

Magic Dataset

Introduction

We are supposed to use this 'magic' dataset to generate a prediction of wether the energy particles are of type gamma (g) or hadron (h).

```
In [58]: import numpy as np
```

Loading data

We are loading the data from file using numpy genfromtext function. Which will load all the features in a hashmap like data structure to help us easily manipulate the data

```
In [59]: data_set = np.genfromtxt('magic04.data', delimiter=',', dtype=None, encoding='utf-8')
data_set

Out[59]: array([[ 28.7967,  16.0021,  2.6449,  0.3918,  0.1982,   27.7004,   22.011 , -8.2027,  40.092 ,   81.882
 8, 'g'],
 [ 31.6036,  11.7235,  2.5185,  0.5303,  0.3773,   26.2722,   23.8238, -9.9574,   6.3609,  205.261
 , 'g'],
 [162.052 ,  136.031 ,   4.0612,  0.0374,  0.0187,  116.741 ,  -64.858 , -45.216 ,   76.96 ,  256.788
 , 'g'],
 ...,
 [ 75.4455,  47.5305,  3.4483,  0.1417,  0.0549,   -9.3561,   41.0562, -9.4662,  30.2987,  256.516
 6, 'h'],
 [120.5135,   76.9018,  3.9939,  0.0944,  0.0683,    5.8043,  -93.5224, -63.8389,  84.6874,  408.316
 6, 'h'],
 [187.1814,   53.0014,  3.2093,  0.2876,  0.1539, -167.3125, -168.4558,  31.4755,  52.731 ,  272.317
 4, 'h']],
      dtype=[('f0', '<f8'), ('f1', '<f8'), ('f2', '<f8'), ('f3', '<f8'), ('f4', '<f8'), ('f5', '<f8'), ('f6', '<f8'), ('f7', '<f8'), ('f8', '<f8'), ('f9', '<f8'), ('f10', '<U1')])
```

Separating Gammas from Hadrons

In order to make both datasets equal in size we will separate Gs from Hs. Knowing that gamma data is in the first 12332 rows we separate the gamma in a variable called g_class

Gamma data

```
In [60]: g_class=data_set[:12332]
g_class

Out[60]: array([[ 28.7967,  16.0021,  2.6449,  0.3918,  0.1982,   27.7004,   22.011 , -8.2027,  40.092 ,   81.8828,
 8, 'g'],
 [ 31.6036,  11.7235,  2.5185,  0.5303,  0.3773,   26.2722,   23.8238, -9.9574,   6.3609,  205.261 ,
 8, 'g'],
 [162.052 ,  136.031 ,   4.0612,  0.0374,  0.0187,  116.741 ,  -64.858 , -45.216 ,   76.96 ,  256.788 ,
 8, 'g'],
 ...,
 [ 22.0913,  10.8949,  2.2945,  0.5381,  0.2919,   15.2776,   18.2296,   7.3975,  21.068 ,  123.281 ,
 8, 'g'],
 [ 56.2216,  18.7019,  2.9297,  0.2516,  0.1393,   96.5758, -41.2969,  11.3764,   5.911 ,  197.209 ,
 8, 'g'],
 [ 31.5125,  19.2867,  2.9578,  0.2975,  0.1515,   38.1833,   21.6729, -12.0726,  17.5809,  171.227 ,
 8, 'g']],
      dtype=[('f0', '<f8'), ('f1', '<f8'), ('f2', '<f8'), ('f3', '<f8'), ('f4', '<f8'), ('f5', '<f8'), ('f6', '<f8'), ('f7', '<f8'), ('f8', '<f8'), ('f9', '<f8'), ('f10', '<U1')])
```

