

Optics_Matrix-Prototype.py

1. Under the main function, we have defined an object "optics" under a class of the same name.

```
if __name__ == '__main__':  
    optics=OPTICS()  
  
    optics.GenNumpyArray("../Cfoil_elastic_p2_slim.root")  
    optics.DefineSectors()  
    optics.DrawHistAllSectors()  
    optics.DrawScatterPlot(optics.secl.gem1_x, optics.secl.gem1_y)  
    optics.SelectOneHole(optics.secl)  
  
    hole_id="12"  
    filename="output/SieveHole_"+hole_id+".csv"  
    optics.GenCSV(hole_id, filename)
```

Within a class we have methods and constructors. Methods are nothing but functions. Constructor is a special kind of method.

The first thing that we define when we create a class is a constructor.

The keyword "self" is used to assign the arguments passed to the method to the attributes. It also communicates among different functions defined within the class.

```
class OPTICS:  
    def __init__(self):  
        self.orig = pd.DataFrame() # data before cut  
        self.secl = pd.DataFrame() # sector1 data before cut  
        self.sec2 = pd.DataFrame() # sector2 data before cut  
        self.sec3 = pd.DataFrame() # sector3 data before cut  
        self.sec4 = pd.DataFrame() # sector4 data before cut  
        self.sec5 = pd.DataFrame() # sector5 data before cut  
        self.sec6 = pd.DataFrame() # sector6 data before cut  
        self.sec7 = pd.DataFrame() # sector7 data before cut  
  
        self.selected = pd.DataFrame() # data of selected holes
```

we have defined 8 empty panda dataframes inside the constructor.

This is the first thing that will be executed when creating an instance of the class.

2. After initializing panda dataframes, we are filling the panda dataframe **self.orig** from data in the slim root files.

```
import uproot
import awkward
import pandas as pd
import math
import matplotlib.pyplot as plt

df=pd.DataFrame()

file = uproot.open(".././../Cfoil elastic.pl.root")
print(file.classnames()) #This statement gives the objects within the file and what class these objects belong to
T = file["neut"]
print(T.keys())          # We want to know what are the branches inside the tree.

branches = T.arrays()
print(branches['main_r']) # prints out the data in the branch main_r for all the events
print(branches['main_r'][0]) #prints out the value of main_r for event 0

geo = T.arrays(["gem1_x", "gem1_y", "gem1_r", "gem1_ph", "gem1_px", "gem1_py", "gem1_pz", "tg_th", "tg_ph", "tg_vz", "tg_p", "rate"], library="pd")
geo.loc[geo['gem1_r'] > 300]
df=geo
print(df.columns)
```

First task would be to open the root file. For this purpose we are using the uproot library.

- The next task is to read the data from the branches. For this purpose we will use the **T.arrays()** method.
 - Within the panda dataframe we only need to store information from a certain number of branches.
 - The "loc" function filters the events according to the condition that " $\text{gem1-}\tau > 300$ ".
3. The next step is to fill the 7 empty panda dataframes one corresponding to each sector. For this purpose we need to have a definition of sectors **That's what the Define Sectors() function does.**

```
def DefineSectors(self):
    rot_angle=0          # put the GEM rotation angle here

    angle_lo=[]
    angle_up=[]

    for i in range(7):
        angle_lo.append(math.pi/14 + i*2*math.pi/7+rot_angle)
        angle_up.append(3*math.pi/14 + i*2*math.pi/7+rot_angle)
        if angle_lo[i]>math.pi:
            angle_lo[i]=angle_lo[i]-2*math.pi
        if angle_up[i]>math.pi:
            angle_up[i]=angle_up[i]-2*math.pi

    geom=self.orig

    self.sect1=geo.loc[(geo['gem1_ph']<angle_lo[0]) & (geo['gem1_ph']>angle_lo[0])]
    self.sect2=geo.loc[(geo['gem1_ph']>angle_lo[0]) & (geo['gem1_ph']>angle_lo[1])]
    self.sect3=geo.loc[(geo['gem1_ph']<angle_lo[2]) & (geo['gem1_ph']>angle_lo[2])]
    self.sect4=geo.loc[(geo['gem1_ph']>angle_lo[2]) & (geo['gem1_ph']>math.pi)]
    self.sect5=geo.loc[(geo['gem1_ph']>angle_lo[0]) & (geo['gem1_ph']>angle_lo[4])]
    self.sect6=geo.loc[(geo['gem1_ph']<angle_lo[5]) & (geo['gem1_ph']>angle_lo[5])]
    self.sect7=geo.loc[(geo['gem1_ph']>angle_lo[5]) & (geo['gem1_ph']>angle_lo[6])]
```



```
def SelectOneHole(self, df):
```

```
fig, ax = plt.subplots(figsize=(10,7))
```

```
pts=ax.scatter(df.geml_x,df.geml_y)
```

→ Producing a scatter plot

```
y_max=df.geml_y.max()
y_min=df.geml_y.min()
dy = (y_max-y_min)*0.1
```

```
x_max=df.geml_x.max()
x_min=df.geml_x.min()
dx = (x_max-x_min)*0.1
```

```
ax.set_ylim(y_min-dy, y_max+dy)
ax.set_xlim(x_min-dx, x_max+dx)
```

→ Setting the min and max values of the plot.

```
selector = SelectFromCollection(ax, pts)
```

```
print("Select points in the figure by enclosing them within a polygon.")
print("Press the 'esc' key to start a new polygon.")
print("Try holding the 'shift' key to move all of the vertices.")
print("Try holding the 'ctrl' key to move a single vertex.")
```

```
plt.show()
```

```
selector.disconnect()
```

```
self.selected=df.loc[df.index[selector.ind]]
```

This is the function we are using to select one particular hole

we are using

from polygone-selector-demo

import SelectFromCollection

Again using the loc function to fill the self.selected dataframe.

