■

**Name: Tushar Panchal**

**En.No: 21162101014**

**Sub: EADC (Enterprise Application Development for Cloud)**

**Branch: CBA**

**Batch:61**

----------------------------PRACTICAL 18----------------------------
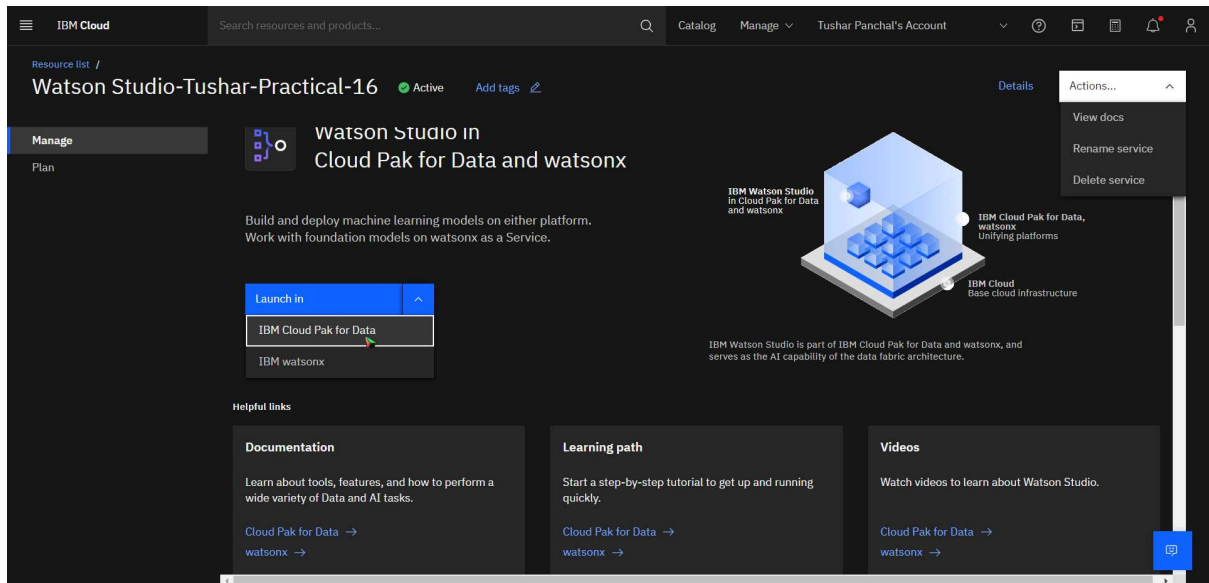
**Assume you are working in a company where you need to extract useful insights from the data collected by organization. Demonstrate how to analyze large datasets with Python data science packages. We'll provide an example use case of analysing hourly air quality data provided by the EPA.**
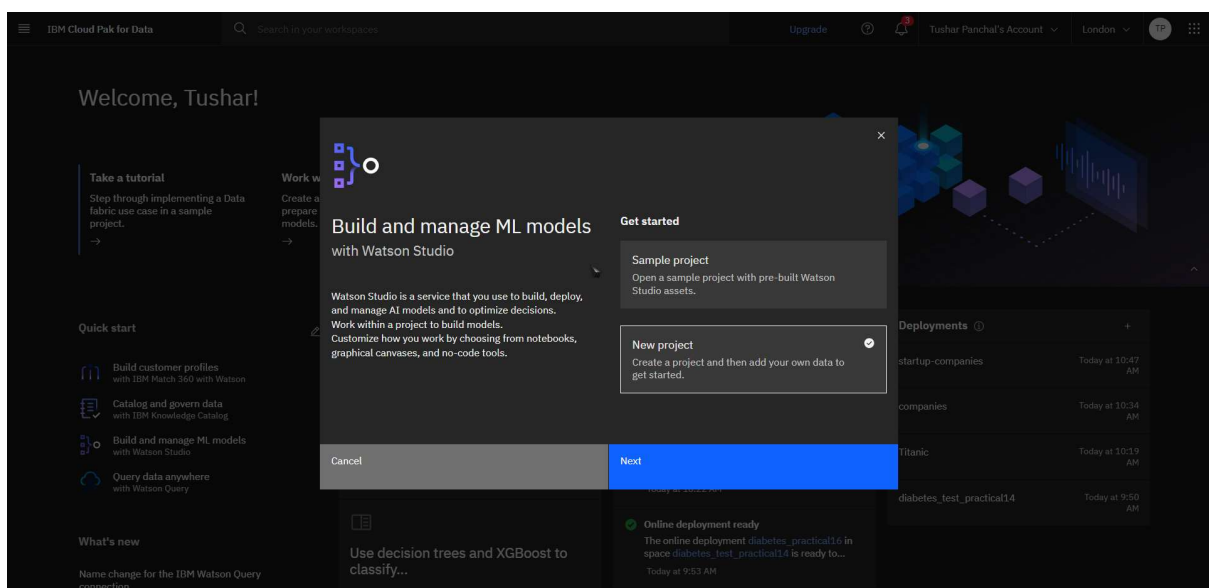
Perform the following Tasks:
1. Create a Juypter notebook in Watson Studio.
2. Extract patterns from datasets using pandas.
3. Visualize data trends via matplotlib graphs.

