# Spark Hands-On

Dr. Villanes

# Once you see this screen.. Type pyspark
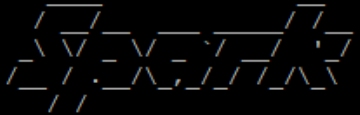
# You should now see something like this…

# Spark Version

>>> spark.version

# Creating Spark Dataframe

>>> df = spark.read.csv("s3a://bucket_name/airlines.csv",header=True)

# View the Dataframe

>>> df.show()

# How many records are in the Dataframe?

>>> df.count()

# Display descriptive stats using the describe function

>>> df.describe().show()

# Select only those observations where location is Ireland or Canada

>>> df[df.location.isin("Ireland","Canada")].show( )

# Create a Spark temporary table (so that we can query the table using SparkSQL)

\>\>\> df.createOrReplaceTempView("airlines")

# Use SparkSQL to query the temporary table that you just loaded in the previous step

>>> sqlDF = spark.sql("SELECT * FROM airlines").show()