Homework due Jul 13, 2021 22:00 +06

Exercise 1

1/1 point (graded)

Our first step is to import the dataset.

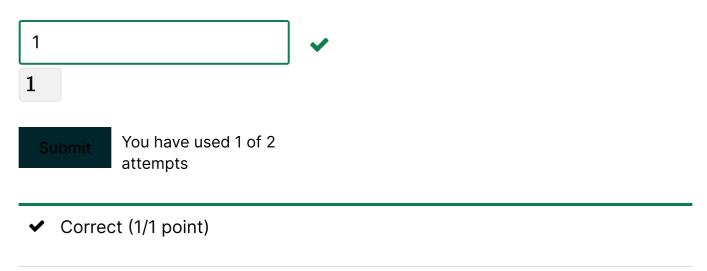
Instructions

Read in the data as a pandas dataframe using pd.read_csv. The data can be found at this link from anywhere and at this link from within the courseware.

This code will get you started:

```
import pandas as pd
# write your code here!
```

Taking a look at the first 5 rows of the dataset, how many wines in those 5 rows are considered high quality?



Exercise 2

1/1 point (graded)

Next, we will inspect the dataset and perform some mild data cleaning.

Instructions

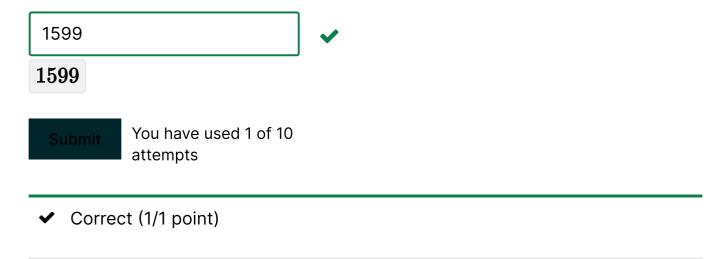
In order to get all numeric data, we will change the color column to an is_red column.

- If color == 'red', we will encode a 1 for is red.
- If color == 'white', we will encode a 0 for is_red.

Create this new column, <code>is_red</code>. Drop the <code>color</code> column as well as <code>quality</code> and <code>high_quality</code>. We will predict the quality of wine using the numeric data in a later exercise

Store this all numeric data in a pandas dataframe called numeric data.

How many red wines are in the dataset?



Exercise 3

1/1 point (graded)

We want to ensure that each variable contributes equally to the kNN classifier, so we will need to scale the data by subtracting the mean of each variable (column) and dividing each variable (column) by its standard deviation. Then, we will use principal components to take a linear snapshot of the data from several different angles, with each snapshot ordered by how well it aligns with variation in the data. In this exercise, we will scale the numeric data and extract the first two principal components.

Instructions

- Scale the data using the sklearn.preprocessing function scale() on numeric data.
- Convert this to a pandas dataframe, and store it as numeric data.
- Include the numeric variable names using the parameter columns = numeric_data.columns.
- Use the sklearn.decomposition module PCA() and store it as pca.
- Use the fit_transform() function to extract the first two principal components from the data, and store them as principal_components.

Note: You may get a DataConversionWarning, but you can safely ignore it.

Fill in this code as you work:

```
import sklearn.preprocessing
scaled_data =
numeric_data =

import sklearn.decomposition
pca =
principal_components =
```

What is the shape of the new dataset?

Enter the first number here.



Enter the second number here.



✓ Correct (1/1 point)

Exercise 4

1/1 point (graded)

In Exercise 4, we will plot the first two principal components of the covariates in the dataset. The high and low quality wines will be colored using red and blue, respectively.

Instructions

The first two principal components can be accessed using principal_components[:,0] and principal_components[:,1]. Store these as x and y respectively, and make a scatter plot of these first two principal components.

Consider how well the two groups of wines are separated by the first two principal components.

Fill in your code where indicated to make the plot:

```
import matplotlib.pyplot as plt
from matplotlib.colors import ListedColormap
from matplotlib.backends.backend_pdf import PdfPages
observation_colormap = ListedColormap(['red', 'blue'])
x = # Enter your code here!
y = # Enter your code here!

plt.title("Principal Components of Wine")
plt.scatter(x, y, alpha = 0.2,
    c = data['high_quality'], cmap = observation_colormap, edgecolors = 'none')
plt.xlim(-8, 8); plt.ylim(-8, 8)
plt.xlabel("Principal Component 1")
plt.ylabel("Principal Component 2")
plt.show()
```

Could you easily draw a linear boundary between the high and low quality wines

using the first two principal components?

✓ Yes

No

✓

Submit You have used 1 of 1 attempt

✓ Correct (1/1 point)

Exercise 5

1/1 point (graded)

In Exercise 5, we will create a function that calculates the accuracy between predictions and outcomes.

Instructions

- Create a function accuracy(predictions, outcomes) that takes two lists of the same size as arguments and returns a single number, which is the percentage of elements that are equal for the two lists.
- Use accuracy to compare the percentage of similar elements in the x and y numpy arrays defined below.
- Print your answer.

Here's the sample code to get you started:

```
import numpy as np
np.random.seed(1) # do not change this!

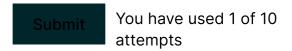
x = np.random.randint(0, 2, 1000)

y = np.random.randint(0 ,2, 1000)

def accuracy(predictions, outcomes):
    # write your code here!
```

What is the accuracy of the x predictions on the "true" outcomes y?





✓ Correct (1/1 point)