

Assignment 1

November 2, 2017

You are currently looking at **version 1.1** of this notebook. To download notebooks and datafiles, as well as get help on Jupyter notebooks in the Coursera platform, visit the [Jupyter Notebook FAQ](#) course resource.

1 Assignment 1 - Introduction to Machine Learning

For this assignment, you will be using the Breast Cancer Wisconsin (Diagnostic) Database to create a classifier that can help diagnose patients. First, read through the description of the dataset (below).

```
In [3]: import numpy as np
import pandas as pd
from sklearn.datasets import load_breast_cancer

cancer = load_breast_cancer()

print(cancer.DESCR) # Print the data set description
print(type(cancer.data))
```

Breast Cancer Wisconsin (Diagnostic) Database
=====

Notes

Data Set Characteristics:

:Number of Instances: 569

:Number of Attributes: 30 numeric, predictive attributes and the class

:Attribute Information:

- radius (mean of distances from center to points on the perimeter)
- texture (standard deviation of gray-scale values)
- perimeter
- area

- smoothness (local variation in radius lengths)
- compactness ($\text{perimeter}^2 / \text{area} - 1.0$)
- concavity (severity of concave portions of the contour)
- concave points (number of concave portions of the contour)
- symmetry
- fractal dimension ("coastline approximation" - 1)

The mean, standard error, and "worst" or largest (mean of the three largest values) of these features were computed for each image, resulting in 30 features. For instance, field 3 is Mean Radius, field 13 is Radius SE, field 23 is Worst Radius.

- class:
 - WDBC-Malignant
 - WDBC-Benign

:Summary Statistics:

	Min	Max
radius (mean):	6.981	28.11
texture (mean):	9.71	39.28
perimeter (mean):	43.79	188.5
area (mean):	143.5	2501.0
smoothness (mean):	0.053	0.163
compactness (mean):	0.019	0.345
concavity (mean):	0.0	0.427
concave points (mean):	0.0	0.201
symmetry (mean):	0.106	0.304
fractal dimension (mean):	0.05	0.097
radius (standard error):	0.112	2.873
texture (standard error):	0.36	4.885
perimeter (standard error):	0.757	21.98
area (standard error):	6.802	542.2
smoothness (standard error):	0.002	0.031
compactness (standard error):	0.002	0.135
concavity (standard error):	0.0	0.396
concave points (standard error):	0.0	0.053
symmetry (standard error):	0.008	0.079
fractal dimension (standard error):	0.001	0.03
radius (worst):	7.93	36.04
texture (worst):	12.02	49.54
perimeter (worst):	50.41	251.2
area (worst):	185.2	4254.0
smoothness (worst):	0.071	0.223
compactness (worst):	0.027	1.058
concavity (worst):	0.0	1.252

concave points (worst):	0.0	0.291
symmetry (worst):	0.156	0.664
fractal dimension (worst):	0.055	0.208
=====	=====	=====

:Missing Attribute Values: None

:Class Distribution: 212 - Malignant, 357 - Benign

:Creator: Dr. William H. Wolberg, W. Nick Street, Olvi L. Mangasarian

:Donor: Nick Street

:Date: November, 1995

This is a copy of UCI ML Breast Cancer Wisconsin (Diagnostic) datasets.
<https://goo.gl/U2Uwz2>

Features are computed from a digitized image of a fine needle aspirate (FNA) of a breast mass. They describe characteristics of the cell nuclei present in the image.

Separating plane described above was obtained using Multisurface Method-Tree (MSM-T) [K. P. Bennett, "Decision Tree Construction Via Linear Programming." Proceedings of the 4th Midwest Artificial Intelligence and Cognitive Science Society, pp. 97-101, 1992], a classification method which uses linear programming to construct a decision tree. Relevant features were selected using an exhaustive search in the space of 1-4 features and 1-3 separating planes.

The actual linear program used to obtain the separating plane in the 3-dimensional space is that described in:

[K. P. Bennett and O. L. Mangasarian: "Robust Linear Programming Discrimination of Two Linearly Inseparable Sets", Optimization Methods and Software 1, 1992, 23-34].

This database is also available through the UW CS ftp server:

```
ftp ftp.cs.wisc.edu
cd math-prog/cpo-dataset/machine-learn/WDBC/
```

References

- W.N. Street, W.H. Wolberg and O.L. Mangasarian. Nuclear feature extraction for breast tumor diagnosis. IS&T/SPIE 1993 International Symposium on Electronic Imaging: Science and Technology, volume 1905, pages 861-870, San Jose, CA, 1993.

