| Structured Streaming | Advanced Analytics | Ecosystem |
| --- | --- | --- |

**Structured APIs**

| Datasets | DataFrames | SQL |
| --- | --- | --- |

**Low level APIs**

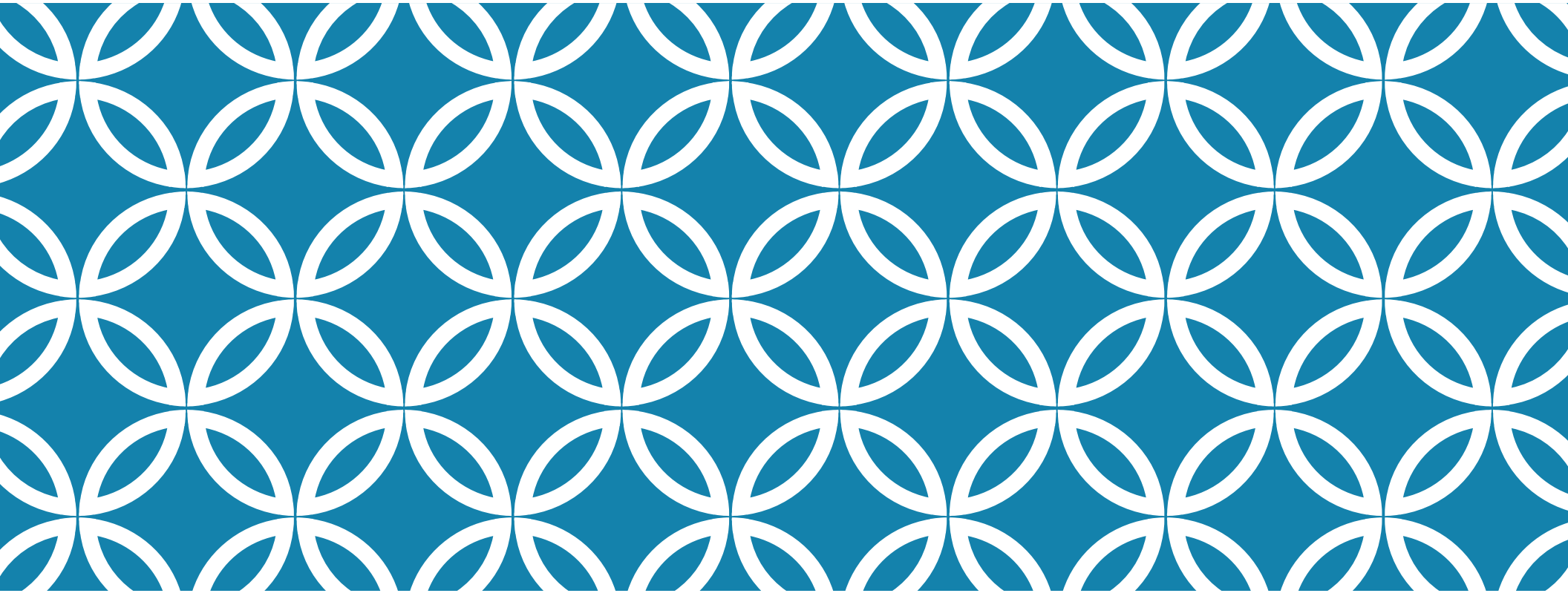| Distributed Variables | | RDDs |
| --- | --- | --- |

# SQL

# DATABASES, VIEWS, AND TABLES

# DATABASES, VIEWS, AND TABLES

Besides the DataFrame API, spark implements the concepts of databases, functions, views and tables. The idea is to enable us to use standard SQL to process data

Spark stores metadata concerning those elements in an internal catalog, which is accessible from a SparkSession

```
[100]:  spark.catalog.listTables()

[100]:  [Table(nam    f  listColumns        function  ^  id
                      f  listDatabases      function
 [21]:  spark.sql(   f  listFunctions      function
                      f  listTables         function
 [21]:  DataFrame[   f  recoverPartitions  function
                      f  refreshByPath      function
 [23]:  usersDf.wr   f  refreshTable       function
                      f  registerFunction   function
 [25]:  spark.cata   i  setCurrentDatabase instance
                      f  uncacheTable       function  v  pt
 [25]:  [Table(nam
                Table(name= pessoas , database=None, descriptio
```

# CREATING A VIEW

```
# Para carregar um DataFrame diretamente de uma fonte de dados externa, pode-se usar:
df = spark.sql("SELECT * FROM csv.`users.csv`",)
df.printSchema()
df.show()
```

```
root
 |-- _c0: string (nullable = true)
 |-- _c1: string (nullable = true)
 |-- _c2: string (nullable = true)

+---+------------------+---+
|_c0|               _c1|_c2|
+---+------------------+---+
|  1|    'Fabio Nogueira'| 47|
|  2|'Andrea Vasconcelos'| 47|
|  3|'Thiago Vasconcel...| 21|
+---+------------------+---+
```

```
# Nao ha ainda tabelas, uma vez que trouxemos direto o conteudo do arquivo para um dataframe
catalog.listTables()
```

```
[]
```

```
# Criando uma tabela temporaria a partir de um dataframe
df.createOrReplaceTempView('pessoas')
```

```
catalog.listTables()
```

```
[Table(name='pessoas', database=None, description=None, tableType='TEMPORARY', isTemporary=True)]
```

