



9. Experiment I

线性分类与回归

Linear Classification and Regression

第一节实验课（第 11 周，3 课时，学院楼 B406）将进行对线性分类与回归的实践能力考察。本小节将具体从 2 个子实验对线性分类与回归的实践要求进行说明。

实验内容 1：Setosa 鸢尾花和 Versicolour 鸢尾花二分类问题

实验内容 2：Diamonds 回归预测问题

本次实验所涉及编程内容不可使用现成函数进行，即分类器（回归器）、损失函数、优化求解过程均需编码实现。

实验报告提交（**请仔细对照**）：

- 格式：四号字体、单倍行距、中文宋体、英文新罗马字体，封面使用“实验报告封面.docx”，电子版编辑完后转为 pdf 格式发送，文件名 20AI+a+ 学号.pdf(例如：20AIa1033200101.pdf)；
- 时间节点：12 周周末（2023/05/07 18:00 前）
- 电子版发送至 raojiyong@stu.jiangnan.edu.cn，同步抄送 tianyang.xu@jiangnan.edu.cn，邮件标题:20AI+a+ 学号（例如：20AIa1033200101）

9.1 Setosa 鸢尾花和 Versicolour 鸢尾花二分类问题

9.1.1 数据：

结合第一章末尾处的内容，利用 ‘conda activate #name’ 进行实验环境激活，再安装 matplotlib、sklearn 等常用工具包。结合下面的代码熟悉鸢尾花 iris 数据集，实现对该数据

集的样本数、维度、类别等属性的分析与可视化。

```
1 # Code source: Gaël Varoquaux
2 # Modified for documentation by Jaques Grobler
3 # License: BSD 3 clause
4
5 import matplotlib.pyplot as plt
6 from mpl_toolkits.mplot3d import Axes3D
7 from sklearn import datasets
8 from sklearn.decomposition import PCA
9
10 # import some data to play with
11 iris = datasets.load_iris()
12 X = iris.data[:, :2] # we only take the first two features.
13 y = iris.target
14
15 x_min, x_max = X[:, 0].min() - 0.5, X[:, 0].max() + 0.5
16 y_min, y_max = X[:, 1].min() - 0.5, X[:, 1].max() + 0.5
17
18 plt.figure(2, figsize=(8, 6))
19 plt.clf()
20
21 # Plot the training points
22 plt.scatter(X[:, 0], X[:, 1], c=y, cmap=plt.cm.Set1, edgecolor="k")
23 plt.xlabel("Sepal length")
24 plt.ylabel("Sepal width")
25
26 plt.xlim(x_min, x_max)
27 plt.ylim(y_min, y_max)
28 plt.xticks(())
29 plt.yticks(())
30
31 # To get a better understanding of interaction of the dimensions
32 # plot the first three PCA dimensions
33 fig = plt.figure(1, figsize=(8, 6))
34 ax = Axes3D(fig, elev=-150, azim=110)
35 X_reduced = PCA(n_components=3).fit_transform(iris.data)
36 ax.scatter(
37     X_reduced[:, 0],
38     X_reduced[:, 1],
39     X_reduced[:, 2],
40     c=y,
41     cmap=plt.cm.Set1,
42     edgecolor="k",
```

```
43     s=40 ,
44 )
45 ax.set_title("First three PCA directions")
46 ax.set_xlabel("1st eigenvector")
47 ax.w_xaxis.set_ticklabels([])
48 ax.set_ylabel("2nd eigenvector")
49 ax.w_yaxis.set_ticklabels([])
50 ax.set_zlabel("3rd eigenvector")
51 ax.w_zaxis.set_ticklabels([])
52
53 plt.show()
```

9.1.2 实验课具体要求：

编写线性分类器模型，编写损失函数、编写梯度反传函数；实现 iris 中 Setosa 鸢尾花和 Versicolour 鸢尾花的二分类学习；训练数据为 Setosa 鸢尾花和 Versicolour 鸢尾花各自前 30 条数据（共 60 条数据）；测试数据为 Setosa 鸢尾花和 Versicolour 鸢尾花各自后 20 条数据（共 40 条数据）。

9.1.3 评分依据：

- ★★★ 完成分类器、损失函数、梯度的编写，调整训练参数实现测试集上的高精度。
- ★★☆☆ 完成分类器、损失函数、梯度的编写，调通训练和测试过程。
- ★☆☆ 利用已有函数的调用，调通训练和测试过程。

9.1.4 实验报告：

- 核心分类器、损失函数、梯度反传的代码及注释
- 训练测试的过程分析（可选：不同学习率、不同损失函数对结果的影响）
- 实验分析（如：哪种特征更加有利于中 Setosa 鸢尾花和 Versicolour 鸢尾花的二分类）

9.1.5 深入学习：

Kaggle 上有大量关于此数据集的代码 (<https://www.kaggle.com/datasets/uciml/iris>)，可以通过阅读他人代码掌握看问题的不同角度。