## » Task 1 : Create a Juypter notebook in Watson Studio.

Navigate watson studio that we created in practical - 16 than click on launch in IBM cloud Pak for Data



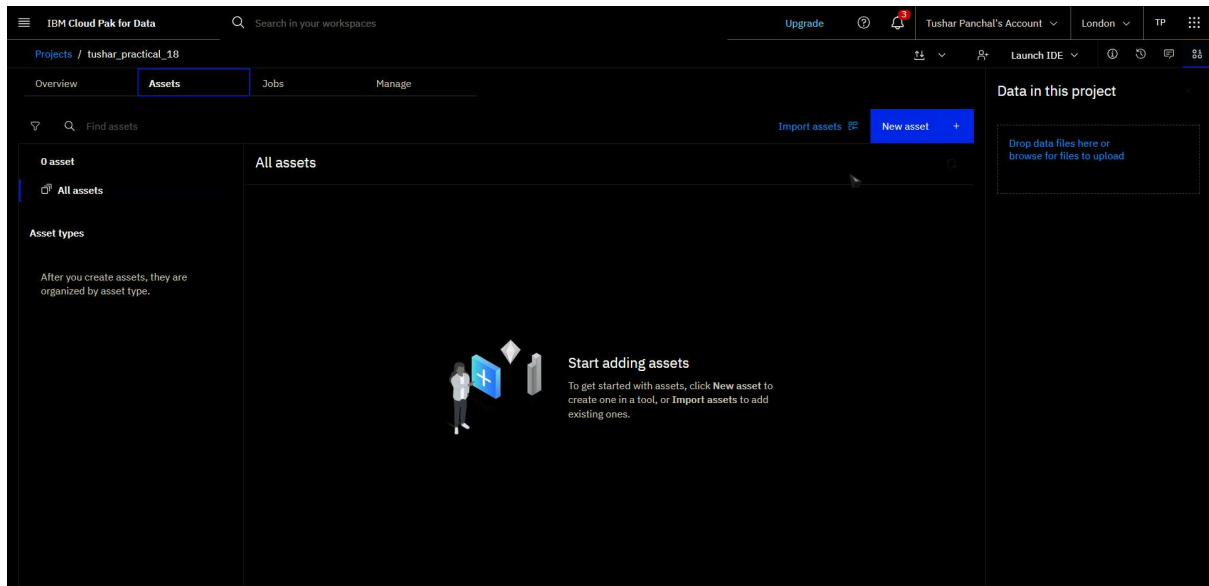Now click on new project and click next

## Give name of the new project and click on create



## Its open this interface

## Select assets then click on new asset



## Now search jupyter than click on "jupyter notebook editor"

Now add name and select "URL" give below URI Then click on create

https://github.com/IBM/smart-city-analytics/blob/master/air_quality_notebook.ipynb



