# SPARK POOL

## Step 1

### Create Spark pool

**Manage → Apache Spark Pool → New**





Under Advanced settings change the idle time to 30 mins.

Click on **create**



## Step 2

Create folder by name "**mydata**" under the container and upload the **basestations_parquet** file.

Go to **Data Hub –> Linked → Expand the Azure data lake storage →**

## Select the container

Click on **New folder** from top menu and create the folder.



Now double click on the folder and click on the **upload** and select the **base_stations.parquet** file to upload.

Select the file and from 3 dots side of more select "properties"



Copy the **ABFSS** path



**Note:** Below is example do not copy the same, copy from your path.

**abfss://myfiles08@mydata08.dfs.core.windows.net/mydata/base_stations.parquet**

## Step 3
### Creating dataframe from parquet file.

Go to **Develop Hub** and select **Notebook**.

Give Notebook name as "**Sparkprog1_dataframe**"

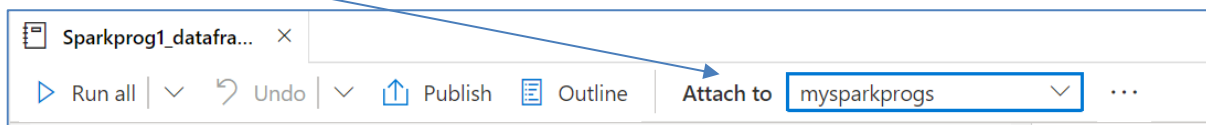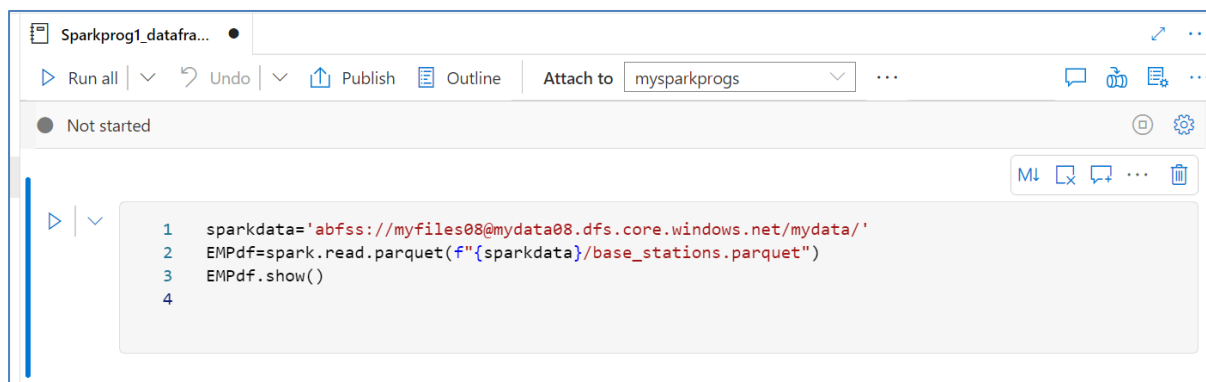Select **Attach to** from top menu and provide your spark pool name.



Enter the below under the code.

**Note:** Please provide your abfss path copied in earlier, under storage folder name.

**sparkdata='abfss://myfiles08@mydata08.dfs.core.windows.net/mydata/'**

**EMPdf=spark.read.parquet(f"{sparkdata}/base_stations.parquet")**

**EMPdf.show()**



```
+---+-----+------------+-----+-------+---------+
| id| zip|        city|state|    lat|      lon|
+---+-----+------------+-----+-------+---------+
|  1|86502|    Chambers|   AZ|35.2375| -109.523|
|  2|86514|Teec Nos Pos|   AZ|36.7797| -109.359|
|  3|85602|      Benson|   AZ|31.9883|-110.2941|
|  4|86011|   Flagstaff|   AZ|35.6308|-112.0524|
|  5|86016|Gray Mountain|  AZ|35.6308|-112.0524|
|  6|86018|       Parks|   AZ|35.2563|  -111.95|
|  7|86336|      Sedona|   AZ|34.8266|-111.7506|
|  8|85547|      Payson|   AZ|34.2575|-111.2878|
|  9|85548|     Safford|   AZ| 32.797|-109.7522|
| 10|85533|     Clifton|   AZ|33.1323|-109.2462|
```

Click on **Publish All** → Publish to save the notebook.

**Publish all**

You are about to publish all pending changes to the live environment. Learn more ⧉

**Pending changes (1)**

| NAME | CHANGE | EXISTING |
|---|---|---|
| ∨ Notebook | | |
| 🗐 Sparkprog1_dataframe | (New) | – |

[ Publish ]  [ Cancel ]

\*\*\* **Happy Learning** \*\*\*