

Titanic visualization assignment.

READY

This assignment shows in a nice way the data set from Titanic passenger list

```
1 %spark
2
3 import org.apache.spark.sql.Dataset;
4 import org.apache.spark.sql.functions.expr
5
6 //import org.apache.spark.sql.types.{StructType, StructField, StringType, IntegerType};
7
8 // import org.apache.commons.io.IOUtils
9 // import java.net.URL
10 // import java.nio.charset.Charset
11 // import org.apache.spark.ml.classification.DecisionTreeClassificationModel;
12 // import org.apache.spark.ml.classification.DecisionTreeClassifier;
13 // import org.apache.spark.ml.feature.StringIndexer;
14 // import org.apache.spark.ml.feature.StringIndexerModel;
15 // import org.apache.spark.ml.feature.VectorAssembler;
16 // import org.apache.spark.mllib.evaluation.MulticlassMetrics;
17 // import org.apache.spark.sql.SparkSession;
18 // import org.apache.spark.sql.Row;
19 // import org.apache.spark.sql.functions;
20 // import org.apache.spark.sql.types.DataTypes;
21 // import org.apache.spark.sql.types.StructType;
22
23 // création du dataframe Spark
24 val df = spark.read.options(Map("inferSchema"->"true","delimiter"->"," ,"header"->"true"))
25   .csv("/Users/FXSUIRE/Google Drive/_00_Dev/Titanic/input.csv")
26
27 // application du schéma de données
28 // colonnes d'origine
29 case class passenger_cls(Survived:Int,Pclass:Int,Name:String,Sex:String,Age:Float,Siblings:Int,ParentAboard:Int,Fare:Float)
30
31 case class passenger_addedcol(Survived:Int,Pclass:Int,Name:String,Sex:String,Age:Float,Siblings:Int,ParentAboard:Int,Fare:Float,Fa
32
33 //Ratio d'actualisation du prix du ticket en liver sterling 2020
```

FINISHED

TitanicV3

```

34 // source: https://www.thisismoney.co.uk/money/bills/article-1633409/Historic-inflation-calculator-value-money-changed-1900.html
35 val UpdatingRatio = 115.8
36
37 // FareActualised = "Fare*UpdatingRatio"
38 val df2 = df.withColumn("FareUpdated", df("Fare") * UpdatingRatio)
39
40 // x0: Survived: Int
41 // x1: Pclass: Int
42 // x2: Name: String
43 // x3: Sex: String
44 // x4: Age: Float
45 // x5: Siblings: Int
46 // x6: ParentAboard: Int
47 // x7: Fare: Float
48
49 // x8: Actualised (2020) Fare: float
50
51 // mapping des colonnes originales
52 val passenger = data1.map(x=>x.split(",")).map(x => passenger_cls(x(0).toInt,x(1).toInt,x(2),x(3),x(4).toFloat,x(5).toInt,x(6).toI
53
54 //enregistrement dans une table temporaire des données des passagers
55 passenger.registerTempTable("PASSENGER")
56
57 //sélection des données que l'on souhaite utiliser dans nos requêtes de visualisations
58
59 // Décès-survie du passager, classe, sex, prix du ticket
60 val SurvivalAgeClassSexFare = sqlContext.sql("select Survived, Pclass, Sex, Fare from PASSENGER")
61 SurvivalAgeClassSexFare.registerTempTable("SMALLPASSTABLE")
62
63 // Statistiques basiques sur le prix des tickets
64 val BasicStats = sqlContext.sql("select min(fare), max(fare), avg(fare), stddev(fare) from PASSENGER where fare>0").toDF
65
66 val BasicStatsRenamed = BasicStats.withColumnRenamed("min(fare)","Minimum fare price")
67     .withColumnRenamed("max(fare)","Maximum fare price")
68     .withColumnRenamed("avg(fare)","Average fare price")
69     .withColumnRenamed("stddev(fare)","Standard deviation")
70
71 BasicStatsRenamed.registerTempTable("BASICFARETABLE")
72
73 // Décès-survie du passager, classe, sex, prix du ticket 2020
74
75 //enregistrement dans une table temporaire des données des passagers avec ajout des colonnes Fareupdated
76 df2.registerTempTable("PASSPROCADATA")
77
78 val SurvivalAgeClassSexFare2020 = sqlContext.sql("select Survived, Pclass, Sex, Fare, FareUpdated from PASSPROCADATA")