9.2 Diamonds 回归预测问题

9.2.1 数据：

利用 seaborn 获取 diamonds 数据集，通过可视化初步了解数据内容，熟悉如下各种数据分析可视化方式。<https://www.kaggle.com/shivam2503/diamonds>

```
1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 sns.get_dataset_names()
5 diamond=sns.load_dataset('diamonds', cache=True, data_home=None)
6 # 网络原因或导致无法直接下载; 则可自行下载后导入
7
8 # Scatter Plot
9 sns.scatterplot(x="carat", y="price", hue="cut", size="depth", data=diamond)
10
11 plt.xlabel('Carat Weight')
12 plt.ylabel('Price')
13 plt.title('Scatter Plot: Carat Vs Price (Based on Cut)')
14
15 # Rug Plot
16 sns.rugplot(x="carat", y="price", hue="cut", data=diamond)
17
18 # Scatter Plot with Rug Plot: Carat Vs Price (Based on Cut)
19 sns.scatterplot(x="carat", y="price", hue="cut", size="depth", data=diamond)
20
21 sns.rugplot(x="carat", y="price", hue="cut", data=diamond)
22
23 # Scatter Plot with Rug Plot: Carat Vs Price (Based on Color)
24 sns.scatterplot(x="carat", y="price", hue="color", size="depth", data=
    diamond)
25
26 sns.rugplot(x="carat", y="price", hue="color", data=diamond)
27
28 # Scatter Plot with Rug Plot: Carat Vs Price (Based on Clarity)
29 sns.scatterplot(x="carat", y="price", hue="clarity", size="depth", data=
    diamond)
30
31 sns.rugplot(x="carat", y="price", hue="clarity", data=diamond)
32
33 # Strip plot
34 fig, st_axis = plt.subplots(2, 3, figsize=(16,8))
35 sns.stripplot(x="cut", y="carat", data=diamond, ax = st_axis[0,0])
36 sns.stripplot(x="color", y="carat", data=diamond, ax = st_axis[0,1])
37 sns.stripplot(x="clarity", y="carat", data=diamond, ax = st_axis[0,2])
38 sns.stripplot(x="cut", y="price", data=diamond, ax = st_axis[1,0])
39 sns.stripplot(x="color", y="price", data=diamond, ax = st_axis[1,1])
40 sns.stripplot(x="clarity", y="price", data=diamond, ax = st_axis[1,2])
41
```

```
42 # Bar Plot
43 fig, ba_axis = plt.subplots(1, 3, figsize=(16,12))
44 sns.barplot(x="cut", y="price", data=diamond, ax = ba_axis[0], color='Red')
45 sns.barplot(x="color", y="price", data=diamond, ax = ba_axis[1], color='
    salmon')
46 sns.barplot(x="clarity", y="price", data=diamond, ax = ba_axis[2], color='
    purple')
47
48 # Pair Plot
49 sns.pairplot(diamond, plot_kws=dict(marker="+", linewidth=1), diag_kws=dict(
    fill=False))
50
51 # Line Plot
52 sns.lineplot(x="carat", y="x", data=diamond, color='salmon')
53 sns.lineplot(x="carat", y="y", data=diamond, color='purple')
54 sns.lineplot(x="carat", y="z", data=diamond, color='Red')
55 sns.lineplot(x="carat", y="depth", data=diamond, sort=False, lw=1, color='
    blue')
56
57 sns.lineplot(x="carat", y="x", data=diamond, hue='cut')
58 sns.lineplot(x="carat", y="x", data=diamond, hue='color')
59 sns.lineplot(x="carat", y="x", data=diamond, hue='clarity')
60 sns.lineplot(x="carat", y="y", data=diamond, hue='cut')
61 sns.lineplot(x="carat", y="y", data=diamond, hue='color')
62 sns.lineplot(x="carat", y="y", data=diamond, hue='clarity')
63 sns.lineplot(x="carat", y="z", data=diamond, hue='cut')
64 sns.lineplot(x="carat", y="z", data=diamond, hue='color')
65 sns.lineplot(x="carat", y="z", data=diamond, hue='clarity')
66
67 # Count Plot
68 sns.histplot(data=diamond, x='carat', color='black')
69 sns.histplot(data=diamond, y='price', color='salmon')
70 sns.histplot(data=diamond, x='depth', hue='cut', color='cyan')
71 sns.histplot(data=diamond, x='table', hue='color', multiple='stack', color='
    red')
72 sns.histplot(data=diamond, x='x', hue='clarity', element='step', color='blue
    ')
73 sns.histplot(data=diamond, x='y', hue='clarity', element='poly', color='pink
    ')
74 sns.histplot(data=diamond, x='z', hue='cut', element="step", stat="density",
    common_norm=False, color='orange')
75
76 # Heat Map
```

```

77 sns.heatmap(diamond.corr(), linecolor='white', linewidths=2,annot=True)
78
79 # ECDF Plot
80 sns.ecdfplot(data=diamond, x="carat", stat="count")
81 sns.ecdfplot(data=diamond, x="price", stat="count")
82 sns.ecdfplot(data=diamond, x="cut", stat="count")
83 sns.ecdfplot(data=diamond, x="color", stat="count")
84 sns.ecdfplot(data=diamond, x="clarity", stat="count")
85 sns.ecdfplot(data=diamond, x="depth", stat="count", complementary=True)
86 sns.ecdfplot(data=diamond, x="table", stat="count", complementary=True)
87
88 # Box Plot
89 sns.boxplot(x="carat", y="cut", data=diamond, color="salmon", width=0.5,
90            linewidth=2.5)
91 sns.boxplot(x="carat", y="color", data=diamond, color="royalblue", width
92            =0.5, linewidth=2.5)
93 sns.boxplot(x="carat", y="clarity", data=diamond, color="cyan", width=0.5,
94            linewidth=2.5)
95
96 # Violin plot
97 sns.violinplot(x="cut", y="carat", data=diamond, color="cyan", width=1)
98 sns.violinplot(x="color", y="carat", data=diamond, palette="muted", color="
99            royalblue", width=1)
100 sns.violinplot(x="clarity", y="carat", data=diamond, palette="muted", split=
101            True, inner="quartile", color="red", width=1)
102
103 # Point plot
104 sns.pointplot(x="cut", y="carat", data=diamond)
105 sns.pointplot(x="cut", y="carat", hue="color", data=diamond, dodge=True)
106 sns.pointplot(x="cut", y="carat", hue="clarity", data=diamond, dodge=True)
107 sns.pointplot(x="color", y="carat", data=diamond)
108 sns.pointplot(x="color", y="carat", hue="cut", data=diamond, join=False)
109 sns.pointplot(x="color", y="carat", hue="clarity", data=diamond, join=False)
110 sns.pointplot(x="carat", y="clarity", data=diamond)
111 sns.pointplot(x="carat", y="clarity", hue="cut", data=diamond)
112 sns.pointplot(x="carat", y="clarity", hue="color", data=diamond)
113
114 # Reg Plot
115 fig, re_axis = plt.subplots(2, 2,figsize=(12,12))
116 sns.regplot(x = 'carat', y = 'price', data=diamond, ax = re_axis[0,0], color
117            ='royalblue')
118 sns.regplot(x = 'x', y = 'price', data=diamond, ax = re_axis[0,1], color='
119            salmon')