# QUERYING VIEWS USING SQL

```
# Creating a temporary table to enable us to submit sql directly
users.createOrReplaceTempView('pessoas')
```

```
pes = spark.sql('select * from pessoas')
pes.show()
```

```
+---+--------------------+---+
|_c0|                 _c1|_c2|
+---+--------------------+---+
|  1|     'Fabio Nogueira'| 47|
|  2|'Andrea Vasconcelos'| 47|
|  3|'Thiago Vasconcel...| 18|
+---+--------------------+---+
```

# TEMPORARY VIEWS

Temporary views in Spark SQL are session-scoped and will disappear if the session that creates it terminates

If you want to have a temporary view that is shared among all sessions and keep alive until the Spark application terminates, you can create a global temporary view

# TEMPORARY VIEWS

```
# Register the DataFrame as a global temporary view
df.createGlobalTempView("people")

# Global temporary view is tied to a system preserved database `global_temp`
spark.sql("SELECT * FROM global_temp.people").show()
# +----+-------+
# | age|   name|
# +----+-------+
# |null|Michael|
# |  30|   Andy|
# |  19| Justin|
# +----+-------+

# Global temporary view is cross-session
spark.newSession().sql("SELECT * FROM global_temp.people").show()
# +----+-------+
# | age|   name|
# +----+-------+
# |null|Michael|
# |  30|   Andy|
# |  19| Justin|
# +----+-------+
```

# PERSISTING DATAFRAMES

DataFrames can also be saved as persistent tables into Hive metastore using the saveAsTable command

These tables can be organized in Databases

Unlike the createOrReplaceTempView command, saveAsTable will materialize the contents of the DataFrame and create a pointer to the data in the Hive metastore.

A DataFrame for a persistent table can be created by calling the table method on a SparkSession with the name of the table.

# PERSISTING DATAFRAMES

For file-based data source, e.g. text, parquet, json, etc. you can specify a custom table path via the path option, e.g. df.write.option("path", "/some/path").saveAsTable("t"). Table persisted that way are referred to as "External", this means that when the table is dropped, the custom table path will not be removed and the table data is still there. If no custom table path is specified, we say the table is "Managed" by Spark. So, it will write data to a default table path under the warehouse directory. When the table is dropped, the default table path will be removed too.

Persistent datasource tables have per-partition metadata stored in the Hive metastore.

```
# 1|24|M|technician|85711
usersDf = spark.read.csv(path='/home/jovyan/work/ml-100k/u.user',sep='|',schema='id int, age int, genre string, occupation string, time int')
usersDf.show(5)
usersDf.printSchema()
```

```
+---+---+-----+----------+-----+
| id|age|genre|occupation| time|
+---+---+-----+----------+-----+
|  1| 24|    M|technician|85711|
|  2| 53|    F|     other|94043|
|  3| 23|    M|    writer|32067|
|  4| 24|    M|technician|43537|
|  5| 33|    F|     other|15213|
+---+---+-----+----------+-----+
only showing top 5 rows

root
 |-- id: integer (nullable = true)
 |-- age: integer (nullable = true)
 |-- genre: string (nullable = true)
 |-- occupation: string (nullable = true)
 |-- time: integer (nullable = true)
```

```
spark.sql("create database ml100k")
```



