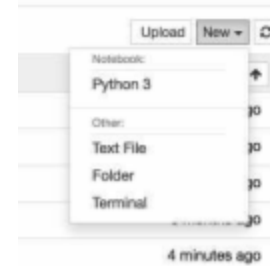


Importing a CMS open data file and creating histograms to analyze data

This exercise uses Jupyter Notebook which can be downloaded for free by installing Anaconda, available at <https://www.continuum.io/downloads> to analyze CMS data sets from the CERN Open Data Portal, <http://opendata.cern.ch/>.

Select a data file to analyze from <http://opendata.cern.ch/record/545>. For the example, download Dimuon_DoubleMu.csv. Note, there are other files available, but the file must be comma separated values (.csv). Download the selected file, noting its exact location and file name.

Open the Anaconda - Navigator. Launch Jupyter notebook.
Select “New” from then “Python 3”.



Boxes that give instruction or explanation but not code, can be entered by selecting “Markdown” from the dropdown menu.



If you are returning to a page after working on it previously, it is recommended that you select "Kernel" followed by "Restart & Clear Output"

To enter code, select "Code" from the dropdown menu.
Code may be copied from other sources, pasted, and edited.

First, you will need to import the packages pandas and matplotlib.pyplot.

```
In [5]: import pandas
import matplotlib.pyplot as plt
%matplotlib inline
```

After entering a code box, press ctrl & enter at the same time to run the code. An asterisk will appear in the In[] while the command is being processed.

Wait until a number appears in that location before proceeding.

Error messages will appear in pink and will indicate information regarding the error.

To use the file you have downloaded, it must be saved into a variable. Type the file location and name exactly as it appears on your computer.

In this example, the file is data/Dimuon_DoubleMu.cvs.

Save the data into the variable dataset and check the contents of the first 5 rows.

```
In [6]: dataset = pandas.read_csv('data/Dimuon_DoubleMu.csv')
        dataset.head()
```

If the command is written properly, a data table will be generated.

For the given example, the invariant mass is of particular interest. To create a histogram of the invariant mass column, save the dataset based on the column heading exactly as it appears in the table (M) as a variable invariant_mass.

Plot a histogram, stating which variable to plot, setting the number of bins and range.

In the example below, we use 50 bins and plot from 0 to 200 GeV.

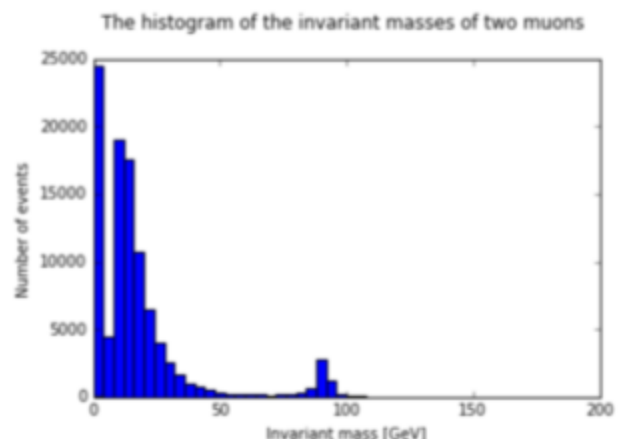
```
In [7]: invariant_mass = dataset['M']
        plt.hist(invariant_mass, bins=50, range=(0,200))
        plt.show()
```

This plot was primarily to check our data and plotting commands. The number of bins and range can be varied in order to more clearly analyze the data.

Below are commands to replot the graph with a title and axis labels.

```
In [6]: plt.xlabel('Invariant mass [GeV]')
        plt.ylabel('Number of events')
        plt.title('The histogram of the invariant masses of two muons \n')
        plt.hist(invariant_mass, bins=50, range=(0,200))
        plt.show()
```

A sample histogram is shown on the right. The graph demonstrates background events below approximately 50 GeV and a peak at approximately 90 GeV.



Mathematical operations, such as addition or subtraction, may be performed on the data by defining new variables which allows the data to be further sorted.

In the example below, the original data is divided into two new datasets based on the energy of the collision. Each dataset is given a name and organized in this case by high energy (>150 GeV), and low energy (<150 GeV).

```
In [7]: newsethighE = dataset[dataset.E1+dataset.E2>150]
        newsetlowE = dataset[dataset.E1+dataset.E2<150]
```

The new dataset can be plotted separately as was done previously above or on one plot.

The two histograms can be overlaid by adjusting the transparency using the alpha command.

Labels for each data set are included in the legend located in the upper right corner.

We can also change the range to focus on the event of interest.

```
In [43]: plt.xlabel('Invariant mass [GeV]')
        plt.ylabel('Number of events')
        plt.title('The invariant masses of two muons comparing high and low energy\n')
        plt.hist(newsetlowE ['M'], bins=50, range=(80,100),alpha=0.5, label='Low E')
        plt.hist(newsethighE ['M'], bins=50, range=(80,100),alpha=0.5, label='High E')
        plt.legend (loc='upper right')
        plt.show()
```

The final output is shown below.