Hadron data

```
In [61]: h_class=data_set[12332:]
h_class

Out[61]: array([[ 93.7035,  37.9432,  3.1454,  0.168 ,  0.1011,   53.2566,   89.0566,  11.8175,  14.1224,  231.9028,
 8, 'h'],
 [102.0005,  22.0017,  3.3161,  0.1064,  0.0724,  -54.0862,   43.0553, -15.0647,  88.4636,  274.9392,
 8, 'h'],
 [100.2775,  21.8784,  3.11 ,  0.312 ,  0.1446,  -48.1834,   57.6547,   -9.6341,  20.7848,  346.433 ,
 8, 'h'],
 ...,
 [ 75.4455,  47.5305,  3.4483,  0.1417,  0.0549,   -9.3561,   41.0562, -9.4662,  30.2987,  256.5166,
 8, 'h'],
 [120.5135,  76.9018,  3.9939,  0.0944,  0.0683,    5.8043,  -93.5224, -63.8389,  84.6874,  408.3166,
 8, 'h'],
 [187.1814,   53.0014,  3.2093,  0.2876,  0.1539, -167.3125, -168.4558,  31.4755,  52.731 ,  272.3174,
 8, 'h']],
      dtype=[('f0', '<f8'), ('f1', '<f8'), ('f2', '<f8'), ('f3', '<f8'), ('f4', '<f8'), ('f5', '<f8'), ('f6', '<f8'), ('f7', '<f8'), ('f8', '<f8'), ('f9', '<f8'), ('f10', '<U1')])
```

Making Gs the same size as Hs

We are taking a random 6688 rows from the gamma array to make both datasets equal

```
In [62]: g_class=np.random.choice(g_class,size=6688,replace=False)
print(g_class)
print(g_class.shape)

[(117.16 , 21.5912, 3.0204, 0.3015, 0.1665, -117.488 , -73.3775, -14.3856, 0.181 , 324.714, 'g')
 ( 20.9316, 15.2379, 2.4857, 0.4575, 0.2402,  25.2371,  14.1903,   5.9376, 36.675 , 137.007, 'g')
 ( 48.5393, 19.3707, 2.4401, 0.3775, 0.1942,   35.1581,  33.3854, -14.009 ,   9.9947, 163.199, 'g')
 ...
 ( 72.8137, 34.9071, 3.2667, 0.2922, 0.1631,   -0.4914,  18.4641,  20.3378, 12.275 , 338.84 , 'g')
 ( 74.7241, 25.517 , 3.3375, 0.1687, 0.0857,   52.4757,  57.9373, -9.0161,  4.5284, 261.352, 'g')
 ( 40.0652, 21.3799, 2.9811, 0.2851, 0.1655,   18.9428,  26.5156,  16.5628, 18.0501, 213.502, 'g')]
(6688,)
```

Constructing a new data array

We will concatenate the new shortened g_class with the h_class to have our new dataset with equal parameters for both predictions.

```
In [63]: data=np.concatenate((g_class,h_class),axis=0)
print(data)
print(data.shape)

[(117.16 , 21.5912, 3.0204, 0.3015, 0.1665, -117.488 , -73.3775, -14.3856, 0.181 , 324.714 , 'g')
 ( 20.9316, 15.2379, 2.4857, 0.4575, 0.2402,  25.2371,  14.1903,   5.9376, 36.675 , 137.007 , 'g')
 ( 48.5393, 19.3707, 2.4401, 0.3775, 0.1942,   35.1581,  33.3854, -14.009 ,   9.9947, 163.199 , 'g')
 ...
 ( 75.4455, 47.5305, 3.4483, 0.1417, 0.0549,   -9.3561,   41.0562, -9.4662,  30.2987,  256.5166, 'h')
 (120.5135, 76.9018, 3.9939, 0.0944, 0.0683,    5.8043,  -93.5224, -63.8389,  84.6874,  408.3166, 'h')
 (187.1814, 53.0014, 3.2093, 0.2876, 0.1539, -167.3125, -168.4558,  31.4755,  52.731 ,  272.3174, 'h')]
(13376,)
```

