# Breast cancer Wisconsin 特徵資訊解析

1. ID number

   患者編號

2. Diagnosis (M = malignant, B = benign)

   診斷結果(M:惡性，B:良性)

Ten real-valued features are computed for each cell nucleus:

以下的 10 種特徵皆是以細胞核為基準計算

1. radius (mean of distances from center to points on the perimeter)

   半徑(癌細胞核的中心至其圓周的長度的平均值)

2. texture (standard deviation of gray-scale values)

   紋理(灰階值的標準差)

3. perimeter

   癌細胞之周長

4. area

   癌細胞所占面積

5. smoothness (local variation in radius lengths)

   平滑度(癌細胞在其半徑長度內的局部變化)

6. compactness (perimeter² / area — 1.0)

   癌細胞緊緻性(半徑 ²/ 癌細胞面積 - 1)

7. concavity (severity of concave portions of the contour)

   凹陷度(輪廓凹陷的嚴重程度)

8. concave points (number of concave portions of the contour)

   凹陷點(輪廓凹陷處的數量)

9. symmetry

   對稱

10.  fractal dimension ("coastline approximation" — 1)

   分形維度(水平近似值 - 1)

# Python code

## Breast cancer Wisconsin

```python
In [1]:   1  import numpy as np
          2  import matplotlib.pyplot as plt
          3  import pandas as pd
          4  from sklearn.neighbors import KNeighborsClassifier
          5  from sklearn import metrics
          6  from sklearn.model_selection import train_test_split
          7  from sklearn.metrics import roc_curve, auc
          8  from sklearn.model_selection import KFold
```

```python
In [2]:   1  # Import cancer data
          2
          3  dataset = pd.read_csv('C:/Users/accel/cancer.csv')
          4  X = dataset.iloc[:, 2:32].values # 從 radius_mean 開始作為資料集
          5  Y = dataset.iloc[:, 0].values  # 將 diagnosis 作為 label
          6
          7  dataset.head()
```

Out[2]:

| | diagnosis | id | radius_mean | texture_mean | perimeter_mean | area_mean | smoothness_mean | compactness_mean | concavity_mean | concave points_mean | ... | ra |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | M | 842302 | 17.99 | 10.38 | 122.80 | 1001.0 | 0.11840 | 0.27760 | 0.3001 | 0.14710 | ... | |
| 1 | M | 842517 | 20.57 | 17.77 | 132.90 | 1326.0 | 0.08474 | 0.07864 | 0.0869 | 0.07017 | ... | |
| 2 | M | 84300903 | 19.69 | 21.25 | 130.00 | 1203.0 | 0.10960 | 0.15990 | 0.1974 | 0.12790 | ... | |
| 3 | M | 84348301 | 11.42 | 20.38 | 77.58 | 386.1 | 0.14250 | 0.28390 | 0.2414 | 0.10520 | ... | |
| 4 | M | 84358402 | 20.29 | 14.34 | 135.10 | 1297.0 | 0.10030 | 0.13280 | 0.1980 | 0.10430 | ... | |

5 rows × 32 columns

```python
In [3]:   1  print('cancer dataset dimensions:{}'.format(dataset.shape))
```

