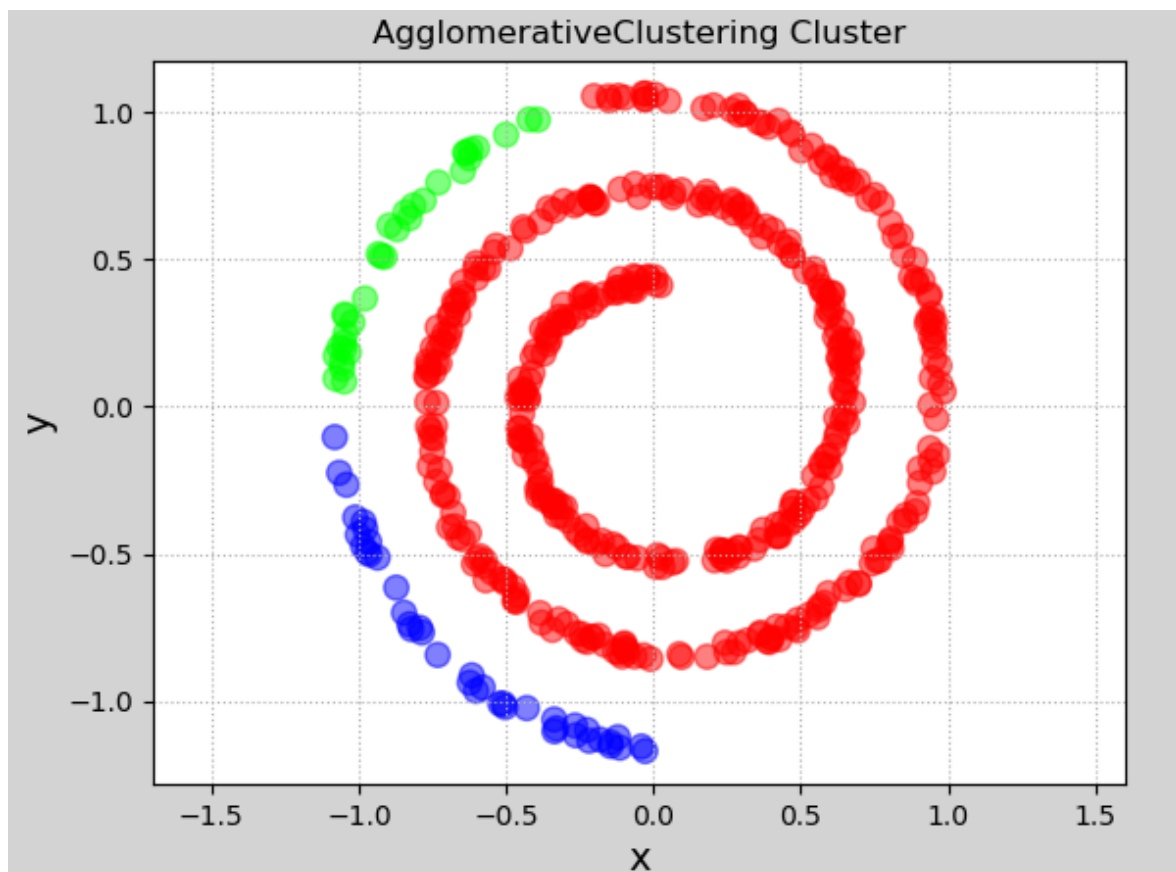


```

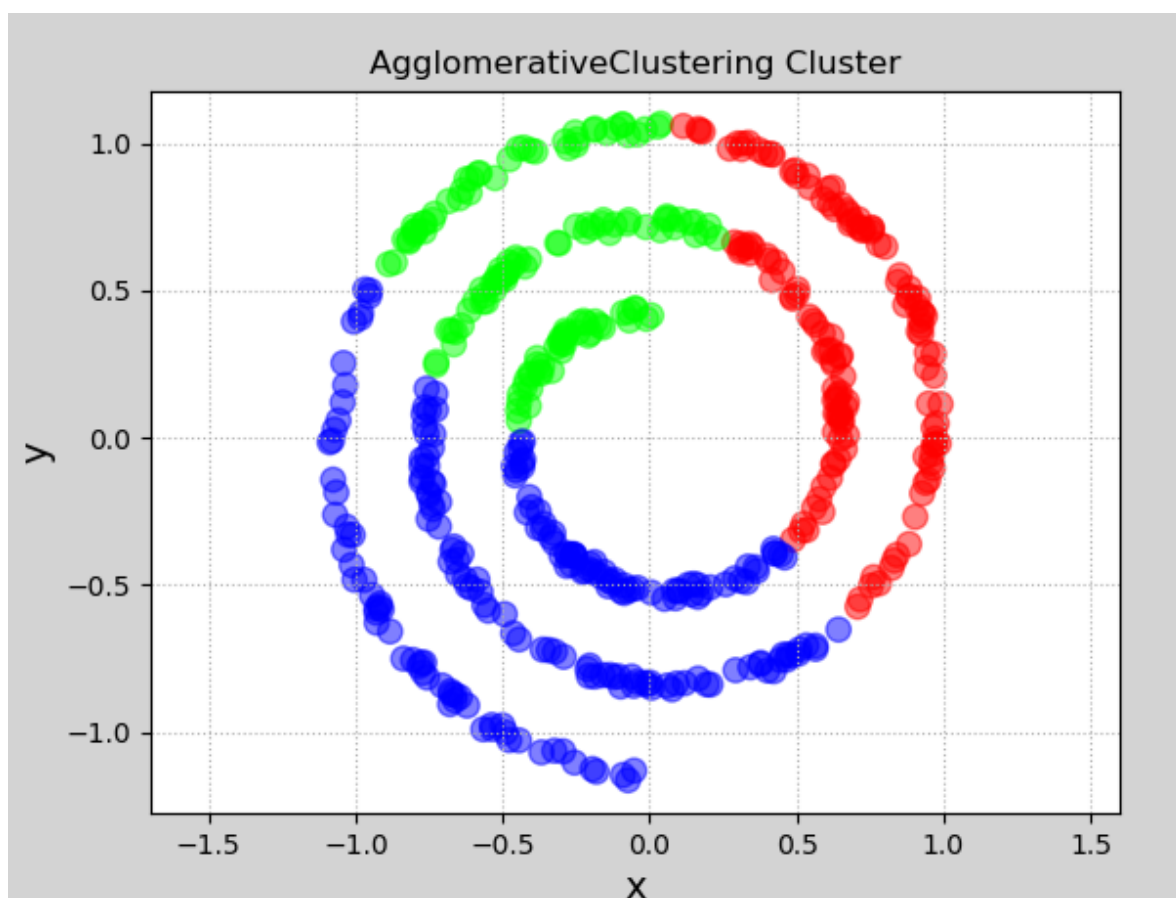
12 y = 0.05 * t * np.sin(t)
13 n = 0.05 * np.random.rand(n_sample, 2) # 产生随机噪声
14
15 x = np.hstack((x, y)) + n # 水平合并
16
17
18 # 无连续性凝聚层次聚类器
19 # model = sc.AgglomerativeClustering(n_clusters=3,
    linkage="average")
20 # model.fit(x) # 训练
21 # pred_y1 = model.labels_ # 聚类标签（聚类结果）
22
23 # 有连续性凝聚层次聚类器
24 conn = nb.kneighbors_graph(x, 10, include_self=False) # 创
    建每个样本的近邻集合
25 model = sc.AgglomerativeClustering(n_clusters=3,
26                                     linkage="average",
27                                     connectivity=conn) # 在凝聚过程中优先选择
    近邻中连续性最好的样本，优先凝聚
28 model.fit(x) # 训练
29 pred_y1 = model.labels_ # 聚类标签（聚类结果）
30
31
32 # 可视化
33 mp.figure("AgglomerativeClustering Cluster",
    facecolor="lightgray")
34 mp.title("AgglomerativeClustering Cluster")
35 mp.xlabel("x", fontsize=14)
36 mp.ylabel("y", fontsize=14)
37 mp.tick_params(labelsize=10)
38 mp.grid(linestyle=":")
39 mp.axis("equal")
40 mp.scatter(x[:, 0], x[:, 1], c=pred_y1, cmap="brg", s=80,
    alpha=0.5)
41 mp.show()

