





### Tecnológico Nacional De México

### Instituto Tecnológico De Tijuana

Subdirección Académica

Departamento de Sistemas y Computación

Semestre Enero - Junio 2022

Ingeniería Informática

**Datos Masivos** 

Práctica Evaluatoria

Unidad 1

Israel López Pablo No.17210585

Perez Ortega Victoria Valeria No.18210718

JOSE CHRISTIAN ROMERO HERNANDEZ

Tijuana, B.C. a 22 de Marzo de 2022.







#### **TECNOLÓGICO NACIONAL DE MÉXICO**

#### INSTITUTO TECNOLÓGICO DE TIJUANA SUBDIRECCIÓN ACADÉMICA

Departamento de Sistemas y Computación **EXAMEN** 

Carrera: Ingeniería En Sistemas Computacionales/ Tecnologías de la información/ Informática Período: **Febrero-Junio 2022** Materia: Datos Masivos Grupo: BDD-1704SC9C Salón: Unidad (es) a evaluar: Unidad 1 Tipo de examen:Practico Fecha: Catedrático: José Christian Romero Hernandez Firma del maestro: Calificación:

Alumnos: Israel Lopez Palo y Perez Ortega Victoria Valeria No. Control: 17210585, 18210718

#### Instrucciones

Responder las siguientes preguntas con Spark DataFrames y Scala utilizando el "CSV" Netflix\_2011\_2016.csv que se encuentra en la carpeta de spark-dataframes.

Comienza una simple sesión Spark.

Para iniciar sesión en Spark, primero abrimos el cmd como administrador, ejecutamos Spark con spark-shell e inmediatamente importamos una biblioteca para poder iniciar sesión

import org.apache.spark.sql.SparkSession

Declaramos la variable session y usamos la librería para crear una nueva sesión

val session = SparkSession.builder().getOrCreate()







2. Cargue el archivo Netflix Stock CSV, haga que Spark infiera los tipos de datos.

```
val netflix =
spark.read.option("header","true").option("inferSchema","tr
ue")csv("Netflix_2011_2016.csv")
```

```
scala> val netflix = spark.read.option("header","true").option("inferSchema","true")csv("Netflix_2011_2016.csv")
netflix: org.apache.spark.sql.DataFrame = [Date: timestamp, Open: double ... 5 more fields]
```

3. ¿Cuáles son los nombres de las columnas?

Mandamos a llamar a nuestra variable netflix que con ese nombre cargamos el archivo CSV de Netflix Stock, y lo que queremos ver son las columnas

#### netflix.columns

```
scala> netflix.columns
res0: Array[String] = Array(Date, Open, High, Low, Close, Volume, Adj Close)
```

#### 4. ¿Cómo es el esquema?

Mandamos a llamar a nuestra variable netflix que con ese nombre cargamos el archivo Netflix Stock CSV, mandamos a imprimir el esquema que es el que queremos ver

#### netflix.printSchema()

```
scala> netflix.printSchema()
root
    |-- Date: timestamp (nullable = true)
    |-- Open: double (nullable = true)
    |-- High: double (nullable = true)
    |-- Low: double (nullable = true)
    |-- Close: double (nullable = true)
    |-- Volume: integer (nullable = true)
    |-- Adj Close: double (nullable = true)
```



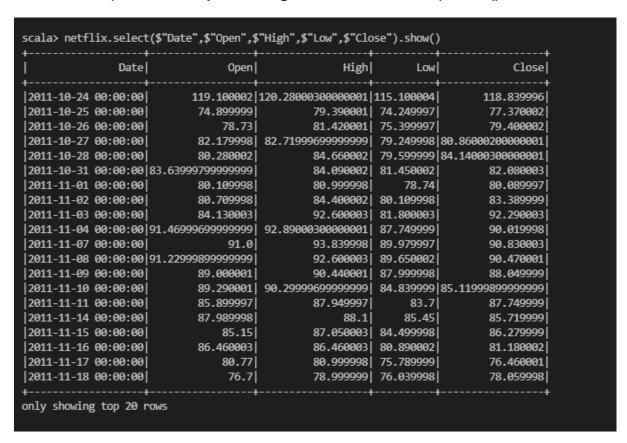




5. Imprime las primeras 5 columnas.

Mandamos a llamar a nuestra variable de netflix que con ese nombre cargamos el archivo CSV de Netflix Stock y seleccionamos las primeras 5 columnas

#### netflix.select(\$"Date",\$"Open",\$"High",\$"Low",\$"Close").show()



6. Usa describe () para aprender sobre el DataFrame.

#### netflix.describe().show()

scala> netflix.describe().show()						
t  summary	Open		Low	Close	Volume	Adj Close
count	1259	1259   233.97320872915006	1259	1259	1259	1259
stddev 164	.37456353264244	165.9705082667129	162.6506358235739	164.40918905512854	2.5634836060365368E7 2.306312683388607E7	35.186669331525486
min    max	53.990001   708.900017	55.480001    716.159996	52.81   697.569984	53.8   707.610001	3531300   315541800	
+						







7. Crea un nuevo dataframe con una columna nueva llamada "HV Ratio" que es la relación que existe entre el precio de la columna "High" frente a la columna "Volumen" de acciones negociadas por un día. Hint - es una operación

Declaramos una nueva variable llamada netflixnew

var netflixnew = netflix.withColumn("HV Ratio", netflix("High")/netflix("Volume"))

```
scala> var netflixnew = netflix.withColumn("HV Ratio", netflix("High")/netflix("Volume"))
netflixnew: org.apache.spark.sql.DataFrame = [Date: timestamp, Open: double ... 6 more fields]
```

8. ¿Qué día tuvo el pico más alto en la columna "Open"?