- O.L. Mangasarian, W.N. Street and W.H. Wolberg. Breast cancer diagnosis and prognosis via linear programming. Operations Research, 43(4), pages 570-577, July-August 1995.
- W.H. Wolberg, W.N. Street, and O.L. Mangasarian. Machine learning techniques to diagnose breast cancer from fine-needle aspirates. Cancer Letters 77 (1994) 163-171.

```
<class 'numpy.ndarray'>
```

The object returned by `load_breast_cancer()` is a scikit-learn Bunch object, which is similar to a dictionary.

```
In [2]: cancer.keys()
```

```
Out[2]: dict_keys(['target_names', 'DESCR', 'data', 'target', 'feature_names'])
```

1.0.1 Question 0 (Example)

How many features does the breast cancer dataset have?

This function should return an integer.

```
In [3]: # You should write your whole answer within the function provided. The auto
# this function and compare the return value against the correct solution v
def answer_zero():
    # This function returns the number of features of the breast cancer dat
    # The assignment question description will tell you the general format
    return len(cancer['feature_names'])

# You can examine what your function returns by calling it in the cell. If
# about the assignment formats, check out the discussion forums for any FAQ
answer_zero()
```

```
Out[3]: 30
```

```
In [33]: def answer_one():
x = pd.DataFrame(cancer.target, columns=['target'])
y = pd.DataFrame(cancer.data, columns=cancer['feature_names'])
return pd.merge(y, x, right_index=True, left_index=True)
```

```
(569, 31)
```

1.0.2 Question 1

Scikit-learn works with lists, numpy arrays, scipy-sparse matrices, and pandas DataFrames, so converting the dataset to a DataFrame is not necessary for training this model. Using a DataFrame does however help make many things easier such as munging data, so let's practice creating a classifier with a pandas DataFrame.

Convert the `sklearn.dataset.cancer` to a `DataFrame`.

This function should return a (569, 31) DataFrame with columns =

```
['mean radius', 'mean texture', 'mean perimeter', 'mean area',  
'mean smoothness', 'mean compactness', 'mean concavity',  
'mean concave points', 'mean symmetry', 'mean fractal dimension',  
'radius error', 'texture error', 'perimeter error', 'area error',  
'smoothness error', 'compactness error', 'concavity error',  
'concave points error', 'symmetry error', 'fractal dimension error',  
'worst radius', 'worst texture', 'worst perimeter', 'worst area',  
'worst smoothness', 'worst compactness', 'worst concavity',  
'worst concave points', 'worst symmetry', 'worst fractal dimension',  
'target']
```

and index =

```
RangeIndex(start=0, stop=569, step=1)
```

1.0.3 Question 2

What is the class distribution? (i.e. how many instances of malignant (encoded 0) and how many benign (encoded 1)?)

This function should return a Series named target of length 2 with integer values and index =

```
['malignant', 'benign']
```

```
In [38]: def answer_two():  
         cancerdf = answer_one()  
         malignant = len(cancerdf.where(cancerdf['target'] == 0).dropna())  
         benign = len(cancerdf) - malignant  
         return pd.Series([malignant, benign], index=['malignant', 'benign'])  
  
         print(answer_two())  
  
malignant    212  
benign       357  
dtype: int64
```

1.0.4 Question 3

Split the `DataFrame` into `X` (the data) and `y` (the labels).

*This function should return a tuple of length 2: (X, y), where * X has shape (569, 30) * y has shape (569,).*

```
In [44]: def answer_three():  
         cancerdf = answer_one()  
         y = cancerdf['target']  
         X = cancerdf.filter(regex=r'^(?!target)')  
  
         return X, y
```

1.0.5 Question 4

Using `train_test_split`, split `X` and `y` into training and test sets (`X_train`, `X_test`, `y_train`, and `y_test`).

Set the random number generator state to 0 using `random_state=0` to make sure your results match the autograder!

*This function should return a tuple of length 4: (`X_train`, `X_test`, `y_train`, `y_test`), where * `X_train` has shape (426, 30) * `X_test` has shape (143, 30) * `y_train` has shape (426,) * `y_test` has shape (143,)*

```
In [43]: from sklearn.model_selection import train_test_split
```

```
def answer_four():
    X, y = answer_three()
    X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=0)
    # Your code here

    return X_train, X_test, y_train, y_test
```

1.0.6 Question 5

Using `KNeighborsClassifier`, fit a k-nearest neighbors (knn) classifier with `X_train`, `y_train` and using one nearest neighbor (`n_neighbors = 1`).

This function should return a `sklearn.neighbors.classification.KNeighborsClassifier`.

```
In [45]: from sklearn.neighbors import KNeighborsClassifier
```

```
def answer_five():
    X_train, X_test, y_train, y_test = answer_four()
    knn = KNeighborsClassifier(n_neighbors=1)
    knn.fit(X_train, y_train)
    return knn
```

1.0.7 Question 6

Using your knn classifier, predict the class label using the mean value for each feature.