DataFrame[]

```
spark.catalog.listTables()
```

```
[Table(name='pessoas', database=None, description=None, tableType='TEMPORARY', isTemporary=True)]
```

```
spark.sql('DROP TABLE IF EXISTS ml100k.users')
```
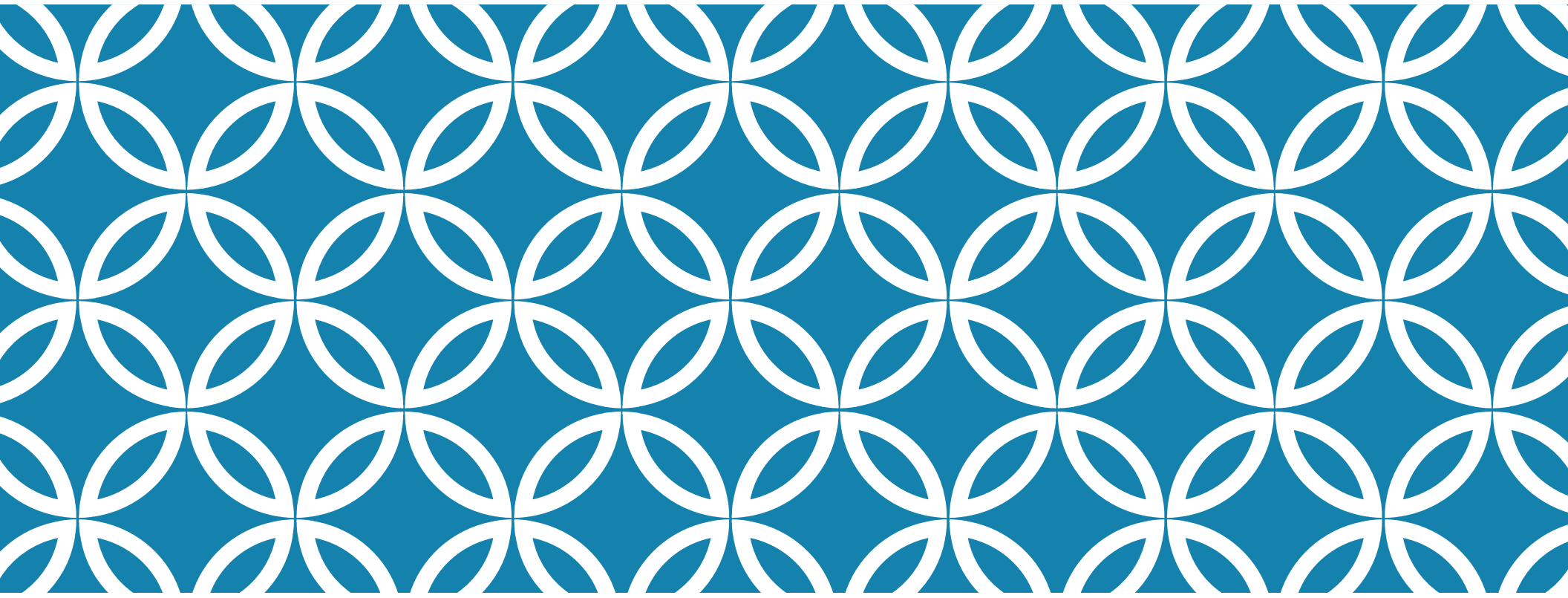
DataFrame[]

```
usersDf.write.saveAsTable('ml100k.users')
```

```
spark.catalog.listTables('ml100k')
```

```
[Table(name='users', database='ml100k', description=None, tableType='MANAGED', isTemporary=False),
 Table(name='pessoas', database=None, description=None, tableType='TEMPORARY', isTemporary=True)]
```
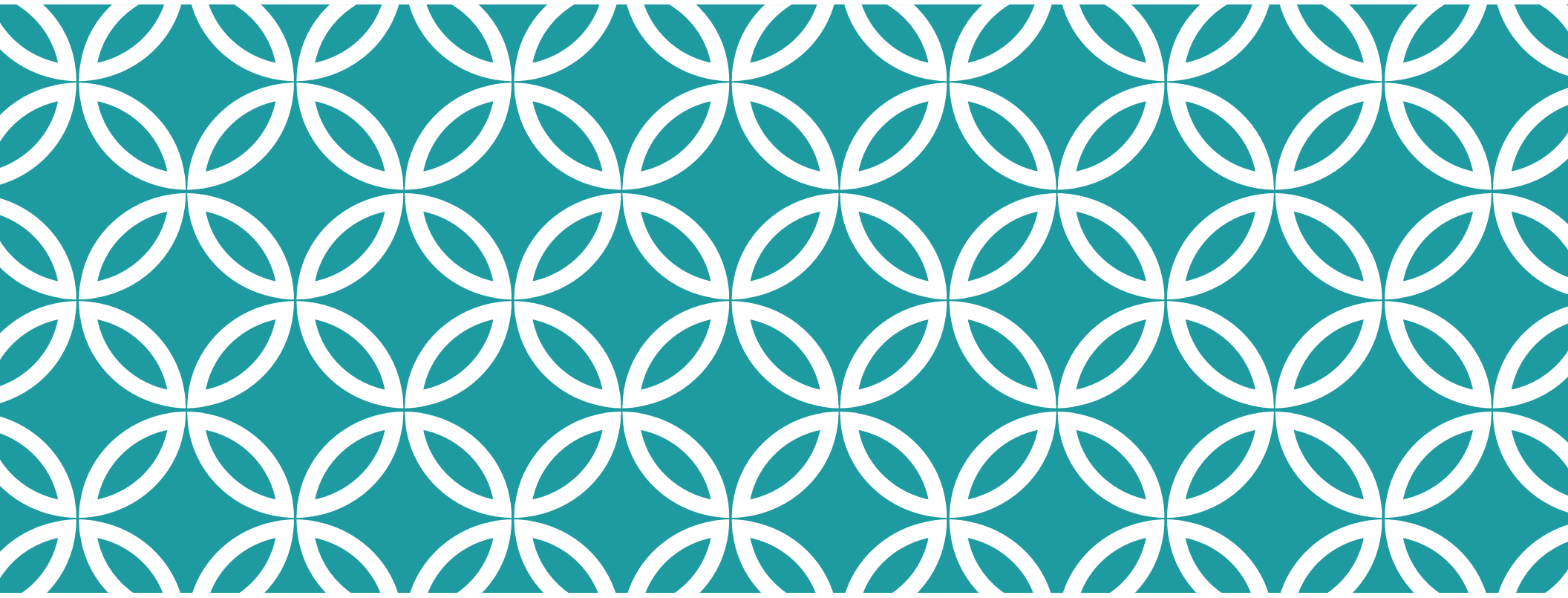
# FUNCTIONS

