Name - Yash Lakhtariya
Enrollment number - 21162101012
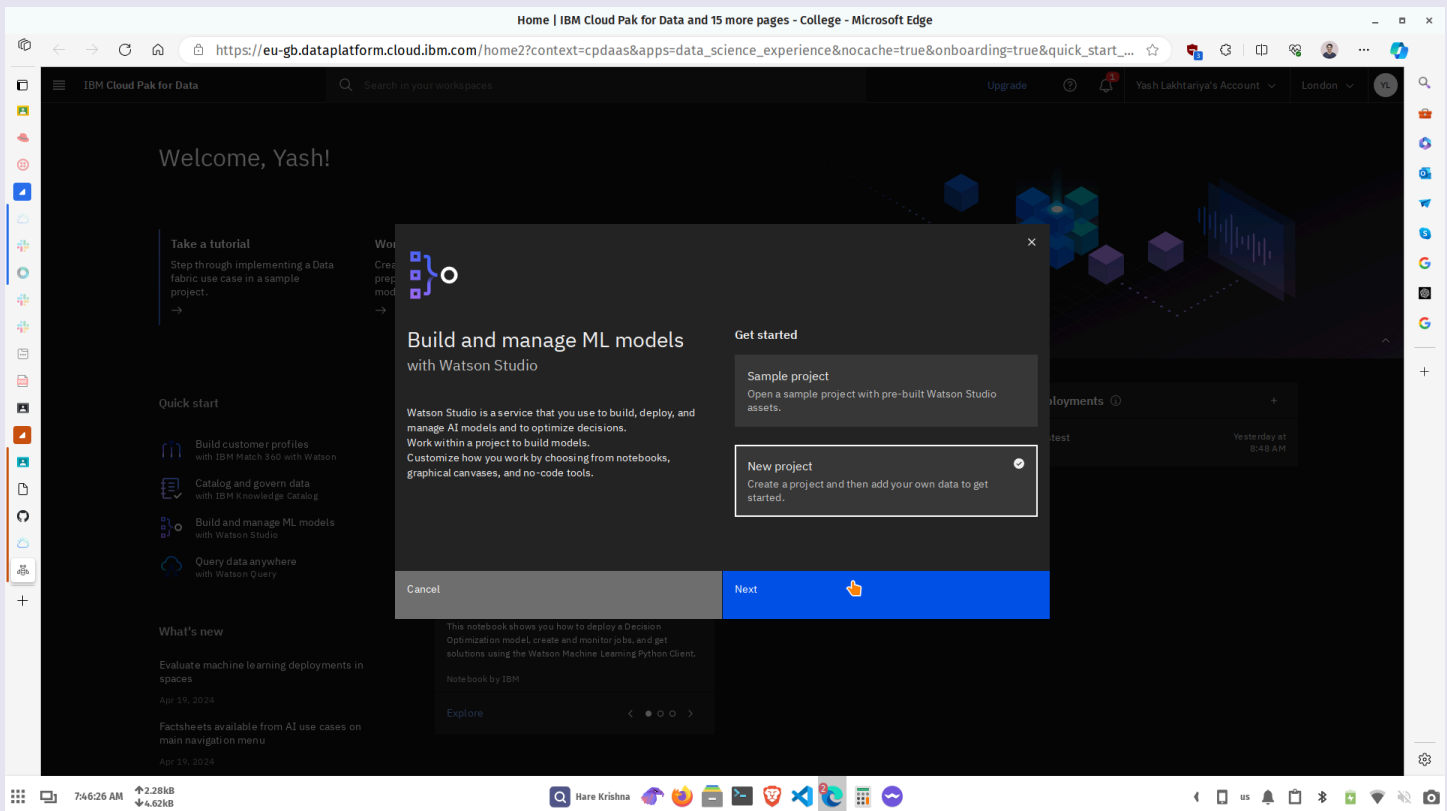Branch - CBA        Batch - 61
EADC Practical 18

<u>Aim</u> : 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 analyzing 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

<u>Steps and Screenshots</u> :