```

TitanicV3

```

79 SurvivalAgeClassSexFare2020.registerTempTable("SMALLPASSTABLE2020")
80
81 // statistiques sur le prix du ticket actualisé, c'est à dire en GBP de 2020
82
83 val BasicUpdFareStats = sqlContext.sql("select min(FareUpdated), max(FareUpdated), avg(FareUpdated), stddev(FareUpdated) from SMAL
84
85 val BasicUpdStatsRenamed = BasicUpdFareStats.withColumnRenamed("min(FareUpdated)","Minimum fare price")
86         .withColumnRenamed("max(FareUpdated)","Maximum fare price")
87         .withColumnRenamed("avg(FareUpdated)","Average fare price")
88         .withColumnRenamed("stddev(FareUpdated)","Standard deviation")
89
90 BasicUpdStatsRenamed.registerTempTable("BASICUPDFARETABLE")
91
92

```

warning: there were 6 deprecation warnings; re-run with -deprecation for details

```

import sqlContext.implicits._
import org.apache.spark.sql.Dataset
import org.apache.spark.sql.functions.expr
df: org.apache.spark.sql.DataFrame = [Survived: int, Pclass: int ... 6 more fields]
defined class passenger_cls
defined class passenger_addedcol
UpdatingRatio: Double = 115.8
df2: org.apache.spark.sql.DataFrame = [Survived: int, Pclass: int ... 7 more fields]
----- org.apache.spark.sql.DataFrame = [Survived: int, Pclass: int ... 6 more fields]

```

Took 6 sec. Last updated by anonymous at February 21 2020, 6:03:56 PM. (outdated)

Extract of Titanic passenger list (source data)

READY

```

1 %spark
2
3 // we display the raw data
4 df.show(false)
5 //df.printSchema()
6
7
8
9
10 11      |Mr. Timothy J McCarthy      |male |54.0|0      |10      |151.8
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100

```

11	13	IMiss. Marguerite Rut Sandstrom	female	14.0	11	110.7
11	11	IMiss. Elizabeth Bonnell	female	158.0	10	126.5
5	1					
10	13	IMr. William Henry Saunderson	male	120.0	10	18.05
10	13	IMr. Anders Johan Andersson	male	139.0	15	131.2
75	1					
10	13	IMiss. Hulda Amanda Adolfina Vestrom	female	14.0	10	17.85

TitanicV3

Extract of Titanic passenger list (source data + added columns)

FINISHED

```

1 %spark
2
3 // we display the raw data
4 df2.show(false)
5 //df2.printSchema()

```

11	11	IMrs. John Bradley (Florence Briggs Thayer) Cumings	female	138.0	11	171.2
833	18254.60614					
11	13	IMiss. Laina Heikkinen	female	126.0	10	17.92
5	1917.7149999999999					
11	11	IMrs. Jacques Heath (Lily May Peel) Futrelle	female	135.0	11	153.1
16	148.98					
10	13	IMr. William Henry Allen	male	135.0	10	18.05
19	32.19					
10	13	IMr. James Moran	male	127.0	10	18.45
83	1979.4711399999999					
10	11	IMr. Timothy J McCarthy	male	154.0	10	151.8
625	16005.6775					
10	13	IMaster. Gosta Leonard Palsson	male	12.0	11	121.0
75	12440.484999999997					
11	13	IMrs. Oscar W (Elisabeth Vilhelmina Berg) Johnson	female	127.0	12	111.1
333	11289.23614					
11	12	IMrs. Nicholas (Adele Achem) Nasser	female	14.0	10	130.0
708	13482.198639999999					

Took 2 sec. Last updated by anonymous at February 21 2020, 6:04:10 PM.

Number of deceased or living passenger by sex after the Shipwreck

FINISHED

```
%spark.sql  
select Survived,Sex, count(Survived) from SMALLPASSTABLE group by Survived, Sex
```



▼

settings ▼

TitanicV3

Survived	Sex
0	female
1	male
1	female
0	male

Took 1 sec. Last updated by anonymous at February 21 2020, 5:50:50 PM.