```



```

113 sns.regplot(x = 'y', y = 'price', data=diamond, ax = re_axis[1,0], color='
    green')
114 sns.regplot(x = 'z', y = 'price', data=diamond, ax = re_axis[1,1], color='
    cyan')
115
116 # Kernel Density Plot
117 sns.kdeplot(data=diamond, x="carat", hue="cut")
118 sns.kdeplot(data=diamond, x="carat", hue="color", multiple="stack")
119 sns.kdeplot(data=diamond, x="carat", hue="clarity", multiple="fill")
120
121 # Jointplot
122 sns.jointplot(data=diamond, x='carat', y='depth', marker='^', hue='cut')
123
124 sns.jointplot(data=diamond, x='carat', y='depth', marker='+', hue='color')
125
126 sns.jointplot(data=diamond, x='carat', y='depth', marker='>', hue='clarity')
127
128 # JointGrid
129 jg_cut = sns.JointGrid(data=diamond, x="carat", y="price", hue="cut")
130 jg_cut.plot(sns.scatterplot, sns.histplot, alpha=.7, edgecolor=".2",
    linewidth=.5)
131
132 jg_color = sns.JointGrid(data=diamond, x="carat", y="price", hue="color")
133 jg_color.plot(sns.scatterplot, sns.boxplot)
134
135 jg_clarity = sns.JointGrid(data=diamond, x="carat", y="price", hue="clarity"
    )
136 jg_clarity.plot(sns.scatterplot, sns.kdeplot)

```

9.2.2 实验课具体要求：

编写线性回归器模型，编写损失函数、编写梯度反传函数；实现 diamonds 特征（carat, cut, color, clarity, depth, table, x, y, z）对价格（price）的预测；训练数据为第 1-40000 条数据中所有合数索引对应的数据；测试数据为第 1-40000 条数据中所有质数索引对应的数据（4203 个）。

[TIPS：求解 n 以内所有质数的函数 ↓]

```

1 import math
2 def func_get_prime(n):
3     return filter(lambda x: not [x%i for i in range(2, int(math.sqrt(x))+1) if
        x%i ==0], range(2,n+1))

```

9.2.3 评分依据:

- ★★★ 完成回归器、损失函数、梯度的编写，调整训练参数实现测试集上的高精度。
- ★★★☆ 完成回归器、损失函数、梯度的编写，调通训练和测试过程。
- ★☆☆ 利用已有函数的调用，调通训练和测试过程。

9.2.4 实验报告:

- 核心回归器、损失函数、梯度反传的代码及注释
- 训练测试的过程分析（可选：不同学习率、不同损失函数对结果的影响）
- 实验分析（如：哪种特征对钻石价格的影响是正相关/负相关的）