

Assignment 03

Pranathi Naga Sai

November 21, 2024

#Load the dataset

```
import pandas as pd
import plotly.express as px
import plotly.io as pio
from pyspark.sql import SparkSession
import re
import numpy as np
import plotly.graph_objects as go
from pyspark.sql.functions import col, split, explode, regexp_replace, transform, when
from pyspark.sql import functions as F
from pyspark.sql.functions import col, monotonically_increasing_id

np.random.seed(42)

pio.renderers.default = "notebook"

# Initialize Spark Session
spark = SparkSession.builder.appName("LightcastData").getOrCreate()

#Load the data
df = spark.read.option("header", "true").option("inferSchema", "true").option("multiLine", "true").load("s3://lightcast-data/lightcast_data.csv")

# Show Schema and Sample Data
#print("---This is Diagnostic check, No need to print it in the final doc---")

#df.printSchema() # comment this line when rendering the submission
#df.show(5)
```

Setting default log level to "WARN".

```
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).  
25/09/24 19:27:53 WARN NativeCodeLoader: Unable to load native-hadoop library for your platf  
java classes where applicable  
[Stage 1:> (0 + 1) / 1]
```