6. The next step is to generate a csv file from the selected hole dataframe.

```
def GenCSV(self, hole_id, filename):
```

```
df=self.selected
```

```
df["geml_rp"]=(df.geml_x*df.geml_px+df.geml_y*df.geml_py)/(df.geml_r*df.geml_pz) # gem 1 r'
```

```
df["geml_php"]=(df.geml_y*df.geml_px+df.geml_x*df.geml_py)/(df.geml_r*df.geml_pz) # gem 1 phi'
```

```
header=["tg_th","tg_ph","tg_vz","tg_p","geml_r","geml_rp","geml_ph","geml_php"]
```

```
df.to_csv(filename,columns=header)
```

There are a lot of columns in the panda dataframe. We do not need to load all of those in the csv file.

For this reason we are passing the argument "columns = header" inside df.to_csv()

7. Now before we can perform a polynomial regression we need to understand how python implements Machine Learning methods in polynomial regression.

- we first divide the dataset into a training set and a testing set. Usually it's an 80:20 ratio but we can change the relative sizes of the two sets.
- if we assume that z depends upon two variables x & y . Machine Learning methods also studies the dependence of z on the correlations of x and y .

$z = f(x,y)$ and we want upto second degree polynomial Regression.

$x^0, x^1, x^2, y^0, y^1, y^2, xy$] we need information on all these.

We use the "sklearn" module to perform the polynomial regression.

```

def PolynomialRegression(x_train, y_train, degree):
    # X is the GM variables in numpy array, y is the target variable
    X_train, y_train, X_test, y_test = train_test_split(X_train, y_train, test_size=0.3, random_state=42)
    # split the data into training set and test set

    poly = PolynomialFeatures(degree)
    X_train_newpoly = poly.fit_transform(X_train)
    X_test_newpoly = poly.fit_transform(X_test)

    # print(X_train_newpoly)
    regression = LinearRegression()
    model = regression.fit(X_train_newpoly, y_train)
    y_pred = regression.predict(X_test_newpoly)

    params = model.coef_
    intercept = model.intercept_

    # ax = plt.axes(projection='3d')
    # ax.scatter(X_train[:,0], X_train[:,1], y_train, c=y_train, cmap='Greens')
    # ax.scatter(X_test[:,0], X_test[:,1], y_pred, c=y_pred, cmap='Reds')
    # plt.show()

    # ax = plt.axes()
    # ax.scatter(X_train[:,0], X_train[:,1], y_train, c=y_train, cmap='Greens')
    # ax.scatter(X_test[:,0], X_test[:,1], y_pred, c=y_pred, cmap='Reds')
    # ax.scatter(X_test[:,0], X_test[:,1], y_pred, c=y_pred, cmap='Reds')
    # plt.show()

    score = model.score(X_test_newpoly, y_test)
    print("score: ", score)

```

The first module that we need is

from sklearn.model_selection
import train_test_split.

train_test_split() is used to split the dataset into a training and a testing set.

Another feature of train_test_split function is that it also randomizes the data set.

Then we need to take into account all the terms in the fit function.

$$\begin{aligned}
 \theta_{tg} &= f(\gamma, \gamma') \\
 &= \gamma^0 \gamma' \gamma^2 \\
 &\quad \gamma'^0 \gamma'^1 \gamma'^2 \\
 &\quad \gamma \gamma'
 \end{aligned}$$

Within the csv file we already have columns for γ and γ' .
Now we need to compute all the other columns.

from sklearn.preprocessing import PolynomialFeatures

We define the model using the training set and we test the efficiency of our model using the testing set.

`np.arange(3)` specifies that you want to slice out a 1D vector of length 3 from a 2D array.
`np.arange(3, 4)` specifies that you want to slice out a 2D sub-array of shape (3, 4) from the 2D array. The square brackets make all the difference here.
`v = np.arange(3).reshape(1, 3)`
`v[0, 0]`
`array([0, 2, 4])`
`v[0, 0]`
`array([0])`
`v[0]`
`array([0, 2, 4])`
`v[0]`
`array([0, 2, 4])`

The fundamental difference is the extra dimension in the latter command (as you've noted). This is intended behaviour, as implemented in `numpy.ndarray.reshape()` and codified in the NumPy documentation.

You can also specify `np.newaxis` to the same effect - a column sub-slice.

```

v[0, 0]
array([0])
v[0, 0]
array([0])

```

For more information, look at the notes on [Advanced Indexing](#) in the NumPy docs.

→ Important Information
When we will show the accuracy of our model using the testing set.