# FUNCTIONS

Structured APIs includes lots of built-in functions that can be applied to columns according to their types

We will show some of them, others can be found here:

- https://spark.apache.org/docs/latest/api/python/pyspark.sql.html#module-pyspark.sql.functions

Moreover, Spark allows users to define their own functions, which can be registered and used in their queries

# NUMERICAL

# GENERATING IDS

```
ratingsDf.printSchema()

root
 |-- userid: integer (nullable = true)
 |-- itemid: integer (nullable = true)
 |-- rating: integer (nullable = true)
 |-- time: integer (nullable = true)
```

```
ratingsWithKey = ratingsDf.select(monotonically_increasing_id().alias('id'),'*')
ratingsWithKey.printSchema()

root
 |-- id: long (nullable = false)
 |-- userid: integer (nullable = true)
 |-- itemid: integer (nullable = true)
 |-- rating: integer (nullable = true)
 |-- time: integer (nullable = true)
```

```
ratingsWithKey.show(5)

+---+------+------+------+---------+
| id|userid|itemid|rating|     time|
+---+------+------+------+---------+
|  0|   196|   242|     3|881250949|
|  1|   186|   302|     3|891717742|
|  2|    22|   377|     1|878887116|
|  3|   244|    51|     2|880606923|
|  4|   166|   346|     1|886397596|
+---+------+------+------+---------+
only showing top 5 rows
```

# POW

```
1  from pyspark.sql.functions import expr, pow
2  fabricatedQuantity = pow(col("Quantity") * col("UnitPrice"), 2) + 5
3  df.select(expr("CustomerId"), fabricatedQuantity.alias("realQuantity")).show(2)
```

```
1  df.selectExpr(
2    "CustomerId",
3    "(POWER((Quantity * UnitPrice), 2.0) + 5) as realQuantity").show(2)
```

# ROUND AND BROUND

```
1  from pyspark.sql.functions import *
2  df.select(round(col("UnitPrice"), 1).alias("rounded"), col("UnitPrice")).show(5)
```
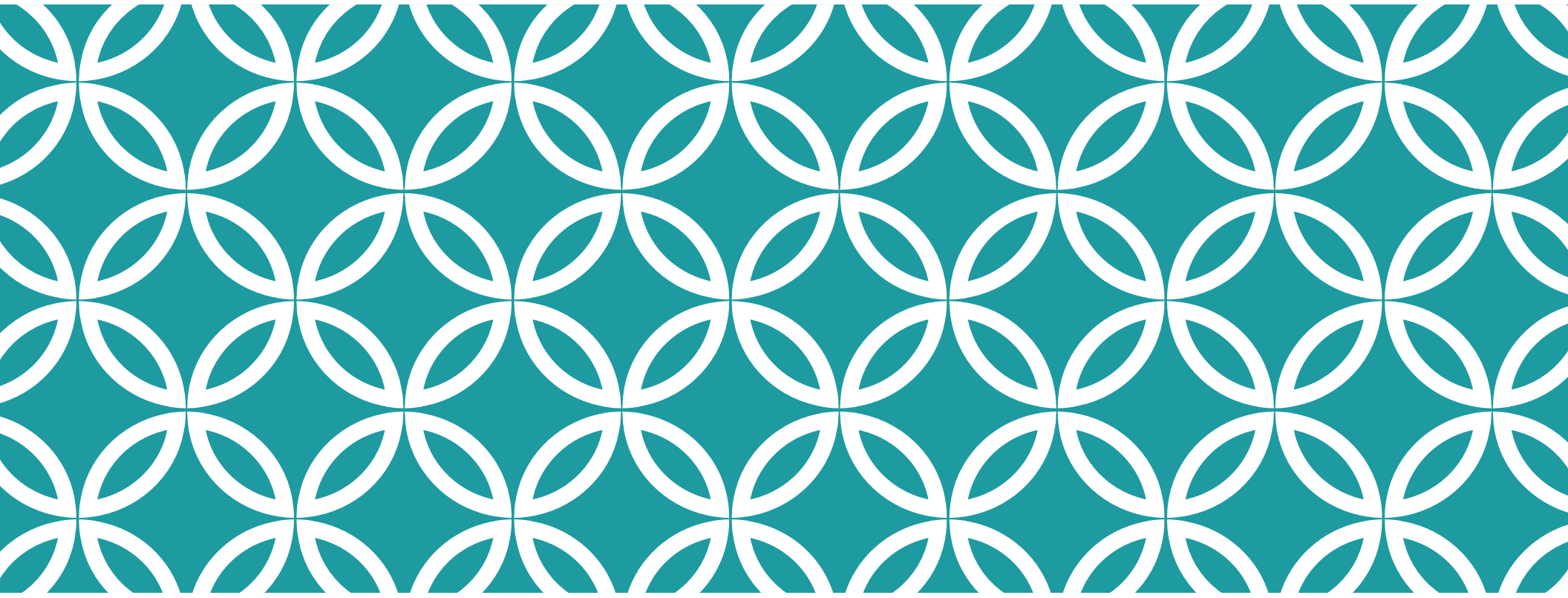
# STATISTICS

```
1  df.describe().show()
2  df.describe('UnitPrice').show()
```

```
count, mean, stddev_pop, min, max
```