cancer dataset dimensions:(569, 32)

```python
In [4]:   1  # Define the score
          2
          3  def TP():
          4      x = 0
          5      for i in range(len(Pred_test)):
          6          if Pred_test[i] == Y_test[i]:
          7              if Pred_test[i] == 'B':
          8                  x += 1
          9      return x
         10
         11  def TN():
         12      x = 0
         13      for i in range(len(Pred_test)):
         14          if Pred_test[i] == Y_test[i]:
         15              if Pred_test[i] == 'M':
         16                  x += 1
         17      return x
         18
         19  def FP():
         20      x = 0
         21      for i in range(len(Pred_test)):
         22          if Pred_test[i] != Y_test[i]:
         23              if Pred_test[i] == 'B':
         24                  x += 1
         25      return x
         26
         27  def FN():
         28      x = 0
         29      for i in range(len(Pred_test)):
         30          if Pred_test[i] != Y_test[i]:
         31              if Pred_test[i] == 'M':
         32                  x += 1
         33      return x
```

```
In [5]:    1  # 10 fold cross validation
           2
           3  KF = KFold(n_splits=10)
           4
           5  XtrainL = []
           6  XtestL = []
           7  YtrainL = []
           8  YtestL = []
           9  # split features
          10  for X_train_index, X_test_index in KF.split(X): # X_train:X_test = 9:1, test依序從左至右
          11      X_train = X[X_train_index]
          12      X_test = X[X_test_index]
          13      XtrainL.append(X_train) # 儲存每次的訓練資料集
          14      XtestL.append(X_test) # 儲存每次的測試資料集
          15
          16  # split labels
          17  for Y_train_index, Y_test_index in KF.split(Y):
          18      Y_train = Y[Y_train_index]
          19      Y_test = Y[Y_test_index]
          20      YtrainL.append(Y_train)
          21      YtestL.append(Y_test)
```

```
In [6]:    1  predL = []
           2  precL = []
           3  recL =[]
           4  f1L = []
           5
           6  # 儲存10次驗證的準確率
           7  for i in range(10):
           8      KNN = KNeighborsClassifier()
           9      KNN.fit(XtrainL[i], YtrainL[i]) # train
          10
          11      Pred_test = KNN.predict(X_test) # Predict
          12
          13      Precision = TP() / (TP() + FP()) # Scoring
          14      Recall = TP() / (TP() + FN())
          15      F1 = (2 * Precision * Recall) / (Precision + Recall)
          16
          17      # store the scores
          18      predL.append(Pred_test)
          19      precL.append(round(Precision*100, 1))
          20      recL.append(round(Recall*100, 1))
          21      f1L.append(round(F1*100, 1))
          22
          23  print('Precision of test set:{}'.format(precL))
          24  print('Recall of test set:{}'.format(recL))
          25  print('F1-measure of test set:{}'.format(f1L))
```

```
Precision of test set:[97.7, 97.7, 97.7, 97.7, 97.7, 97.7, 97.7, 97.7, 97.7, 97.7]
Recall of test set:[100.0, 100.0, 100.0, 100.0, 100.0, 100.0, 100.0, 100.0, 100.0, 97.7]
F1-measure of test set:[98.9, 98.9, 98.9, 98.9, 98.9, 98.9, 98.9, 98.9, 98.9, 97.7]
```

```
In [7]:    1  # change label into binary
           2  for i in range(len(YtestL[9])):
           3      if YtestL[9][i] == 'M':
           4          YtestL[9][i] = 0
           5      elif YtestL[9][i] == 'B':
           6          YtestL[9][i] = 1
           7
           8  for i in range(len(predL[9])):
           9      if predL[9][i] == 'M':
          10          predL[9][i] = 0
          11      elif predL[9][i] == 'B':
          12          predL[9][i] = 1
          13
          14  # 為符合 roc_curve 的參數規格, 將 array 轉為 list
          15  Y_true = list(YtestL[9])
          16  Y_pred = list(predL[9])
          17
          18  print(Y_true)
          19  print(Y_pred)
```

```
[1, 0, 1, 0, 0, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1]
[1, 0, 1, 0, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 0, 1, 0, 1, 1, 1, 1, 1, 0, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1, 1,
1, 1, 1, 1, 1, 1, 0, 0, 0, 0, 0, 0, 1]
```

```
In [8]:    1  fpr, tpr, threshold = roc_curve(Y_true, Y_pred, pos_label=1)
```

```python
# roc curve & auc

roc_auc = auc(fpr,tpr) ###计算auc的值

plt.figure()
lw = 2
plt.figure(figsize=(10,10))
plt.plot(fpr, tpr, color='darkorange',
         lw=lw, label='ROC curve (area = %0.3f)' % roc_auc) ###FPR:row，TPR:column
plt.plot([0, 1], [0, 1], color='navy', lw=lw, linestyle='--')
plt.xlim([0.0, 1.0])
plt.ylim([0.0, 1.05])
plt.xlabel('False Positive Rate')
plt.ylabel('True Positive Rate')
plt.title('Receiver operating characteristic example')
plt.legend(loc="lower right")

plt.show()
```

<Figure size 432x288 with 0 Axes>