# Lab 10 06-10-2022 k-NN cancer

October 6, 2022

## 0.1 Breast Cancer

```
[1]: import numpy as np
     import pandas as pd
     !pip install scikit-learn
```

```
Defaulting to user installation because normal site-packages is not writeable
Requirement already satisfied: scikit-learn in
/opt/anaconda3/lib/python3.9/site-packages (1.0.2)
Requirement already satisfied: numpy>=1.14.6 in
/opt/anaconda3/lib/python3.9/site-packages (from scikit-learn) (1.21.5)
Requirement already satisfied: threadpoolctl>=2.0.0 in
/opt/anaconda3/lib/python3.9/site-packages (from scikit-learn) (2.2.0)
Requirement already satisfied: joblib>=0.11 in
/opt/anaconda3/lib/python3.9/site-packages (from scikit-learn) (1.1.0)
Requirement already satisfied: scipy>=1.1.0 in
/opt/anaconda3/lib/python3.9/site-packages (from scikit-learn) (1.7.3)
```

```
[3]: ds = pd.read_csv('BreastCancer.csv')
     print(ds)
```

```
           id diagnosis  radius_mean  texture_mean  perimeter_mean  area_mean  \
0      842302         M        17.99         10.38          122.80     1001.0
1      842517         M        20.57         17.77          132.90     1326.0
2    84300903         M        19.69         21.25          130.00     1203.0
3    84348301         M        11.42         20.38           77.58      386.1
4    84358402         M        20.29         14.34          135.10     1297.0
..        ...       ...          ...           ...             ...        ...
564    926424         M        21.56         22.39          142.00     1479.0
565    926682         M        20.13         28.25          131.20     1261.0
566    926954         M        16.60         28.08          108.30      858.1
567    927241         M        20.60         29.33          140.10     1265.0
568     92751         B         7.76         24.54           47.92      181.0

     smoothness_mean  compactness_mean  concavity_mean  concave points_mean  \
0            0.11840           0.27760         0.30010              0.14710
1            0.08474           0.07864         0.08690              0.07017
2            0.10960           0.15990         0.19740              0.12790
3            0.14250           0.28390         0.24140              0.10520
```

1

|     |         |         |         |         |
|-----|---------|---------|---------|---------|
| 4   | 0.10030 | 0.13280 | 0.19800 | 0.10430 |
| ..  | …       | …       | …       | …       |
| 564 | 0.11100 | 0.11590 | 0.24390 | 0.13890 |
| 565 | 0.09780 | 0.10340 | 0.14400 | 0.09791 |
| 566 | 0.08455 | 0.10230 | 0.09251 | 0.05302 |
| 567 | 0.11780 | 0.27700 | 0.35140 | 0.15200 |
| 568 | 0.05263 | 0.04362 | 0.00000 | 0.00000 |

|     |     | texture_worst | perimeter_worst | area_worst | smoothness_worst \ |
|-----|-----|---------------|-----------------|------------|--------------------|
| 0   | …   | 17.33         | 184.60          | 2019.0     | 0.16220            |
| 1   | …   | 23.41         | 158.80          | 1956.0     | 0.12380            |
| 2   | …   | 25.53         | 152.50          | 1709.0     | 0.14440            |
| 3   | …   | 26.50         | 98.87           | 567.7      | 0.20980            |
| 4   | …   | 16.67         | 152.20          | 1575.0     | 0.13740            |
| ..  | …   | …             | …               | …          | …                  |
| 564 | …   | 26.40         | 166.10          | 2027.0     | 0.14100            |
| 565 | …   | 38.25         | 155.00          | 1731.0     | 0.11660            |
| 566 | …   | 34.12         | 126.70          | 1124.0     | 0.11390            |
| 567 | …   | 39.42         | 184.60          | 1821.0     | 0.16500            |
| 568 | …   | 30.37         | 59.16           | 268.6      | 0.08996            |

|     | compactness_worst | concavity_worst | concave points_worst | symmetry_worst \ |
|-----|-------------------|-----------------|----------------------|------------------|
| 0   | 0.66560           | 0.7119          | 0.2654               | 0.4601           |
| 1   | 0.18660           | 0.2416          | 0.1860               | 0.2750           |
| 2   | 0.42450           | 0.4504          | 0.2430               | 0.3613           |
| 3   | 0.86630           | 0.6869          | 0.2575               | 0.6638           |
| 4   | 0.20500           | 0.4000          | 0.1625               | 0.2364           |
| ..  | …                 | …               | …                    | …                |
| 564 | 0.21130           | 0.4107          | 0.2216               | 0.2060           |
| 565 | 0.19220           | 0.3215          | 0.1628               | 0.2572           |
| 566 | 0.30940           | 0.3403          | 0.1418               | 0.2218           |
| 567 | 0.86810           | 0.9387          | 0.2650               | 0.4087           |
| 568 | 0.06444           | 0.0000          | 0.0000               | 0.2871           |