# STRINGS

# CASE

To modify a string column in order to change the case of its contents, we can use *initcap, lower,* and *upper* functions

```
1  from pyspark.sql.functions import col, initcap, lower, upper
2  df.select(col("Description"), \
3            initcap(col("Description")), \
4            lower(col("Description")), \
5            upper(col("Description"))).show()
```

```
1  +--------------------+--------------------+--------------------+-----------------
2  |         Description|initcap(Description)|  lower(Description)|  upper(Description
3  +--------------------+--------------------+--------------------+-----------------
4  |WHITE HANGING HEA...|White Hanging Hea...|white hanging hea...|WHITE HANGING HEA..
5  | WHITE METAL LANTERN| White Metal Lantern| white metal lantern| WHITE METAL LANTER
6  |CREAM CUPID HEART...|Cream Cupid Heart...|cream cupid heart...|CREAM CUPID HEART..
7  |KNITTED UNION FLA...|Knitted Union Fla...|knitted union fla...|KNITTED UNION FLA..
8  |RED WOOLLY HOTTIE...|Red Woolly Hottie...|red woolly hottie...|RED WOOLLY HOTTIE..
```

# DEALING WITH SPACES

```
from pyspark.sql.functions import lit, ltrim, rtrim, rpad, lpad, trim
df.select(
    ltrim(lit("    HELLO    ")).alias("ltrim"),
    rtrim(lit("    HELLO    ")).alias("rtrim"),
    trim(lit("   HELLO    ")).alias("trim"),
    lpad(lit("HELLO"), 3, "#").alias("lp"),
    rpad(lit("HELLO"), 10, "#").alias("rp")).show()
```

```
+---------+---------+-----+---+----------+
|    ltrim|    rtrim| trim| lp|        rp|
+---------+---------+-----+---+----------+
|HELLO    |    HELLO|HELLO|HEL|HELLO#####|
|HELLO    |    HELLO|HELLO|HEL|HELLO#####|
|HELLO    |    HELLO|HELLO|HEL|HELLO#####|
+---------+---------+-----+---+----------+
```

But, wait!!! Why do I see 3 lines here???

```
df.show()

+---+-------------------+---+
|_c0|                _c1|_c2|
+---+-------------------+---+
|  1|     'Fabio Nogueira'| 47|
|  2|'Andrea Vasconcelos'| 47|
|  3|'Thiago Vasconcel...| 25|
+---+-------------------+---+
```