1.  Create new project on Watson Studio IBM instance

**Name - Yash Lakhtariya**
**Enrollment number - 21162101012**
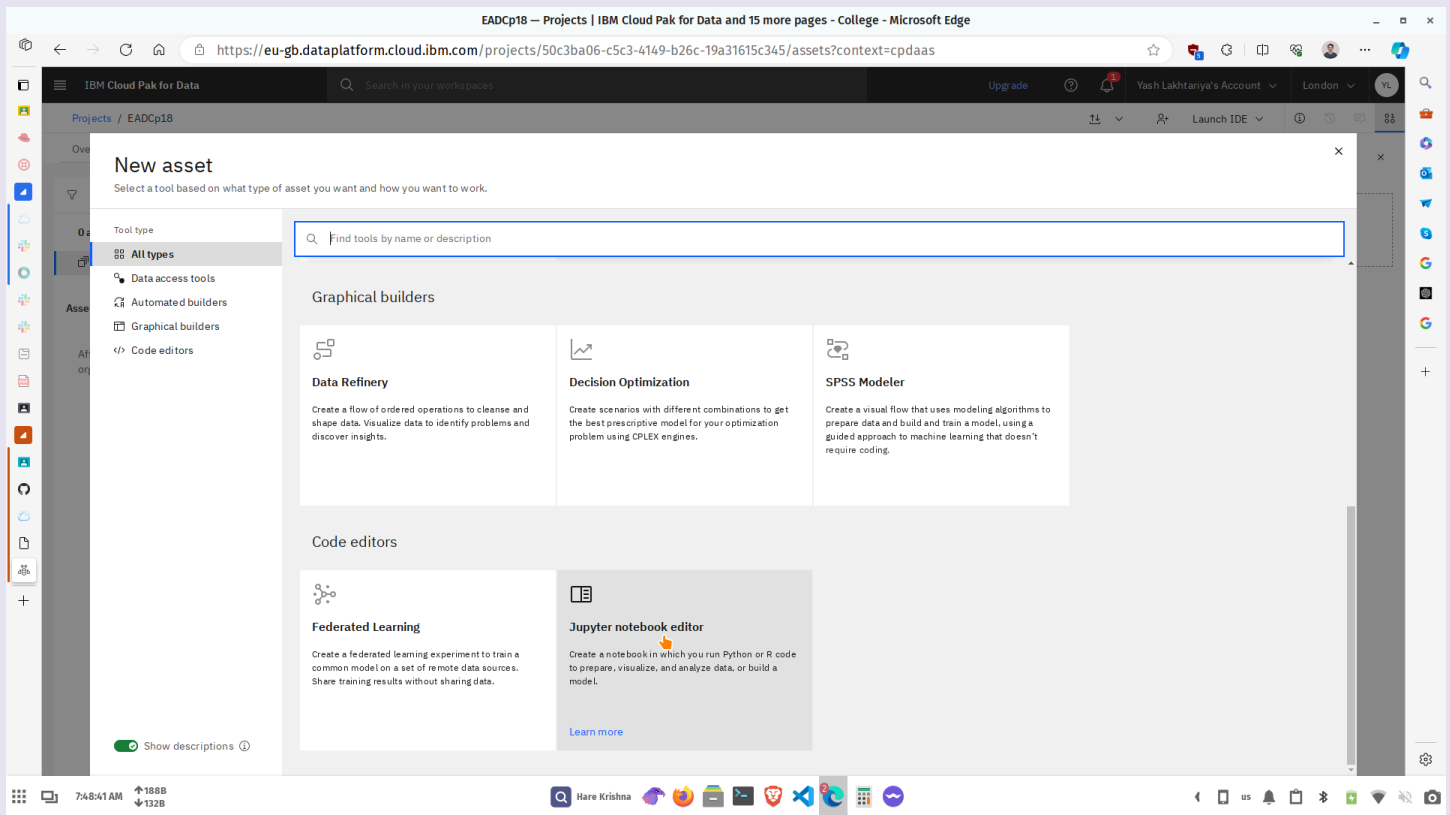**Branch - CBA        Batch - 61**
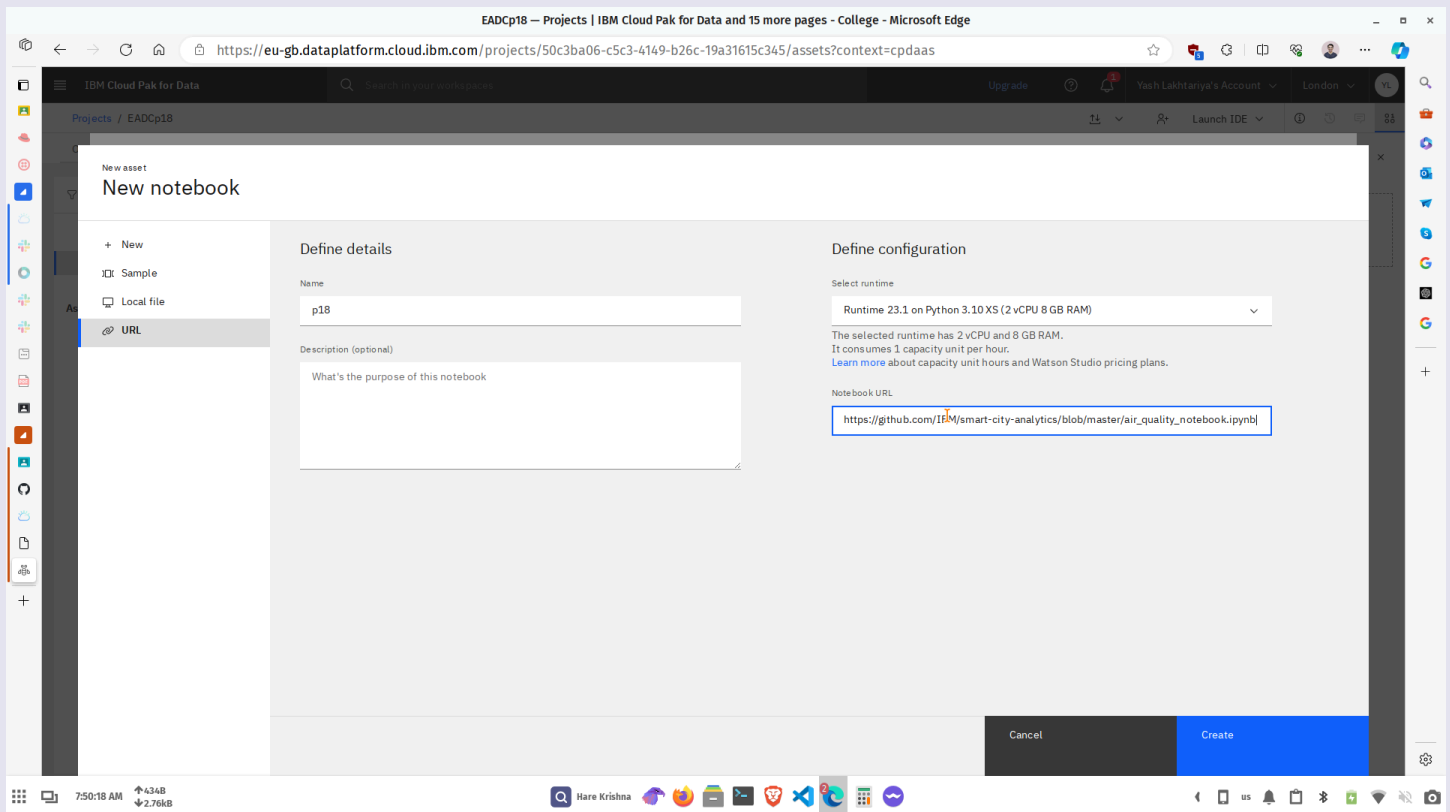**EADC Practical 18**