|     | fractal_dimension_worst | Unnamed: 32 |
|-----|-------------------------|-------------|
| 0   | 0.11890                 | NaN         |
| 1   | 0.08902                 | NaN         |
| 2   | 0.08758                 | NaN         |
| 3   | 0.17300                 | NaN         |
| 4   | 0.07678                 | NaN         |
| ..  | …                       | …           |
| 564 | 0.07115                 | NaN         |
| 565 | 0.06637                 | NaN         |
| 566 | 0.07820                 | NaN         |
| 567 | 0.12400                 | NaN         |
| 568 | 0.07039                 | NaN         |

[569 rows x 33 columns]

```
[5]: ds.shape
     cols=ds.columns
     print(cols)
     ds.value_counts("diagnosis")
```

```
Index(['id', 'diagnosis', 'radius_mean', 'texture_mean', 'perimeter_mean',
       'area_mean', 'smoothness_mean', 'compactness_mean', 'concavity_mean',
       'concave points_mean', 'symmetry_mean', 'fractal_dimension_mean',
       'radius_se', 'texture_se', 'perimeter_se', 'area_se', 'smoothness_se',
       'compactness_se', 'concavity_se', 'concave points_se', 'symmetry_se',
       'fractal_dimension_se', 'radius_worst', 'texture_worst',
       'perimeter_worst', 'area_worst', 'smoothness_worst',
       'compactness_worst', 'concavity_worst', 'concave points_worst',
       'symmetry_worst', 'fractal_dimension_worst', 'Unnamed: 32'],
      dtype='object')
```

```
[5]: diagnosis
     B    357
     M    212
     dtype: int64
```

### 0.1.1 Data preprocessing

```
[8]: y = ds['diagnosis']
     ds.drop('diagnosis', axis = 1,inplace = True)
     ds.drop('Unnamed: 32', axis = 1,inplace=True)
     ds.drop('id', axis = 1,inplace=True)
     cols=ds.columns
     print(cols)
     x = ds
```

```
Index(['radius_mean', 'texture_mean', 'perimeter_mean', 'area_mean',
       'smoothness_mean', 'compactness_mean', 'concavity_mean',
       'concave points_mean', 'symmetry_mean', 'fractal_dimension_mean',
       'radius_se', 'texture_se', 'perimeter_se', 'area_se', 'smoothness_se',
       'compactness_se', 'concavity_se', 'concave points_se', 'symmetry_se',
       'fractal_dimension_se', 'radius_worst', 'texture_worst',
       'perimeter_worst', 'area_worst', 'smoothness_worst',
       'compactness_worst', 'concavity_worst', 'concave points_worst',
       'symmetry_worst', 'fractal_dimension_worst'],
      dtype='object')
```

```
[10]: ds.describe()
```

```
[10]:        radius_mean  texture_mean  perimeter_mean    area_mean  \
     count   569.000000    569.000000      569.000000   569.000000
     mean     14.127292     19.289649       91.969033   654.889104
```

```
std      3.524049      4.301036      24.298981   351.914129
min      6.981000      9.710000      43.790000   143.500000
25%     11.700000     16.170000      75.170000   420.300000
50%     13.370000     18.840000      86.240000   551.100000
75%     15.780000     21.800000     104.100000   782.700000
max     28.110000     39.280000     188.500000  2501.000000


       smoothness_mean  compactness_mean  concavity_mean  concave points_mean  \
count       569.000000        569.000000      569.000000           569.000000
mean          0.096360          0.104341        0.088799             0.048919
std           0.014064          0.052813        0.079720             0.038803
min           0.052630          0.019380        0.000000             0.000000
25%           0.086370          0.064920        0.029560             0.020310
50%           0.095870          0.092630        0.061540             0.033500
75%           0.105300          0.130400        0.130700             0.074000
max           0.163400          0.345400        0.426800             0.201200


       symmetry_mean  fractal_dimension_mean  …  radius_worst  \
count     569.000000              569.000000  …    569.000000
mean        0.181162                0.062798  …     16.269190
std         0.027414                0.007060  …      4.833242
min         0.106000                0.049960  …      7.930000
25%         0.161900                0.057700  …     13.010000
50%         0.179200                0.061540  …     14.970000
75%         0.195700                0.066120  …     18.790000
max         0.304000                0.097440  …     36.040000


       texture_worst  perimeter_worst   area_worst  smoothness_worst  \
count     569.000000       569.000000   569.000000        569.000000
mean       25.677223       107.261213   880.583128          0.132369
std         6.146258        33.602542   569.356993          0.022832
min        12.020000        50.410000   185.200000          0.071170
25%        21.080000        84.110000   515.300000          0.116600
50%        25.410000        97.660000   686.500000          0.131300
75%        29.720000       125.400000  1084.000000          0.146000
max        49.540000       251.200000  4254.000000          0.222600


       compactness_worst  concavity_worst  concave points_worst  \
count         569.000000       569.000000            569.000000
mean            0.254265         0.272188              0.114606
std             0.157336         0.208624              0.065732
min             0.027290         0.000000              0.000000
25%             0.147200         0.114500              0.064930
50%             0.211900         0.226700              0.099930
75%             0.339100         0.382900              0.161400
max             1.058000         1.252000              0.291000
```

```
        symmetry_worst  fractal_dimension_worst
count       569.000000               569.000000
mean          0.290076                 0.083946
std           0.061867                 0.018061
min           0.156500                 0.055040
25%           0.250400                 0.071460
50%           0.282200                 0.080040
75%           0.317900                 0.092080
max           0.663800                 0.207500

[8 rows x 30 columns]
```