After create its loading the interface

As we can see in below screenshot all query



Now run one by one all codes

## >> Task 2 : Extract patterns from datasets using pandas.

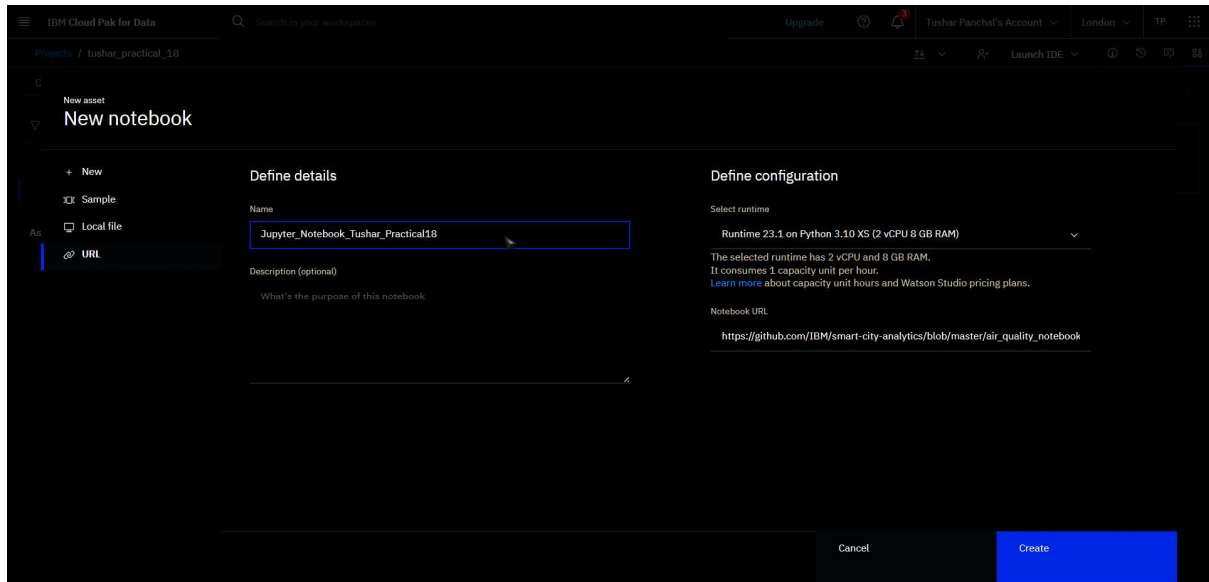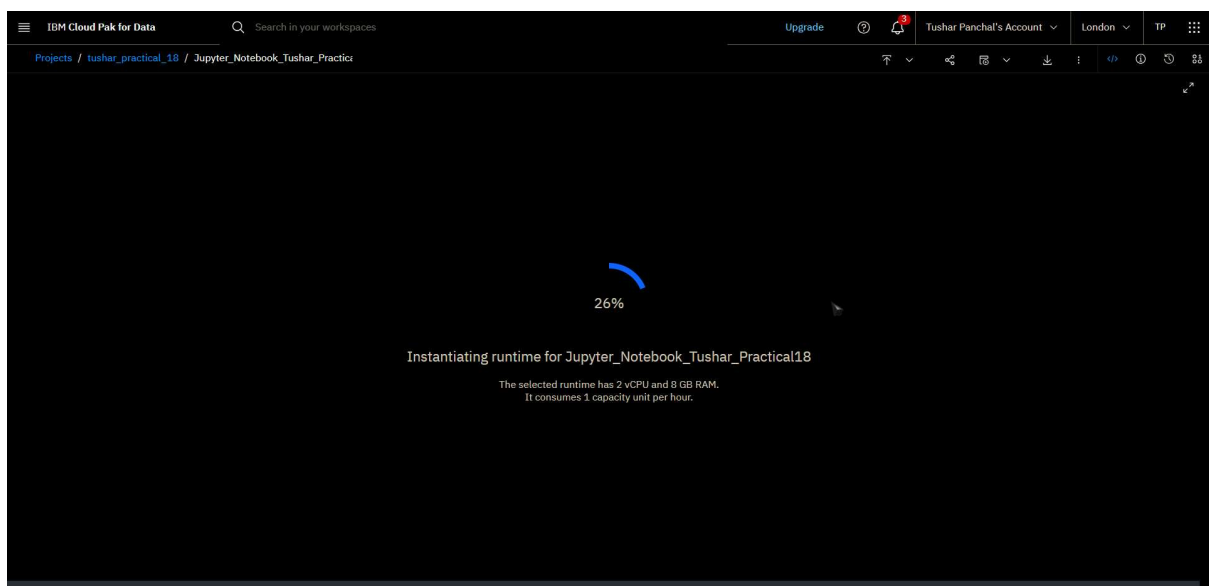**IBM Cloud Pak for Data**  Q Search in your workspaces   Upgrade  ? ⊘³  Tushar Panchal's Account ⌄   London ⌄  TP  ⠿