Randomizing data by shuffle and splitting

We are shuffling all our data using numpy

Spliting data to Training, Testing, Validation data sets

```
In [64]: rng = np.random.default_rng()
rng.shuffle(data)
train,test_validate=np.array_split(data,[int(0.70 * len(data))])
test,validation=np.array_split(test_validate,[int(0.50 * len(test_validate))])

In [65]: print(f"train: {train.shape[0]}\n"
          f"test: {test.shape[0]}\n"
          f"validation: {validation.shape[0]}")

train: 9363
test: 2006
validation: 2007
```

Separating features from the prediction

We are taking the first ten features and putting them in array x_data and the last feature (the prediction) in array y_data

```
In [66]: def return_x_y_arrays(data_set_to_be_sliced):
        buf=data_set_to_be_sliced.tolist()
        x_data= []
        y_data=[]
        for d in buf:
            x=d[:10]
            x_data.append(x)
            y=np.array([d[10]])
            y_data.append(y)
        x_data=np.array(x_data)
        y_data=np.array(y_data).ravel()
        return x_data,y_data

In [67]: x_train,y_train=return_x_y_arrays(train)
x_validation,y_validation=return_x_y_arrays(validation)
x_test,y_test=return_x_y_arrays(test)
```

Fitting and training the model

```
In [68]: from sklearn.neighbors import KNeighborsClassifier as knn
model=knn(n_neighbors=1)
model.fit(x_train, y_train)
```

```
Out[68]: KNeighborsClassifier
KNeighborsClassifier(n_neighbors=1)
```

Scoring the model using test and validation datasets

We will use various K values to score the model

k=1 k=3 k=10 k=23 -> Best k=600

```
In [69]: print(model.score(x_test,y_test))
print(model.score(x_validation,y_validation))

0.7318045862412762
0.7409068261086198
```

```
In [70]: model.n_neighbors=3
print(model.score(x_test,y_test))
print(model.score(x_validation,y_validation))

0.7582253240279162
0.7513702042850025
```

```
In [71]: model.n_neighbors=10
print(model.score(x_test,y_test))
print(model.score(x_validation,y_validation))

0.7681954137587238
0.7678126557050324
```

```
In [72]: model.n_neighbors=23
print(model.score(x_test,y_test))
print(model.score(x_validation,y_validation))

0.765702891326022
0.7688091679123069
```

```
In [73]: model.n_neighbors=600
print(model.score(x_test,y_test))
print(model.score(x_validation,y_validation))

0.7288135593220338
0.7159940209267563
```

```
In [74]: from sklearn.neighbors import KNeighborsClassifier
from sklearn.model_selection import GridSearchCV
KN = KNeighborsClassifier()
k_range = list(range(1, 26, 2)) # 1, 3, 5, ..., 25
param_grid = dict(n_neighbors=k_range)
print(param_grid)
grid = GridSearchCV(KN, param_grid, cv=10, scoring='accuracy', return_train_score=True)

# Training data
g = grid.fit(x_train, y_train)
print(grid.best_params_)
accuracy=grid.best_score_ * 100
print("Accuracy for our training dataset with tuning is : {:.2f}%".format(accuracy) )

{'n_neighbors': [1, 3, 5, 7, 9, 11, 13, 15, 17, 19, 21, 23, 25]}
{'n_neighbors': 17}
Accuracy for our training dataset with tuning is : 77.28%
```

```
In [75]: best_k = grid.best_params_['n_neighbors']

# Calculation test accuracy score
knn = KNeighborsClassifier(n_neighbors=best_k)
knn.fit(x_train,y_train)
y_test_pred = knn.predict(x_test)
test_accuracy = knn.score(x_test, y_test)
print(test_accuracy)

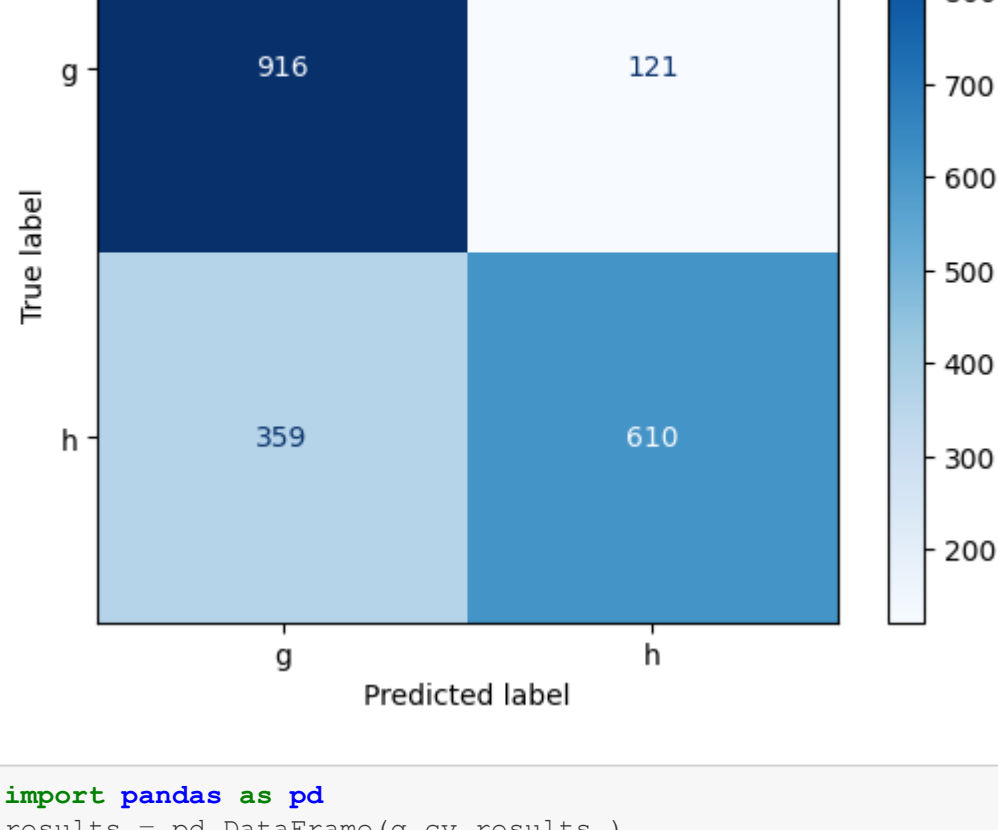
# Calculation validation accuracy score
knn = KNeighborsClassifier(n_neighbors=best_k)
knn.fit(x_train,y_train)
y_validation_pred = knn.predict(x_validation)
test_accuracy = knn.score(x_validation, y_validation)
print(test_accuracy)

0.7607178464606181
0.7802690582959642
```