[11]:
```python
from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test= train_test_split(x, y, test_size= 0.25,
 ↪random_state=0)
```

### 0.1.2 Feature Scaling

[13]:
```python
from sklearn.preprocessing import MinMaxScaler
st_x= MinMaxScaler()
x_train= st_x.fit_transform(x_train)
x_test=st_x.fit_transform(x_test)
print(x_train)
```

```
[[0.23044157 0.32157676 0.21940433 … 0.31484671 0.30277942 0.09858323]
 [0.20062473 0.42116183 0.19452699 … 0.06965208 0.34042973 0.06677161]
 [0.62232003 0.76929461 0.60403566 … 0.56079917 0.19850187 0.07431457]
 …
 [0.11619102 0.35726141 0.11077327 … 0.17402687 0.17524147 0.17263545]
 [0.12963226 0.35311203 0.11706171 … 0.          0.06780997 0.06919848]
 [0.21434995 0.59004149 0.21235575 … 0.33251808 0.10782574 0.21172767]]
```

### 0.1.3 Fitiing k-NN

[15]:
```python
from sklearn.neighbors import KNeighborsClassifier
classifier = KNeighborsClassifier(n_neighbors = 5)
classifier.fit(x_train, y_train)
```

[15]: KNeighborsClassifier()

[16]:
```python
y_pred= classifier.predict(x_test)
print(y_pred)
```

```
['M' 'B' 'B' 'B' 'B' 'B' 'B' 'B' 'B' 'B' 'M' 'B' 'B' 'M' 'M' 'M' 'B' 'M'
 'M' 'M' 'M' 'M' 'B' 'B' 'M' 'B' 'B' 'M' 'B' 'M' 'B' 'M' 'B' 'M' 'B' 'M'
 'B' 'M' 'B' 'M' 'M' 'B' 'M' 'B' 'B' 'M' 'B' 'B' 'B' 'M' 'M' 'M' 'M' 'B'
 'B' 'B' 'B' 'B' 'B' 'M' 'M' 'M' 'B' 'B' 'M' 'B' 'M' 'M' 'M' 'B' 'M' 'M'
```

```
'B' 'M' 'M' 'B' 'B' 'M' 'B' 'B' 'M' 'M' 'M' 'B' 'M' 'B' 'B' 'B' 'M' 'M'
'B' 'M' 'M' 'M' 'B' 'B' 'M' 'B' 'B' 'B' 'B' 'B' 'B' 'B' 'M' 'B' 'M' 'B'
'M' 'M' 'B' 'M' 'M' 'M' 'B' 'B' 'B' 'B' 'B' 'B' 'B' 'B' 'M' 'B' 'M' 'M'
'M' 'B' 'B' 'B' 'M' 'B' 'B' 'M' 'B' 'B' 'M' 'M' 'M' 'B' 'B' 'B' 'M']
```
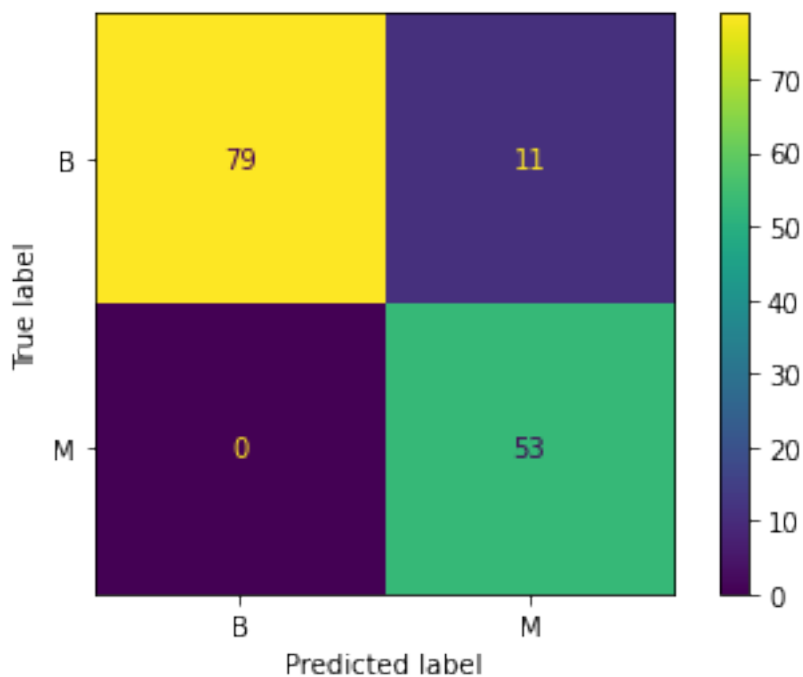
### 0.1.4 Creating confusion matrix

```python
[17]: from sklearn.metrics import confusion_matrix
      cm= confusion_matrix (y_test, y_pred,labels=classifier.classes_)
      print(cm)
```

```
[[79 11]
 [ 0 53]]
```

```python
[20]: from sklearn.metrics import ConfusionMatrixDisplay
      disp = ConfusionMatrixDisplay(confusion_matrix=cm,