Projects / tushar_practical_18 / Tushar_Practical18   ↑ ⌄   ⇄   ☐ ⌄   ⬇   ⋮   </> ⓘ ↺ ⠿

File  Edit  View  Insert  Cell  Kernel  Help                          Not Trusted | Python 3.10 ◯ ⤢

◉ ⊕ ✛ ⎘ ⎗ ↑ ↓ ▶ Run ◉ ⟳ ⏭ Format Code ⌄ ▦        Memory:359.1 MB / 8 GB

```
In [76]: !ls
         LICENSE                  hourly_42602_2017.csv
         README.md                hourly_42602_2017.zip
         air_quality_notebook.ipynb  images

In [77]: # Load dataset
         aq_data = pd.read_csv('./hourly_42602_2017.csv')
         # aq_data = pd.read_csv('/Users/kkbankol@us.ibm.com/Downloads/hourly_42602_2017.csv')

         /usr/local/lib/python3.7/site-packages/IPython/core/interactiveshell.py:2785: DtypeWarning: Columns (17) have mixed types. Specify dtype option on import or set low_memory=False.
           interactivity=interactivity, compiler=compiler, result=result)

In [78]: # View first 5 rows
         aq_data.head()
```

Out[78]:

| | State Code | County Code | Site Num | Parameter Code | POC | Latitude | Longitude | Datum | Parameter Name | Date Local | ... | Units of Measure | MDL | Uncertainty | Qualifier | Method Type | Method Code | Method Name | State Name | County Name | Date of Last Change |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 73 | 23 | 42602 | 1 | 33.553056 | -86.815 | WGS84 | Nitrogen dioxide (NO2) | 2017-01-01 | ... | Parts per billion | 0.1 | NaN | NaN | FEM | 200 | Teledyne-API Model 200EUP or T200UP - Photolyt... | Alabama | Jefferson | 2017-04-19 |
| 1 | 1 | 73 | 23 | 42602 | 1 | 33.553056 | -86.815 | WGS84 | Nitrogen dioxide (NO2) | 2017-01-01 | ... | Parts per billion | 0.1 | NaN | NaN | FEM | 200 | Teledyne-API Model 200EUP or T200UP - Photolyt... | Alabama | Jefferson | 2017-04-19 |
| 2 | 1 | 73 | 23 | 42602 | 1 | 33.553056 | -86.815 | WGS84 | Nitrogen dioxide (NO2) | 2017-01-01 | ... | Parts per billion | 0.1 | NaN | NaN | FEM | 200 | Teledyne-API Model 200EUP or T200UP - Photolyt... | Alabama | Jefferson | 2017-04-19 |
| 3 | 1 | 73 | 23 | 42602 | 1 | 33.553056 | -86.815 | WGS84 | Nitrogen dioxide (NO2) | 2017-01-01 | ... | Parts per billion | 0.1 | NaN | NaN | FEM | 200 | Teledyne-API Model 200EUP or T200UP - Photolyt... | Alabama | Jefferson | 2017-04-19 |
| 4 | 1 | 73 | 23 | 42602 | 1 | 33.553056 | -86.815 | WGS84 | Nitrogen dioxide (NO2) | 2017-01-01 | ... | Parts per billion | 0.1 | NaN | NaN | FEM | 200 | Teledyne-API Model 200EUP or T200UP - Photolyt... | Alabama | Jefferson | 2017-04-19 |

