# Hands on Lab: Getting Started with Apache Spark on Watson Studio (20 mins)

## Objectives

After completing this lab you will be able to:

- Use your IBM Cloud account to explore and create resources.
- Create a Watson Studio Service instance.
- Create a Jupyter Notebook on Watson Studio with a Apache Spark + Python kernel
- Run the notebook and inspect the outputs

> **Note:** If you already have an IBM Cloud account, please skip Exercise 1. Additionally if you also have a Watson Studio service created, skip Exercise 2 as well.

## Exercise 1: Create an IBM Cloud Account

Follow the steps in [Hands on Lab: IBM Cloud Service Creation](#) to create an IBM cloud account.

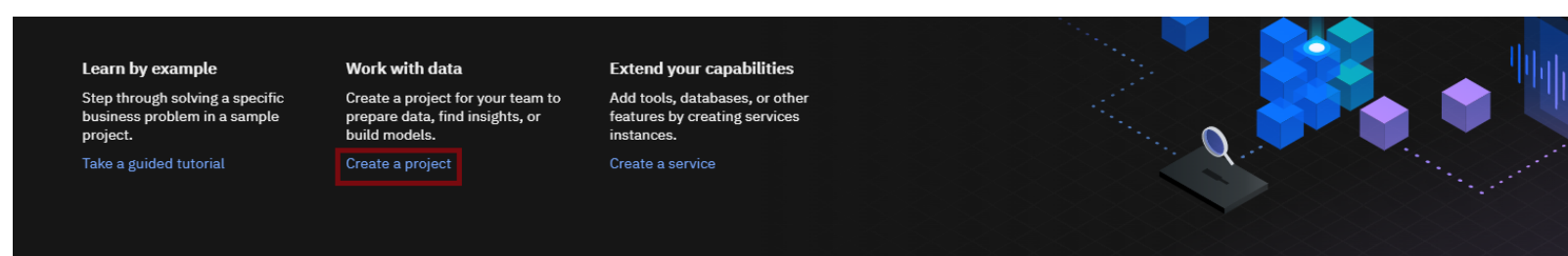## Exercise 2: Create an instance of Watson Studio service

Follow the steps in [Hands on Lab: IBM Watson Setup](#) to create a Watson Studio service and launch it.

## Exercise 3: Create a Spark + Python Jupyter notebook on Watson Studio

Once the Watson Studio service has been created and Watson Studio has been launched via the Cloud Pak for Data dashboard.

## Step 1: Creating a Watson Studio Project:

Click on **Create a project:**



On the Create a project page, click **Create an empty project**

Provide a **Project Name** and **Description**, as shown below:



You must also create storage for the project.

Click **Add**



On the Cloud Object Storage page, Select the 'Lite' plan and then click on **Create.** at the bottom.

You will be redirected to the Object storage page. If you do not see your instance active, please click on **Refresh** as below:

On the New project page, note that the storage has been added, and then click **Create.**

After creating the project you will need to add a Jupyter notebook to your project.

# Step 2: Adding a Notebook to the Project:

You need to add a Notebook to your project. Go to the **Assets** tab & Click on **New asset.**



Under **All types** select **Jupyter Notebook Editor**

On the New Notebook page, enter a name and description for the notebook, and then click From URL as shown below.



Important: **Select** `"Default Spark 3.0 & Python 3.9"` **as the runtime.**

This will initiate a kernel with Spark installed. Copy and paste the notebook URL - https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-BD0225EN-SkillsNetwork/labs/SparkIntro.ipynb for the **Apache Spark Python Intro** from this course into the **Notebook URL** box, and then click **Create Notebook**.

*Note*: For future Watson Studio labs that involve Jupyter notebooks, please replace the above notebook link with the relevant link or upload the notebook manually if needed.

## New notebook

| Blank | From file | **From URL** |
|---|---|---|

**Name**

Apache Spark Fundamentals

**Description (optional)**

Intro notebook to Apache Spark & IBM cloud

**Select runtime**

Default Spark 3.0 & Python 3.9 (Driver: 1 vCPU 4 GB RAM, 2 Executors: ⌄

The selected runtime uses 1 driver with 1 vCPU and 4 GB RAM, and 2 executors each with 1 vCPU and 4 GB RAM.
It consumes 1.5 capacity units per hour.
Learn more about capacity unit hours and Watson Studio pricing plans.

**Notebook URL**

https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBM-BI

| Cancel | Create |
|---|---|

You should see a loading screen like this:

Projects / Apache Spark Fundamentals / Apache Spark Fundamentals

32%

### Instantiating runtime for Apache Spark Fundamentals
The selected runtime uses 1 driver with 1 vCPU and 4 GB RAM, and 2 executors each with 1 vCPU and 4 GB RAM.
It consumes 1.5 capacity units per hour.

Click on **Set Kernel**

## Kernel not found

Could not find a kernel matching Python 3. Please select a kernel:

Python 3.7 with Spark ⌄

Continue Without Kernel     Set Kernel

Once the kernel has been initiated you will see the notebook like this. Please run all the cells to complete the lab.

File Edit View Insert Cell Kernel Help

Not Trusted | Python 3.7 with Spark

Format Markdown


cognitiveclass.ai logo

**Getting Started With Spark using Python**

Estimated time needed: **15** minutes



**The Python API**

Spark is written in Scala, which compiles to Java bytecode, but you can write python code to communicate to the java virtual machine through a library called py4j. Python has the richest API, but it can be somewhat limiting if you need to use a method that is not available, or if you need to write a specialized piece of code. The latency associated with communicating back and forth to the JVM can sometimes cause the code to run slower. An exception to this is the SparkSQL library, which has an execution planning engine that precompiles the queries. Even with this optimization, there are cases where the code may run slower than the native scala version. The general recommendation for PySpark code is to use the "out of the box" methods available as much as possible and avoid overly frequent (iterative) calls to Spark methods. If you need to write high-performance or specialized code, try doing it in scala. But hey, we know Python rules, and the plotting libraries are way better. So, it's up to you!

**Objectives**

In this lab, we will go over the basics of Apache Spark and PySpark. We will start with creating the SparkContext and SparkSession. We then create an RDD and apply some basic transformations and actions. Finally we demonstrate the basics dataframes and SparkSQL.

After this lab you will be able to:

- Create the SparkContext and SparkSession
- Create an RDD and apply some basic transformations and actions to RDDs

# Changelog

| Date | Version | Changed by | Change Description |
| --- | --- | --- | --- |
| 2021-07-15 | 1.0 | Karthik | Initial draft |
| 2021-08-17 | 1.1 | Karthik | Post Beta feedback |
| 2022-02-22 | 1.2 | K Sundararajan | Instructions Updated |
| 2022-04-06 | 1.3 | Sourabh | Images Updated |