Hint: You can use `cancerdf.mean()[:-1].values.reshape(1, -1)` which gets the mean value for each feature, ignores the target column, and reshapes the data from 1 dimension to 2 (necessary for the `predict` method of `KNeighborsClassifier`).

This function should return a numpy array either `array([0.])` or `array([1.])`

```
In [46]: def answer_six():
    cancerdf = answer_one()
    means = cancerdf.mean()[ :-1].values.reshape(1, -1)
    knn = answer_five()

    return knn.predict(means)
```

1.0.8 Question 7

Using your knn classifier, predict the class labels for the test set `X_test`.

This function should return a numpy array with shape (143,) and values either 0.0 or 1.0.

[illegible]

[illegible]

[illegible]

[illegible]

```

/opt/conda/lib/python3.5/site-packages/sklearn/utils/validation.py:395: Deprecation
  DeprecationWarning)
/opt/conda/lib/python3.5/site-packages/sklearn/utils/validation.py:395: Deprecation
  DeprecationWarning)
/opt/conda/lib/python3.5/site-packages/sklearn/utils/validation.py:395: Deprecation
  DeprecationWarning)
/opt/conda/lib/python3.5/site-packages/sklearn/utils/validation.py:395: Deprecation
  DeprecationWarning)
/opt/conda/lib/python3.5/site-packages/sklearn/utils/validation.py:395: Deprecation
  DeprecationWarning)
/opt/conda/lib/python3.5/site-packages/sklearn/utils/validation.py:395: Deprecation
  DeprecationWarning)

```

```

Out[59]: 512      1
          457      1
          439      1
          298      0
           37      1
          515      1
          382      1
          310      1
          538      1
          345      1
          421      0
           90      1
          412      1
          157      1
           89      0
          172      0
          318      1
          233      0
          389      0
          250      0
           31      0
          283      1
          482      1
          211      1
          372      0
          401      1
          159      1
           14      1
          364      1
          337      0
          ..
          500      1
          338      1
          427      1

```

```

406    0
96     1
490    1
384    1
281    1
325    1
190    1
380    1
366    0
469    1
225    1
271    1
547    1
550    1
492    0
185    1
306    1
208    1
242    1
313    1
542    1
514    0
236    0
113    1
527    1
76     1
162    0
Name: mean radius, dtype: int64

```

1.0.9 Question 8

Find the score (mean accuracy) of your knn classifier using `X_test` and `y_test`.
This function should return a float between 0 and 1

```

In [48]: def answer_eight():
          X_train, X_test, y_train, y_test = answer_four()
          knn = answer_five()

          return knn.score(X_test, y_test)

```

1.0.10 Optional plot

Try using the plotting function below to visualize the differet prediction scores between training and test sets, as well as malignant and benign cells.

```

In [ ]: def accuracy_plot():
          import matplotlib.pyplot as plt

          %matplotlib notebook

```

```

X_train, X_test, y_train, y_test = answer_four()

# Find the training and testing accuracies by target value (i.e. malignant/benign)
mal_train_X = X_train[y_train==0]
mal_train_y = y_train[y_train==0]
ben_train_X = X_train[y_train==1]
ben_train_y = y_train[y_train==1]

mal_test_X = X_test[y_test==0]
mal_test_y = y_test[y_test==0]
ben_test_X = X_test[y_test==1]
ben_test_y = y_test[y_test==1]

knn = answer_five()

scores = [knn.score(mal_train_X, mal_train_y), knn.score(ben_train_X, ben_train_y),
          knn.score(mal_test_X, mal_test_y), knn.score(ben_test_X, ben_test_y)]

plt.figure()

# Plot the scores as a bar chart
bars = plt.bar(np.arange(4), scores, color=['#4c72b0', '#4c72b0', '#55a862', '#55a862'])

# directly label the score onto the bars
for bar in bars:
    height = bar.get_height()
    plt.gca().text(bar.get_x() + bar.get_width()/2, height*.90, '{0:.1f}'.format(score))
    ha='center', color='w', fontsize=11)

# remove all the ticks (both axes), and tick labels on the Y axis
plt.tick_params(top='off', bottom='off', left='off', right='off', labelbottom='off')

# remove the frame of the chart
for spine in plt.gca().spines.values():
    spine.set_visible(False)

plt.xticks([0,1,2,3], ['Malignant\nTraining', 'Benign\nTraining', 'Malignant\nTest', 'Benign\nTest'])
plt.title('Training and Test Accuracies for Malignant and Benign Cells')

```

```

In [ ]: # Uncomment the plotting function to see the visualization,
        # Comment out the plotting function when submitting your notebook for grading

        #accuracy_plot()

```

```

In [ ]:

```