# REPLACING STRINGS

```python
from pyspark.sql.functions import regexp_replace, col

regex1 = 'BLACK|WHITE|RED|GREEN|BLUE'
regex2 = r'\bBLACK\b|\bWHITE\b|\bRED\b|\bGREEN\b|\bBLUE\b'
regex3 = r'((?i)\bBLACK\b)|\bWHITE\b|\bRED\b|\bGREEN\b|\bBLUE\b'

mList = [('BLUEBLUE',),('REDO it',),('BLACK Black line',)]

memoryDf = spark.createDataFrame(mList, ['Desc',])

replacedDf = memoryDf.select('Desc',
                             regexp_replace(col('Desc'),regex1,'COLOR').alias('Pater
                             regexp_replace(col('Desc'),regex2,'COLOR').alias('Patte
                             regexp_replace(col('Desc'),regex3,'COLOR').alias('Patte
                             )
replacedDf.show(truncate=False)
```

```
+----------------+---------------+---------------+---------------+
|Desc            |Patern 1       |Pattern 2      |Pattern 3      |
+----------------+---------------+---------------+---------------+
|BLUEBLUE        |COLORCOLOR     |BLUEBLUE       |BLUEBLUE       |
|REDO it         |COLORO it      |REDO it        |REDO it        |
|BLACK Black line|COLOR Black line|COLOR Black line|COLOR COLOR line|
+----------------+---------------+---------------+---------------+
```
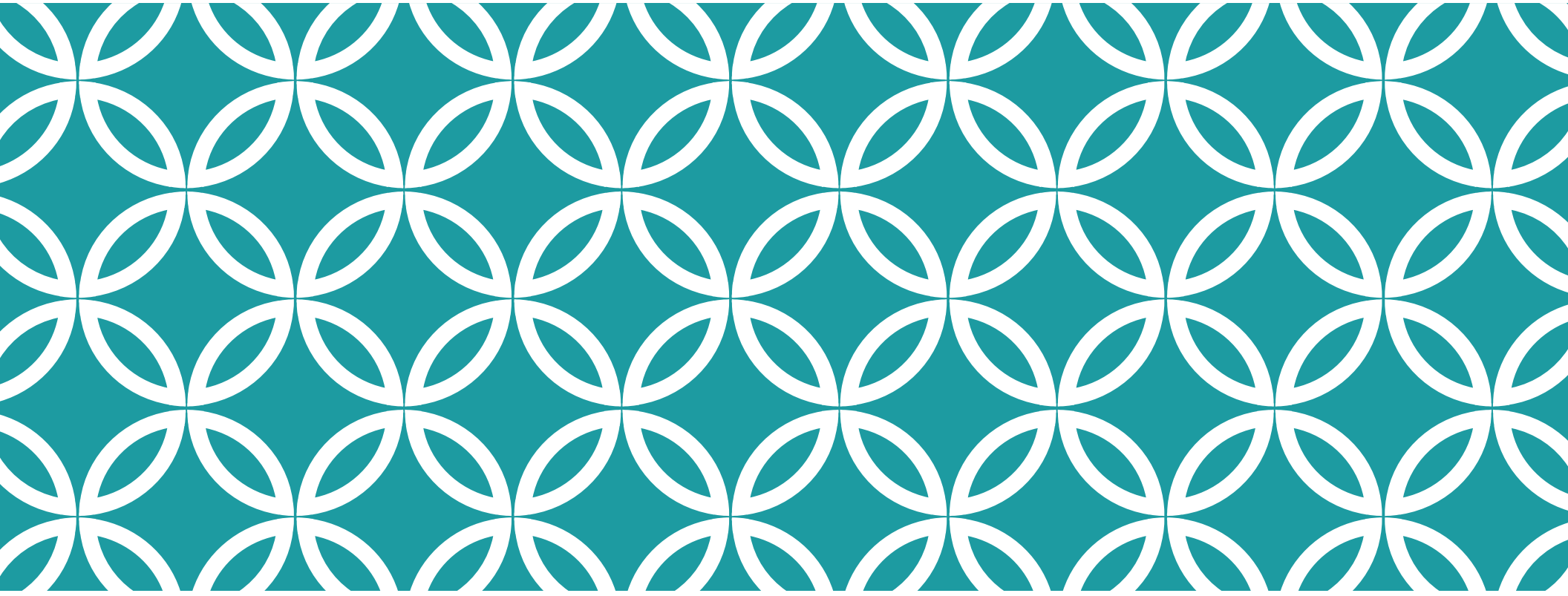
# MODIFYING INDIVIDUAL CHARACTERS

```
1  from pyspark.sql.functions import translate
2  df.select(translate(col("Description"), "LEET", "1337"),col("Description"))\
3    .show(2)
```

```
1  +-------------------------------+--------------------+
2  |translate(Description, LEET, 1337)|         Description|
3  +-------------------------------+--------------------+
4  |          WHI73 HANGING H3A...|WHITE HANGING HEA...|
5  |          WHI73 M37A1 1AN73RN| WHITE METAL LANTERN|
6  +-------------------------------+--------------------+
```

# VERIFYING IF A STRING IS IN A COLUMN

```
1  from pyspark.sql.functions import instr
2  containsBlack = instr(col("Description"), "BLACK") >= 1
3  containsWhite = instr(col("Description"), "WHITE") >= 1
4  df.withColumn("hasSimpleColor", containsBlack | containsWhite)\
5    .where("hasSimpleColor")\
6    .select("Description").show(3, False)
```

```
1  +----------------------------------+
2  |Description                       |
3  +----------------------------------+
4  |WHITE HANGING HEART T-LIGHT HOLDER|
5  |WHITE METAL LANTERN               |
6  |RED WOOLLY HOTTIE WHITE HEART.    |
7  +----------------------------------+
```

# DATES AND TIME

# CURRENT DATE AND TIMESTAMP

```python
from pyspark.sql.functions import current_date, current_timestamp
# spark.range(start, end=None, step=1, numPartitions=None)
# Create a DataFrame with single pyspark.sql.types.LongType column named id, contai

dateDF = spark.range(10)\
  .withColumn("today", current_date())\
  .withColumn("now", current_timestamp())
dateDF.createOrReplaceTempView("dateTable")
dateDF.show(truncate=False)
dateDF.printSchema()
```

exit: Ctrl+⏎

```
+---+----------+-----------------------+
|id |today     |now                    |
+---+----------+-----------------------+
|0  |2019-10-31|2019-10-31 14:54:29.498|
|1  |2019-10-31|2019-10-31 14:54:29.498|
|2  |2019-10-31|2019-10-31 14:54:29.498|
|3  |2019-10-31|2019-10-31 14:54:29.498|
|4  |2019-10-31|2019-10-31 14:54:29.498|
|5  |2019-10-31|2019-10-31 14:54:29.498|
|6  |2019-10-31|2019-10-31 14:54:29.498|
|7  |2019-10-31|2019-10-31 14:54:29.498|
|8  |2019-10-31|2019-10-31 14:54:29.498|
|9  |2019-10-31|2019-10-31 14:54:29.498|
+---+----------+-----------------------+
```

