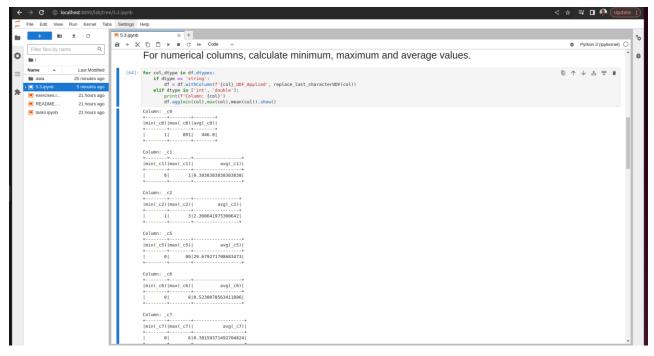# Talha Khan (2303.009.KHI.DEG)

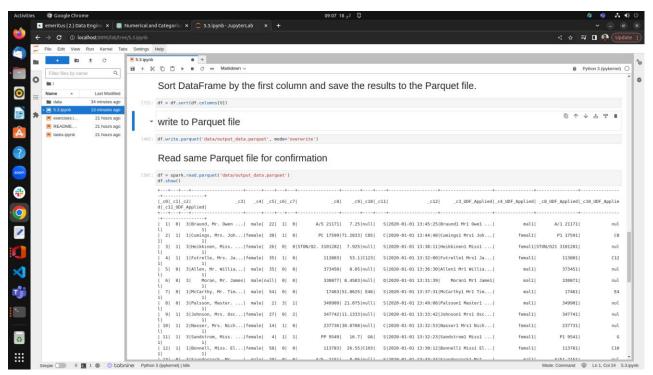# Muhammad Moiz Khan (2303.022.KHI.DEG)

# ASSIGNMNET 5.3





> ➢ We define a UDF called replace_last_character to replace the last character of a string column. We also create the UDF using udf and specify the return type as StringType().

➢ Next, we iterate over the columns of the DataFrame using df.dtypes. If the column type is 'string', we apply the UDF to that column using withColumn and create a new column with the suffix "_UDF_Applied".

➢ For numerical columns, we calculate statistics by using the agg function. The example code demonstrates how to print the minimum, maximum, and average values for each numerical column.

> For categorical columns, create and apply UDF that will change the last letter of every word to "1".

```
[47]: def replace_last_character(s):
          try:
              return ' '.join([word[:-1]+'1'
          for word in s.split(' ')])

          except:
              return s

      replace_last_characterUDF = udf(lambda z: replace_last_character(z),StringType())
```

> Sort DataFrame by the first column and save the results to the Parquet file.

➢ This code defines the replace_last_character function that takes a string s as input. It splits the string into words, removes the last character from each word, and appends '1'. It then joins the modified words back into a string.

➢ The replace_last_character_udf is created using a lambda function that calls replace_last_character and wraps it with the udf function. It specifies the return type as StringType().

➢ Finally, we sort the DataFrame by the first column and save the sorted DataFrame to a Parquet file, overwriting it if it already exists, using the write method and mode('overwrite').

➢ Then, we use the spark.read.parquet() method to read the Parquet file located at 'data/output_data.parquet' into a DataFrame named df. Finally, we use df.show() to display the contents of the DataFrame.