Declaramos una variable llamada columnday y del archivo de netflix tomamos los datos de la columna day

val columnday = netflix.withColumn("Day",dayofmonth(netflix("Date")))

Declaramos una variable llamada netflixdaycolumn y vamos a agrupar todos los valores máximos del día

val netflixdaycolumn = columnday.groupBy("Day").max()

llamamos a la variable que declaramos para mostrar la tabla

netflixdaycolumn.show()

mandamos a llamar a la variable que declaramos y seleccionamos el dia maximo de la columna maxopen

netflixdaycolumn.select(\$"Day",\$"max(Open)").show()

```
scala> val columnday = netflix.withColumn("Day",dayofmonth(netflix("Date")))
columnday: org.apache.spark.sql.DataFrame = [Date: timestamp, Open: double ... 6 more fields]
scala> val netflixdaycolumn = columnday.groupBy("Day").max()
netflixdaycolumn: org.apache.spark.sql.DataFrame = [Day: int, max(Open): double ... 6 more fields]
```







```
scala> netflixdaycolumn.select($"Day",$"max(Open)").show()
|Day|
            max(Open)
 31
           430.259995
 28
            628.000008
 26
            667.100021
 27
            617.000023
  12 664.4100269999999
 22
            665.29998
  1 663.6400219999999
 13 686.6900019999999
            659.700012
           654.309982
  6
           624.700012
  20
            617.500008
            624.500008
 19
            673.700012
            649.999977
            664.300011
  17
            665.930008
            618.649994
            654.299995
            674.350014
only showing top 20 rows
```

9. ¿Cuál es el significado de la columna Cerrar "Close" en el contexto de información financiera, explíquelo que no hay que codificar nada?

La columna cerrar significa que la empresa finalizó el día o en este caso que netflix finalizó el día

10. ¿Cuál es el máximo y mínimo de la columna "Volumen"?Del archivo netflix mandamos a llamar el valor maximo y minimo de la

```
scala> netflix.select(max("Volume"), min("Volume")).show()
+-----+
|max(Volume)|min(Volume)|
+----+
| 315541800| 3531300|
+-----+
```







### 11.Con Scala/Spark Syntax \$ responda lo siguiente:

## a. ¿Cuántos días estuvo la columna "Close" menos de \$600?

Para esto solo necesitamos un filtro "cerrar" con una condición y esa condición es menor.

```
val result =netflix.filter($"Close" < 600 ).count()</pre>
```

#### Result:

aquí podemos ver el resultado

```
scala> val result =netflix.filter($"Close" < 600 ).count()
result: Long = 1218</pre>
```

## b. ¿¿Qué porcentaje de veces la columna "High" superó los \$500?

Este es el caso contrario que la pregunta anterior.

```
val result = (netflix.filter( $"High" >
500).count()*1.0/netflix.count())*100
```

#### Result:

Aquí podemos ver el resultado.

```
scala> val result = (netflix.filter( $"High" > 500).count()*1.0/netflix.count())*100 result: Double = 4.924543288324067
```

# c. ¿Cuál es la correlación de Pearson entre la columna "High" y la columna "Volume"?

El coeficiente de correlación de Pearson es una prueba que mide la relación estadística entre dos variables continuas. así que en este caso específicamente







#### estamos usando para dos columnas High y Volume

```
netflix.select(corr($"High",$"Volume")).show()
```

#### Result:

Este es el resultado de usar la correlación de Pearson.

```
scala> netflix.select(corr($"High",$"Volume")).show()
+-----+
| corr(High, Volume)|
+-----+
|-0.20960233287942157|
+-----+
```

### d. ¿Cuál es el máximo de la columna "High" por año?

En este caso necesitamos obtener un máximo por lo que usamos una función max para obtener el resultado

```
netflix.select(max($"High")).show()
```

#### Result:

Estos son los resultados

```
scala> netflix.select(max($"High")).show()
+----+
| max(High)|
+----+
|716.159996|
+----+
```







## e. ¿Cuál es el promedio de columna "Close" para cada mes del calendario?

Este es diferente a todos los demás ya que debemos obtener ciertos valores los cuales necesitan ser modificados varias veces, cómo agregar una nueva columna que almacene los meses y usar métodos como obtener el mes, usar el promedio y agrupar para obtener el resultado deseado.

```
val meses = netflix.withColumn("Meses", month(netflix("Date")))
val porcentajemes=
meses.select($"Meses",$"Close").groupBy("Meses").mean()
porcentajemes.select($"Meses",$"avg(Close)").show()
```

#### Result:

```
scala> val meses = netflix.withColumn("Meses",month(netflix("Date")))
meses: org.apache.spark.sql.DataFrame = [Date: timestamp, Open: double ... 6 more fields]
scala> val porcentajemes= meses.select($"Meses",$"Close").groupBy("Meses").mean()
porcentajemes: org.apache.spark.sql.DataFrame = [Meses: int, avg(Meses): double ... 1 more field]
scala> porcentajemes.select($"Meses",$"avg(Close)").show()
Meses
               avg(Close)
    12 | 199.3700942358491 |
    1 212.22613874257422
    6 295.1597153490566
     3 | 249.5825228971963 |
     5 | 264.37037614150944 |
    9 | 206.09598121568627 |
     4 246.97514271428562
    8 | 195.25599892727263
     7 243.64747528037387
    10 | 205.93297300900903 |
    11 | 194.3172275445545
    2 254.1954634020619
```

URL: <a href="https://www.youtube.com/watch?v=02E6qq6Crpc">https://www.youtube.com/watch?v=02E6qq6Crpc</a>