      display_labels=classifier.classes_)

      disp.plot()
```

```
[20]: <sklearn.metrics._plot.confusion_matrix.ConfusionMatrixDisplay at
      0x7eff2c581f10>
```



**Accuracy**

```
[21]: training_score = classifier.score(x_train, y_train)
      test_score = classifier.score(x_test, y_test)
      print(training_score)
      print(test_score)
```

0.9765258215962441
0.9230769230769231

**New set**

```
[25]: K = []
      training = []
      test = []
      scores = {}
      for k in range(2, 22):
          clf = KNeighborsClassifier(n_neighbors = k)
          clf.fit(x_train, y_train)
          training_score = clf.score(x_train, y_train)
          test_score = clf.score(x_test, y_test)
          K.append(k)
          training.append(training_score)
          test.append(test_score)
          scores[k] = [training_score, test_score]
      for keys, values in scores.items():
          print(keys, ':', values)
```

2 : [0.9765258215962441, 0.9230769230769231]
3 : [0.9812206572769953, 0.8951048951048951]
4 : [0.9835680751173709, 0.916083916083916]
5 : [0.9765258215962441, 0.9230769230769231]
6 : [0.9788732394366197, 0.9090909090909091]
7 : [0.9788732394366197, 0.916083916083916]
8 : [0.9812206572769953, 0.951048951048951]
9 : [0.9765258215962441, 0.9230769230769231]
10 : [0.9741784037558685, 0.9370629370629371]
11 : [0.9765258215962441, 0.9230769230769231]
12 : [0.9694835680751174, 0.9300699300699301]
13 : [0.9741784037558685, 0.9300699300699301]
14 : [0.9694835680751174, 0.9370629370629371]
15 : [0.9694835680751174, 0.9370629370629371]
16 : [0.9647887323943662, 0.9440559440559441]
17 : [0.9694835680751174, 0.9370629370629371]
18 : [0.9671361502347418, 0.9440559440559441]
19 : [0.971830985915493, 0.9440559440559441]
20 : [0.9624413145539906, 0.9440559440559441]
21 : [0.9647887323943662, 0.9440559440559441]

**Visualisation**

```
[27]: import matplotlib.pyplot as plt
      plt.scatter(K, training, color ='r')
      plt.scatter(K, test, color ='g')
      plt.show()
```