# ADD AND SUB DATES

```
1  from pyspark.sql.functions import col, date_add, date_sub
2  dateDF.select('id', 'today', date_sub(col("today"), 5), date_add(col("today"), 5)).
```

```
1  +---+----------+-----------------+-----------------+
2  | id|     today|date_sub(today, 5)|date_add(today, 5)|
3  +---+----------+-----------------+-----------------+
4  |  0|2019-10-31|       2019-10-26|       2019-11-05|
5  +---+----------+-----------------+-----------------+
```

```
1  from pyspark.sql.functions import datediff, months_between, to_date, lit
2  dateDF.withColumn("week_ago", date_sub(col("today"), 7)) \
3    .withColumn('next_week', date_add(col('today'),7)) \
4    .select(datediff(col("week_ago"), col("today")).alias('DaysFromWeekAgo'), \
5          datediff(col('next_week'), col('today')).alias('DaysFromNextWeek')).show(
6
7  dateDF.select(
8      to_date(lit("2016-01-01")).alias("start"),
9      to_date(lit("2017-05-22")).alias("end"))\
10   .select(months_between(col("start"), col("end"))).show(1)
```

```
1  +---------------+----------------+
2  |DaysFromWeekAgo|DaysFromNextWeek|
3  +---------------+----------------+
4  |            -7|               7|
5  +---------------+----------------+
6  only showing top 1 row
7
8  +-------------------------------+
9  |months_between(start, end, true)|
10 +-------------------------------+
11 |                   -16.67741935|
12 +-------------------------------+
13 only showing top 1 row
```
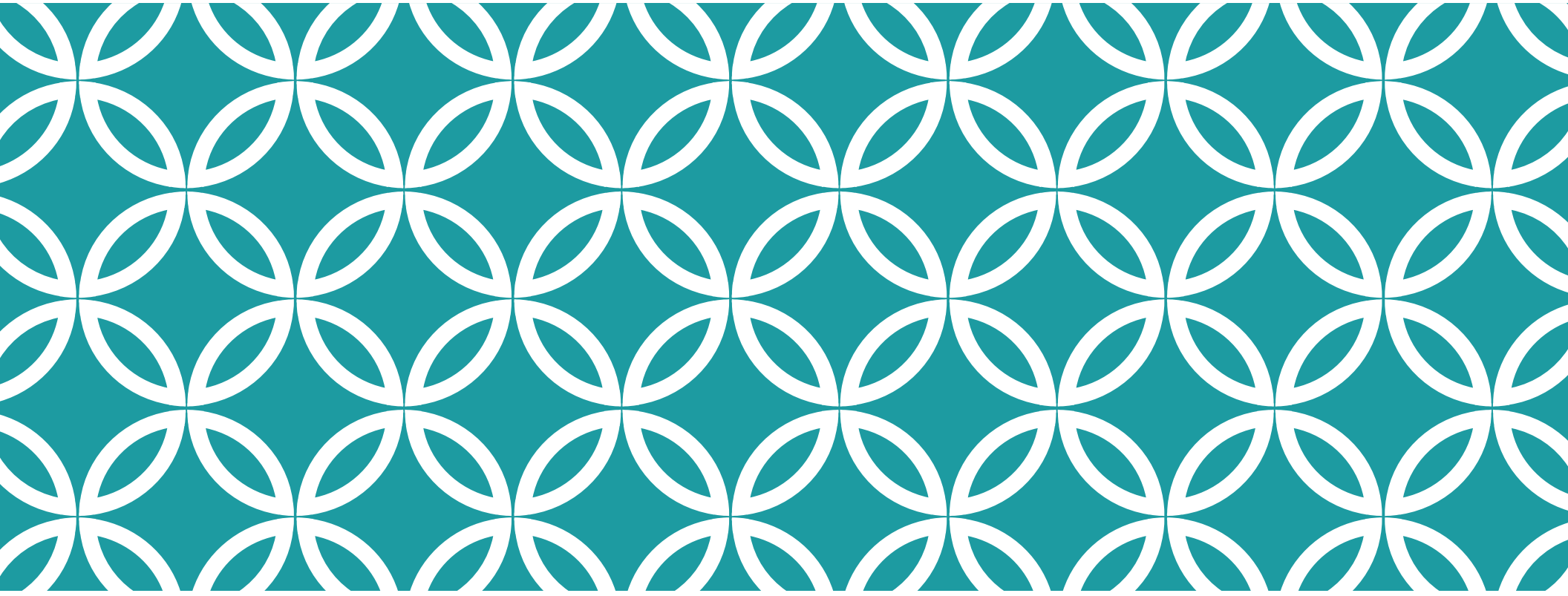
# CREATING DATES