Using confusion matrix to display the results in human-readable graphic

```
In [76]: from sklearn.metrics import confusion_matrix,classification_report
from sklearn.metrics import ConfusionMatrixDisplay
import matplotlib.pyplot as plt
matix = confusion_matrix(y_test,y_test_pred)
print(classification_report(y_test,y_test_pred))
disp = ConfusionMatrixDisplay(confusion_matrix=matix,display_labels=['g','h'])
disp = disp.plot(cmap=plt.cm.Blues)
```

		precision	recall	f1-score	support
	g	0.72	0.88	0.79	1037
	h	0.83	0.63	0.72	969
accuracy				0.76	2006
macro avg		0.77	0.76	0.76	2006
weighted avg		0.77	0.76	0.76	2006



```
In [77]: import pandas as pd
results = pd.DataFrame(g.cv_results_)
needed_results = results[['param_n_neighbors', 'mean_train_score', 'mean_test_score']]
needed_results
```

```
Out[77]:
```

	param_n_neighbors	mean_train_score	mean_test_score
0	1	1.000000	0.740256
1	3	0.870163	0.759692
2	5	0.837777	0.764817
3	7	0.822956	0.767594
4	9	0.813201	0.768556
5	11	0.806983	0.771227
6	13	0.802295	0.769518
7	15	0.798308	0.768877
8	17	0.795899	0.772829
9	19	0.792374	0.766313
10	21	0.789645	0.765244
11	23	0.786951	0.765993
12	25	0.786180	0.768022