Share of Passenger who survived or not by sex

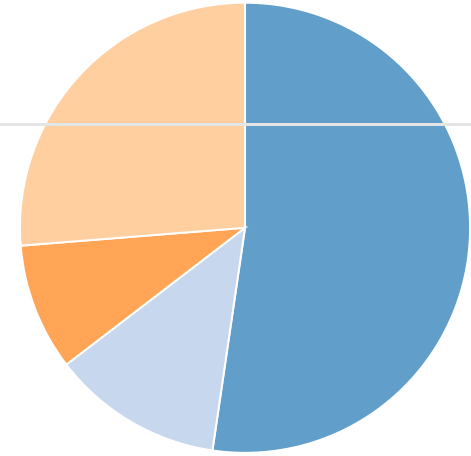
FINISHED



▼

settings ▼

TitanicV3

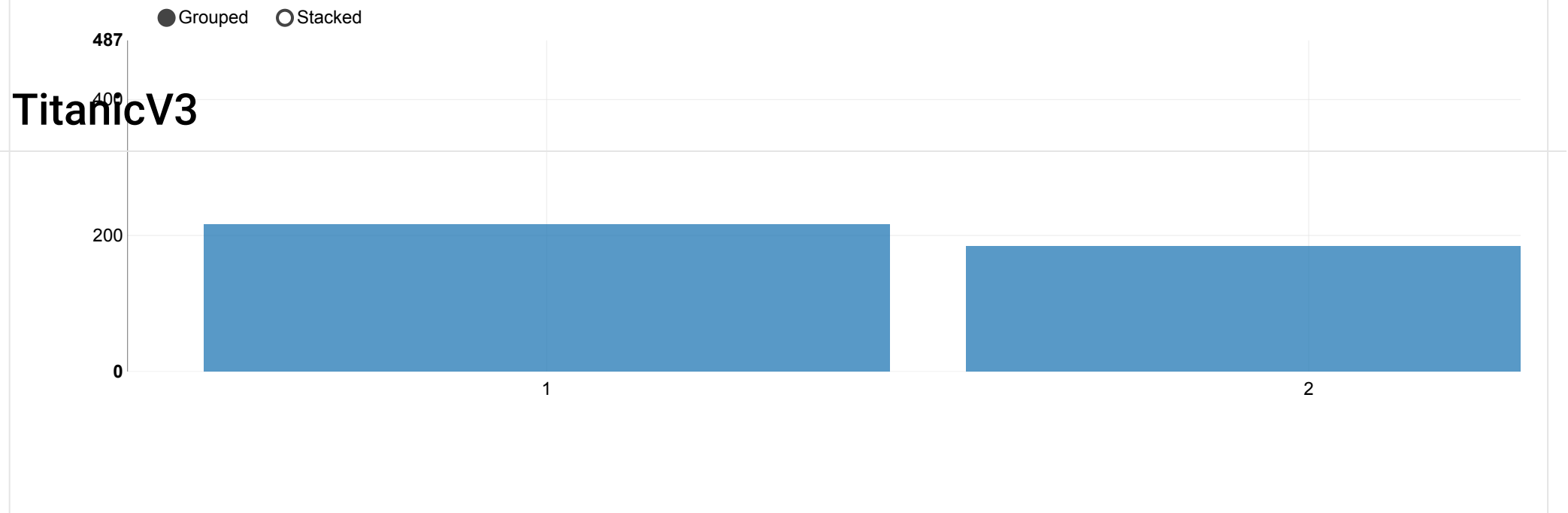


Took 1 sec. Last updated by anonymous at February 21 2020, 6:06:22 PM.

