# Lab 2 BDAT 1008
# Loading Data with Schema

**Nikhilesh Ramabhadraiah Vashisht**
**200512808**

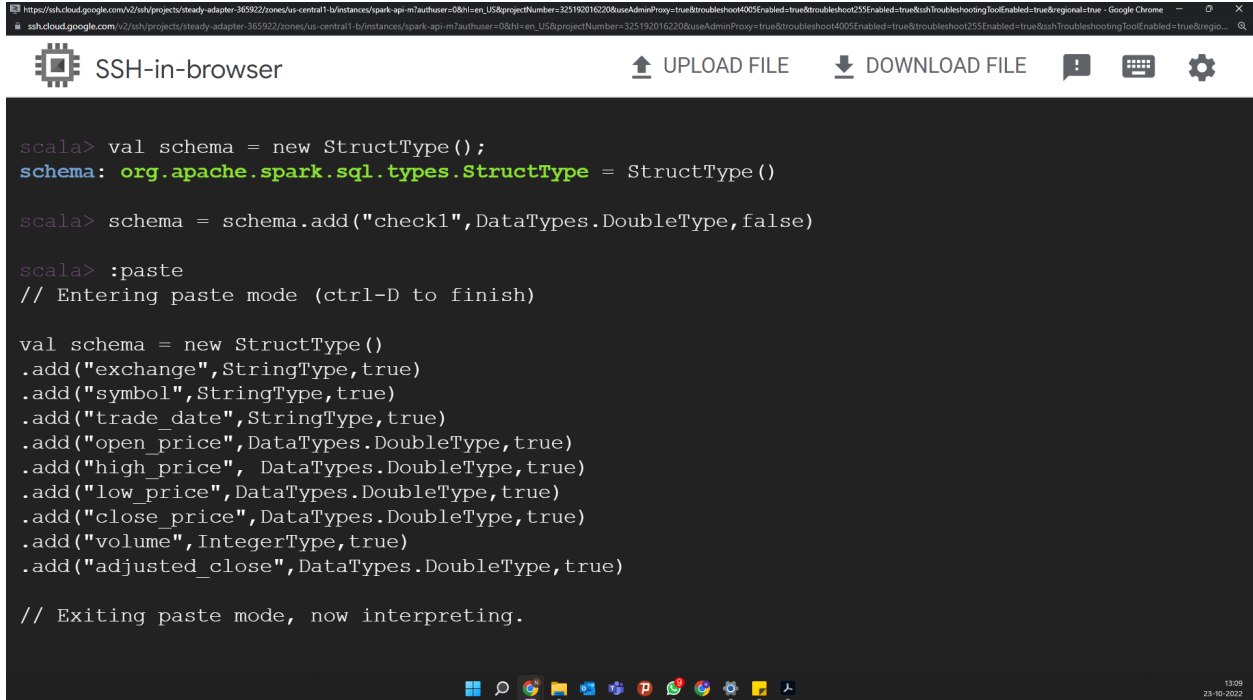# Create a schema for the data frame using an object of class StructType

Importing the class:
**import org.apache.spark.sql.Row;**
**import org.apache.spark.sql.types.StructField;**
**import org.apache.spark.sql.types.StructType;**
**import org.apache.spark.sql.types.StringType;**
**import org.apache.spark.sql.types.DataTypes._;**