```

执行结果（有连续层次）：



因为是随机产生数据，该程序每次执行结果都不一样. 可以将代码22~24行注释打开，27~30行注释，就是一个非连续凝聚层次聚类. 执行结果：



## 示例5：利用SVM实现图像分类

- 数据集：包含两个目录train和test，每个目录下三个类别水果，apple、banana、grape
- 代码：

```
1  # -*- coding: utf-8 -*-
2  import os
3  import numpy as np
4  import cv2 as cv
5  import sklearn.metrics as sm
6  import sklearn.preprocessing as sp
7  import sklearn.svm as svm
8
9  name_dict = {"apple": 0, "banana": 1, "grape": 2}
10
11
12  # 读取图片、类别，并且存入字典
13  def search_samples(dir_path):
14      img_samples = {}
15
16      dirs = os.listdir(dir_path)
17      for d in dirs:
18          sub_dir_path = dir_path + "/" + d  # 拼接子目录完整
           路径
19          if not os.path.isdir(sub_dir_path):  # 不是子目录
20              continue
21
22          imgs = os.listdir(sub_dir_path)  # 列出子目录中所有
           文件
23          for img_file in imgs:
24              img_path = sub_dir_path + "/" + img_file  #
           拼接完整路径
25
26              if d in img_samples:  # 该类别已经在字典中
27                  img_samples[d].append(img_path)
28              else:
29                  img_list = []  # 定义空列表
30                  img_list.append(img_path)  # 将图像加入列表
31                  img_samples[d] = img_list
```

```

32
33     return img_samples
34
35
36 train_samples =
    search_samples('../data/fruits_tiny/train') # 搜索所有图像
    样本
37 train_x, train_y = [], []
38
39 # 加载训练集样本数据，训练模型，模型存储
40 for label, img_list in train_samples.items():
41     desc = np.array([])
42
43     for img_file in img_list:
44         # 读取原始图像，并转为灰度图像
45         print("读取样本:", img_file)
46         im = cv.imread(img_file)
47         im_gray = cv.cvtColor(im, cv.COLOR_BGR2GRAY)
48
49         # 调整大小
50         h, w = im_gray.shape[:2] # 取出高度、宽度
51         f = 200 / min(h, w) # 计算缩放比率
52         im_gray = cv.resize(im_gray, None, fx=f, fy=f) #
    图像缩放
53
54         # 计算特征矩阵
55         sift = cv.xfeatures2d.SIFT_create()
56         keypoints = sift.detect(im_gray)
57         _, desc = sift.compute(im_gray, keypoints)
58
59         # 添加到样本、输出数组
60         # print("desc.shape:", desc.shape)
61         desc = np.sum(desc, axis=0) # 0-列方向
62         train_x.append(desc) # 图像数据特征
63         train_y.append(name_dict[label]) # 标签
64
65 train_x = np.array(train_x)
66 train_y = np.array(train_y)

```

```
67 # print("train_y.shape:", train_y.shape)
68
69 # 定义模型、训练
70 print("开始训练.....")
71
72 model = svm.SVC(kernel='poly', degree=2)
73 model.fit(train_x, train_y)
74
75 print("训练结束.")
76
77 # 测试模型
78 test_samples = search_samples('../data/fruits_tiny/test')
79 test_x, test_y = [], []
80
81 # 读取测试数据，并计算特征值
82 for label, filenames in test_samples.items():
83     descs = np.array([])
84
85     for img_file in filenames:
86         print("读取测试样本:", img_file)
87
88         # 读取原始图像，并转为灰度图像
89         image = cv.imread(img_file)
90         gray = cv.cvtColor(image, cv.COLOR_BGR2GRAY)
91
92         # 调整大小
93         h, w = gray.shape[:2]
94         f = 200 / min(h, w)
95         gray = cv.resize(gray, None, fx=f, fy=f)
96
97         # 计算特征矩阵
98         sift = cv.xfeatures2d.SIFT_create()
99         keypoints = sift.detect(gray)
100         _, desc = sift.compute(gray, keypoints)
101
102         # 添加测试输入、输出数组
103         desc = np.sum(desc, axis=0) # 0-列方向
104         test_x.append(desc)
```

```

105 |         test_y.append(name_dict[label]) # 标签
106
107 | # 执行预测
108 | print("开始预测.....")
109 | pred_test_y = model.predict(test_x)
110 | print("预测结束.")
111
112 | # 打印分类报告
113 | print(sm.classification_report(test_y, pred_test_y))

```

执行结果：

```

1 | 中间打印省略.....
2
3 |           precision    recall  f1-score   support
4
5 |         0          1.00      0.80      0.89         10
6 |         1          1.00      1.00      1.00         10
7 |         2          0.83      1.00      0.91         10
8
9 |    accuracy                    0.93         30
10 |   macro avg          0.94      0.93      0.93         30
11 | weighted avg          0.94      0.93      0.93         30

```