5 rows × 24 columns

```
In [79]: # View titles of all columns
         aq_data.columns
```

Out[79]: Index(['State Code', 'County Code', 'Site Num', 'Parameter Code', 'POC',
       'Latitude', 'Longitude', 'Datum', 'Parameter Name', 'Date Local',

---

**IBM Cloud Pak for Data**  Q Search in your workspaces   Upgrade  ? ⊘³  Tushar Panchal's Account ⌄   London ⌄  TP  ⠿

Projects / tushar_practical_18 / Tushar_Practical18   ↑ ⌄   ⇄   ☐ ⌄   ⬇   ⋮   </> ⓘ ↺ ⠿

File  Edit  View  Insert  Cell  Kernel  Help                          Not Trusted | Python 3.10 ◯ ⤢

◉ ⊕ ✛ ⎘ ⎗ ↑ ↓ ▶ Run ◉ ⟳ ⏭ Format Code ⌄ ▦        Memory:359.2 MB / 8 GB

```
In [80]: aq_data['Sample Measurement'].describe()

Out[80]: count    3.558683e+06
         mean     8.250732e+00
         std      9.166534e+00
         min     -5.000000e+00
         25%      2.000000e+00
         50%      5.000000e+00
         75%      1.130000e+01
         max      1.296000e+02
         Name: Sample Measurement, dtype: float64

In [81]: # print list of all unique monitoring site numbers
         aq_data['Site Num'].unique()

Out[81]: array([  23, 2059,   34,   19, 3002, 4011, 4019, 4020, 9997, 1011, 1028,
            5,    7,    9,   11,   12,   13, 2005,    8,    2, 1002, 1004,
         2007,  242, 2016, 4001, 5001, 1005, 1003,   14, 2012, 6001,   16,
          113, 1103, 1201, 1302, 1602, 1701, 4006, 4008, 5005, 6012, 9033,
            4,    1,    3,    6, 1016, 8001, 8005, 9001,   10,   15,   26,
           27,  306, 1234, 2002, 9004, 1006, 1008, 1014, 1017, 1022, 3005,
         4002, 8002, 1001, 1013, 1018, 1021, 1025, 2004, 2011, 4003, 3001,
         7004,   28, 7001, 7003, 9003,   25,   41,   43,   50,   51,   35,
           32,  108, 1065,   21,   18,   56,   63,   76, 3103,   22,   78,
           87,   30,   17,   67,   75, 1024, 1100,   29,   40, 2006, 4005,
         3003,   42,   44,   93,   94,   95, 1010,   20,  423,  480,  962,
           85,   86,  760,  761,  762,  540,  561, 1501, 1502, 1233,  110,
          133,  124,  125,   45,   60,   73,   37,   38,   48,  101,   65,
           33,   97, 9021, 1127,   80, 1376,  100, 4000, 5200,   46,   59,
         1069,   69, 1067, 1044,   55, 1034,   24,   47,  416, 1015, 1035,
         1039, 1050, 1052, 1066, 1037, 1051, 1053, 3009, 3011, 1068, 7011,
         3006, 3013, 2003, 7022,   31,  123,  456,  892,   99,  232, 2601,
          700,  200,  300])

In [82]: # get number of aq sites in a single state, "California"
         aq_data.loc[aq_data['State Name'] == "California"]['Site Num'].unique().shape
```

## » Task 3 : Visualize data trends via matplotlib graphs.