Specifying schema:
**val schema = new StructType()**
**.add("exchange",StringType,true)**
**.add("symbol",StringType,true)**
**.add("trade_date",StringType,true)**
**.add("open_price",DataTypes.DoubleType,true)**
**.add("high_price", DataTypes.DoubleType,true)**
**.add("low_price",DataTypes.DoubleType,true)**
**.add("close_price",DataTypes.DoubleType,true)**
**.add("volume",IntegerType,true)**
**.add("adjusted_close",DataTypes.DoubleType,true)**

```scala
scala> val schema = new StructType();
schema: org.apache.spark.sql.types.StructType = StructType()

scala> schema = schema.add("check1",DataTypes.DoubleType,false)

scala> :paste
// Entering paste mode (ctrl-D to finish)

val schema = new StructType()
.add("exchange",StringType,true)
.add("symbol",StringType,true)
.add("trade_date",StringType,true)
.add("open_price",DataTypes.DoubleType,true)
.add("high_price", DataTypes.DoubleType,true)
.add("low_price",DataTypes.DoubleType,true)
.add("close_price",DataTypes.DoubleType,true)
.add("volume",IntegerType,true)
.add("adjusted_close",DataTypes.DoubleType,true)

// Exiting paste mode, now interpreting.
```
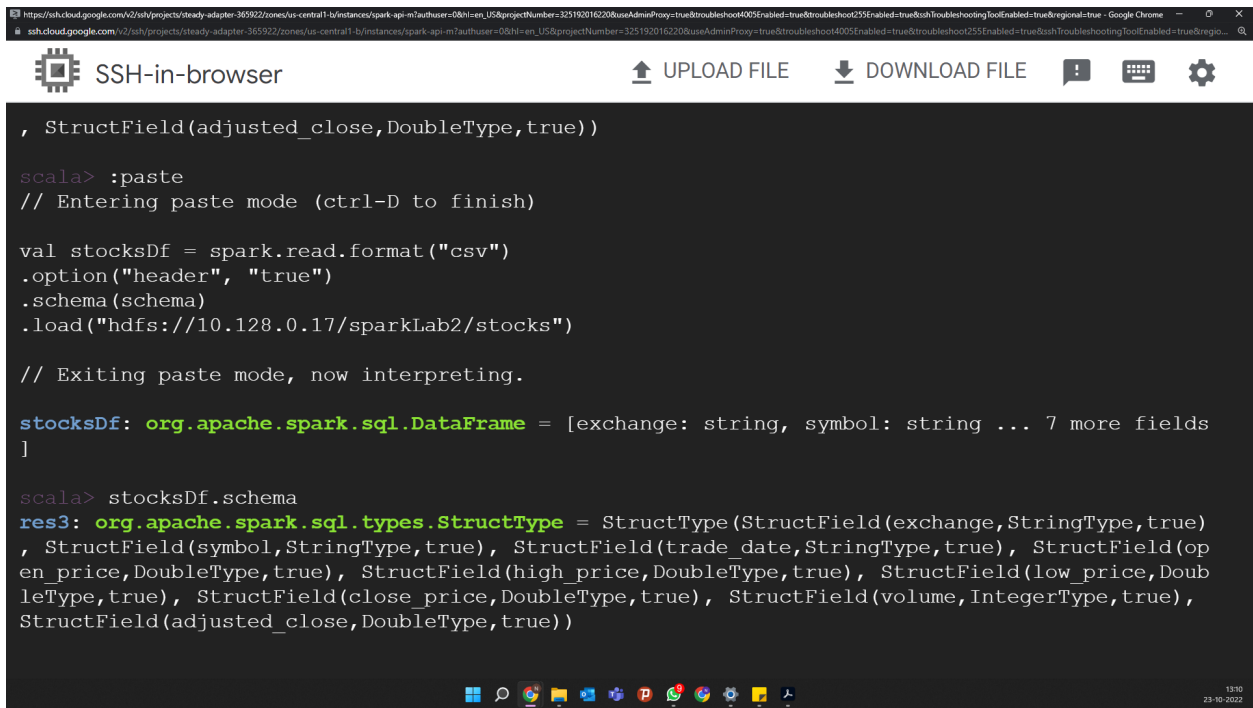
Loading the stocks data
**val stocksDf = spark.read.format("csv")**
**.option("header", "true")**
**.schema(schema)**
**.load("hdfs://10.128.0.17/sparkLab2/stocks")**



```scala
, StructField(adjusted_close,DoubleType,true))

scala> :paste
// Entering paste mode (ctrl-D to finish)

val stocksDf = spark.read.format("csv")
.option("header", "true")
.schema(schema)
.load("hdfs://10.128.0.17/sparkLab2/stocks")

// Exiting paste mode, now interpreting.

stocksDf: org.apache.spark.sql.DataFrame = [exchange: string, symbol: string ... 7 more fields
]

scala> stocksDf.schema
res3: org.apache.spark.sql.types.StructType = StructType(StructField(exchange,StringType,true)
, StructField(symbol,StringType,true), StructField(trade_date,StringType,true), StructField(op
en_price,DoubleType,true), StructField(high_price,DoubleType,true), StructField(low_price,Doub
leType,true), StructField(close_price,DoubleType,true), StructField(volume,IntegerType,true),
StructField(adjusted_close,DoubleType,true))
```

# Question 1

**Write a command to find the stocks with average daily volume larger than 1 million shares**

val Q1 =
stocksDf.select(col("symbol"),col("trade_date"),col("volume")).groupBy(col("symbol")).agg(avg(col("volume"))).filter(avg(col("volume")) > 1000000).show()

# Question 2

**Write a Scala DataFrame query to find the top 3 stocks by volume for the year 2004.**

val Q2 =
stocksDf.select(col("exchange"),col("symbol"),col("trade_date"),col("volume")).filter(col("trade_date").contains( "2004")).orderBy(col("volume").desc).show(3)

# Question 3

**Write a Scala DataFrame query to find the top 3 stocks by volume and whose symbol start with the first letter of your name (example for Saber, it is symbols starting with "S").**

val Q3 =
stocksDf.select(col("exchange"),col("symbol"),col("volume")).filter(col("symbol").like("N%")).orderBy(col("volume").desc).show(3)



There were no symbols starting with the first letter of my name (N)

# Question 4

**Write a Scala DataFrame to find all the stocks symbols whose closing price is larger than your age.**

val Q4 = stocksDf.select(col("exchange"),col("symbol"),col("close_price")).filter(col("close_price")>24).show()

# Question 5

**Write a Scala DataFrame to find the top 10 stocks with largest intraday price change (difference between high and low price during a trading day) and also display the amount of the change.**

val Q5 =
stocksDf.select(col("exchange"),col("symbol"),col("high_price"),col("low_price"),(col("high_price")-col("low_price")).alias("price_change_in_a_day")).orderBy(col("price_change_in_a_day").desc).show(10)