Titanic passenger class distribution

READY

settings ▼



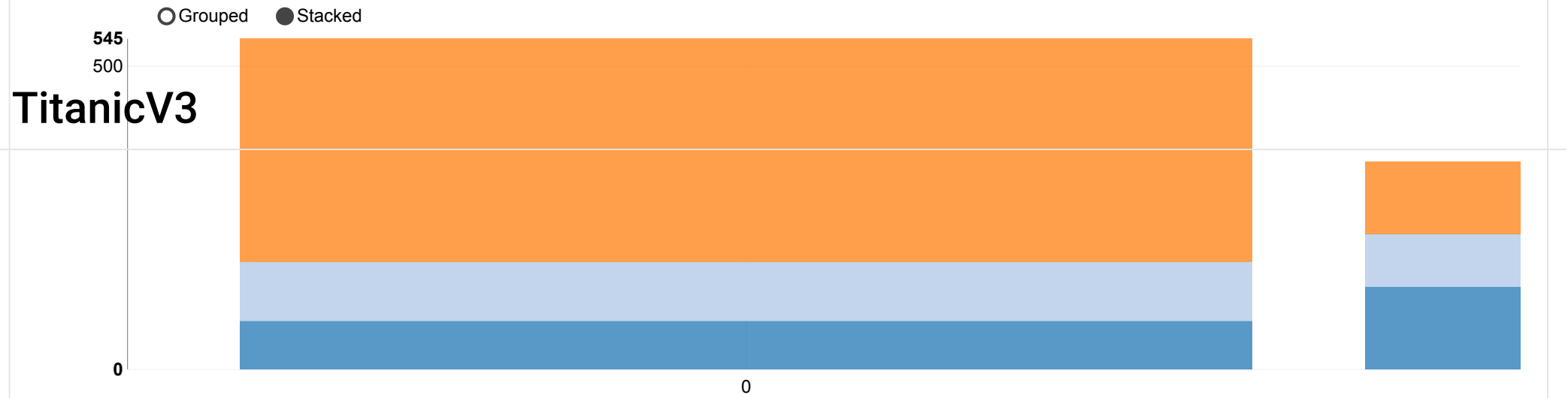
```
%spark.sql  
select Survived,Pclass from SMALLPASSTABLE order by Pclass asc
```

READY





settings ▼

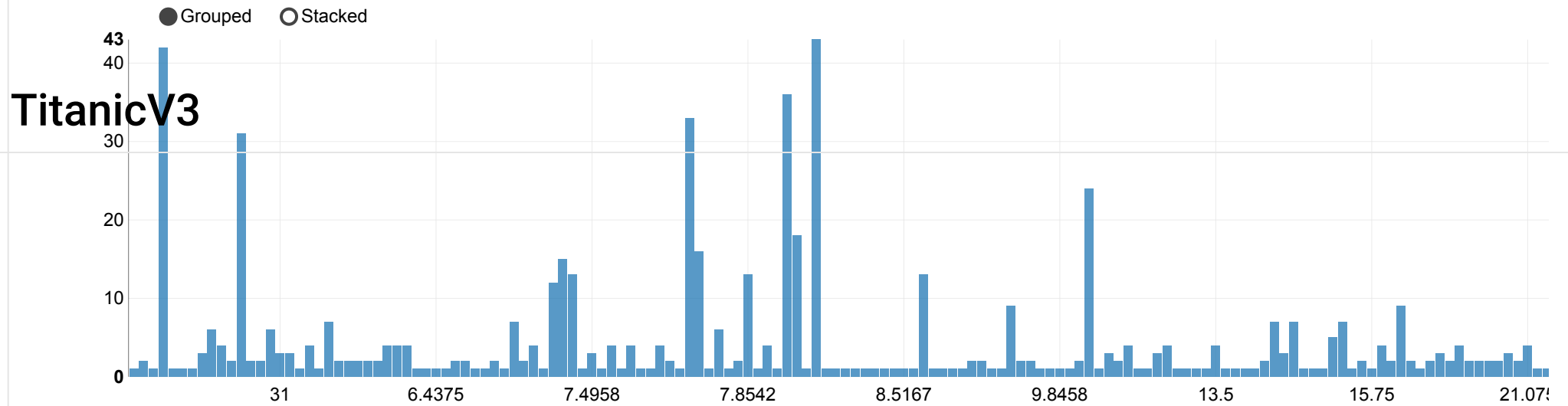


Titanic fare price frequency distribution

READY

```
%spark.sql
select fare, count(1) nInBin
from SMALLPASSTABLE
where fare>0
group by fare
order by fare
```





Basic statistics about Titanic fare price (in year 1912 GBP)

READY

```
%spark.sql  
select * from BASICFARETABLE
```

settings ▼

Minimum fare price ▼	Maximum fare price ▼	Average fare price ≡
4.0125	512.3292	32.86113275737937

TitanicV3

Basic statistics about titanic fare price (in GBP actualised as of 2020)

FINISHED

```
%spark.sql  
select * from BASICUPDFARETABLE
```





settings ▼

Minimum fare price ▼	Maximum fare price ▼	Average fare price ≡
464.647500000000004	59327.72136	3805.319164747697

Took 0 sec. Last updated by anonymous at February 21 2020, 6:05:00 PM. (outdated)

```
%spark.sql
```

READY