# HOME ASSIGNMENT 4
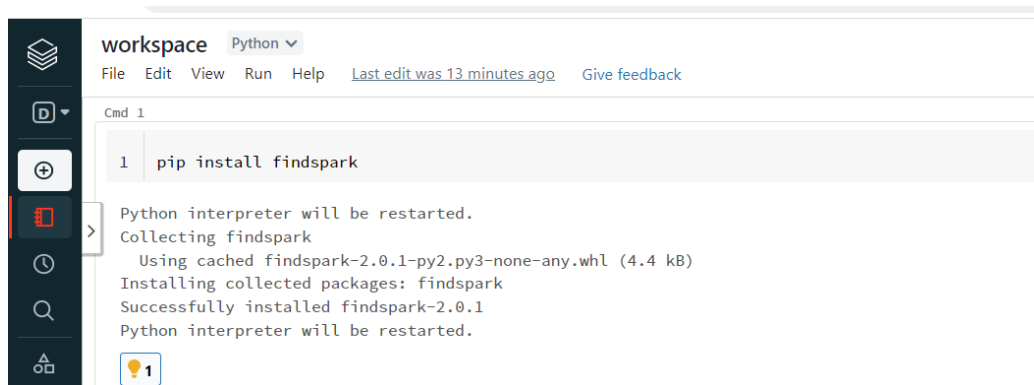
**Roll No**: 208W1A1299

**Name**: MOHAMMAD RIZWANULLAH

**Aim**: To do analysis using Apache spark in data bricks.

**Execution**:

1)Open databricks and create a new cluster with latest configuration version.

2) Now create a Notebook to work with spark.

3) Import csv file or any file for using it and do analysis using pyspark.

4) Now write code to import dataset and create some RDD to read input in spark.

5) Run the cell.

## Below are the screenshots of execution:

Twitter followers count and protected data analysis.

```
1  df.show()
```

▸ (3) Spark Jobs

```
+-------+-----+
|    _1|   _2|
+-------+-----+
|  funny|10202|
| comedy| 8911|
| [none]| 5513|
|  music| 3807|
| how to| 3783|
|   2017| 3750|
| makeup| 3635|
|trailer| 3578|
|   news| 3430|
|  humor| 3376|
+-------+-----+
```

```python
1  import matplotlib.pyplot as plt
2
3  my_data = [10202, 8911, 3635,3807,3578,3430,3376]
4  my_labels = 'funny', 'comedy', 'makeup','music','trailer','news','humor'
5  plt.pie(my_data, labels=my_labels, autopct='%1.1f%%')
6  plt.title('My Tasks')
7  plt.axis('equal')
8  plt.show()
```



My Tasks