**Name - Yash Lakhtariya**
**Enrollment number - 21162101012**
**Branch - CBA          Batch - 61**
**EADC Practical 18**
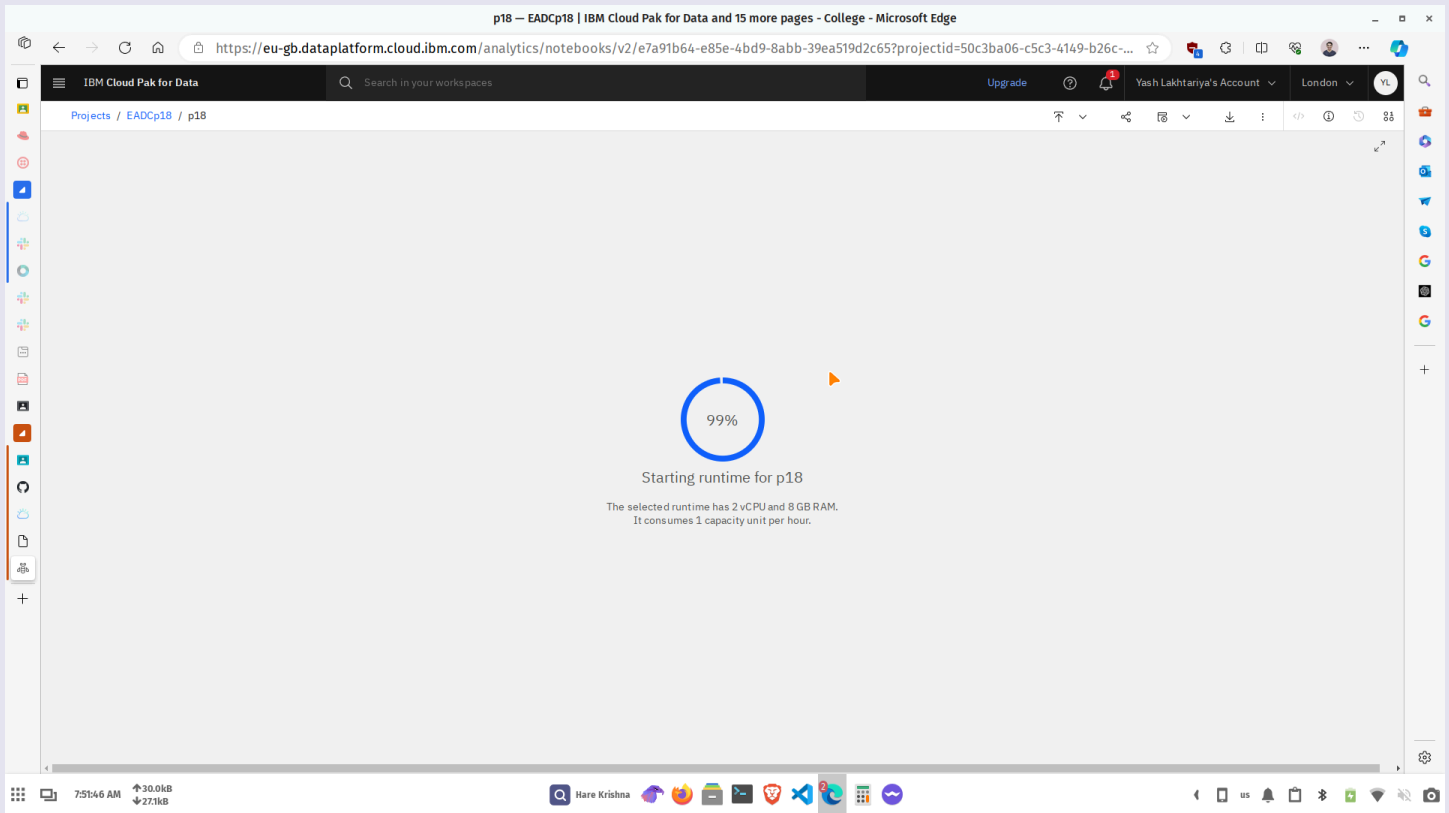
2. Add asset of Jupyter notebook

**Name - Yash Lakhtariya**
**Enrollment number - 21162101012**
**Branch - CBA        Batch - 61**
**EADC Practical 18**

3. Add new notebook from URL source and provide URL for ipynb file, here from Github

**Name - Yash Lakhtariya**
**Enrollment number - 21162101012**
**Branch - CBA        Batch - 61**
**EADC Practical 18**

4. Wait for creation to be completed

**Name - Yash Lakhtariya**
**Enrollment number - 21162101012**
**Branch - CBA          Batch - 61**
**EADC Practical 18**

5. Now explore the notebook, first clear all outputs

**Name - Yash Lakhtariya**
**Enrollment number - 21162101012**
**Branch - CBA          Batch - 61**
**EADC Practical 18**

## 6. Check data imported from csv file downloaded via wget

7. Try adding columns derived from given columns of data, like Weekday column using Date given and getting day by datetime library, and plotting data day and timewise