We can convert a string or timestamp column to a *date* using *to_date* function. Besides the column, it can receive a format, specified according to *SimpleDateFormats* ( https://docs.oracle.com/javase/tutorial/i18n/format/simpleDateFormat.html). If a value in the string column can not be recognized as a *date*, Spark will just return *null* as the corresponding value.

```
1  from pyspark.sql.functions import to_date
2  dateFormat = "yyyy-dd-MM"
3  cleanDateDF = spark.range(1).select(
4      to_date(lit("2017-12-11"), dateFormat).alias("date"),
5      to_date(lit("2017-20-12"), dateFormat).alias("date2")).show()
```

```
1  from pyspark.sql.functions import to_timestamp
2  cleanDateDF.select(to_timestamp(col("date"), dateFormat)).show()
```

# WORKING WITH NULL

# COALESCE

```
1   cDf = spark.createDataFrame([(None, None), (1, None), (None, 2)], ("a", "b"))
2   cDf.select('*',
3       lit('0.0').alias('lit(0.0)'),
4       coalesce(cDf["a"], cDf["b"]),
5       coalesce(cDf["a"], lit(0.0)),
6       coalesce(lit('0.0'), cDf['a'])
7   ).show()
8   +----+----+--------+-------------+---------------+---------------+
9   |   a|   b|lit(0.0)|coalesce(a, b)|coalesce(a, 0.0)|coalesce(0.0, a)|
10  +----+----+--------+-------------+---------------+---------------+
11  |null|null|     0.0|         null|            0.0|            0.0|
12  |   1|null|     0.0|            1|            1.0|            0.0|
13  |null|   2|     0.0|            2|            0.0|            0.0|
14  +----+----+--------+-------------+---------------+---------------+
```

# DROP

```
1  cDf = spark.createDataFrame([(None, None), (1, None), (None, 2), (1,1)], ("a", "b")
2  cDf.show()
3  +----+----+
4  |   a|   b|
5  +----+----+
6  |null|null|
7  |   1|null|
8  |null|   2|
9  |   1|   1|
10 +----+----+
11 cDf.dropna().show()
12 +---+---+
13 |  a|  b|
14 +---+---+
15 |  1|  1|
16 +---+---+
17 cDf.na.drop().show()
18 +---+---+
19 |  a|  b|
20 +---+---+
21 |  1|  1|
22 +---+---+
23 cDf.na.drop('any').show()
24 +---+---+
25 |  a|  b|
26 +---+---+
27 |  1|  1|
28 +---+---+
```

# DROP (CONT)

```
29
30  cDf.na.drop('all').show()
31  +----+----+
32  |   a|   b|
33  +----+----+
34  |   1|null|
35  |null|   2|
36  |   1|   1|
37  +----+----+
38  cDf.na.drop(subset=['a']).show()
39  +---+----+
40  |  a|   b|
41  +---+----+
42  |  1|null|
43  |  1|   1|
44  +---+----+
45  cDf.na.drop(subset=['a','b']).show()
46  +---+---+
47  |  a|  b|
48  +---+---+
49  |  1|  1|
50  +---+---+
```

# FILL

```
1  df4.show()
2  +----+------+-----+
3  | age|height| name|
4  +----+------+-----+
5  |  10|    80|Alice|
6  |   5|  null|  Bob|
7  |null|  null|  Tom|
8  |null|  null| null|
9  +----+------+-----+
```

```
1  df4.na.fill(50).show()
2  +---+------+-----+
3  |age|height| name|
4  +---+------+-----+
5  | 10|    80|Alice|
6  |  5|    50|  Bob|
7  | 50|    50|  Tom|
8  | 50|    50| null|
9  +---+------+-----+
```

# FILL (CONT)

```
1  df4.na.fill(50, subset=['age']).show()
2  +---+------+-----+
3  |age|height| name|
4  +---+------+-----+
5  | 10|    80|Alice|
6  |  5|  null|  Bob|
7  | 50|  null|  Tom|
8  | 50|  null| null|
9  +---+------+-----+
```

```
1  df4.na.fill({'age': 50, 'name': 'unknown'}).show()
2  +---+------+-------+
3  |age|height|   name|
4  +---+------+-------+
5  | 10|    80|  Alice|
6  |  5|  null|    Bob|
7  | 50|  null|    Tom|
8  | 50|  null|unknown|
9  +---+------+-------+
```

exit: Ctrl+↵

# REPLACE

```
df4.show()
+----+------+-----+
| age|height| name|
+----+------+-----+
|  10|    80|Alice|
|   5|  null|  Bob|
|null|  null|  Tom|
|null|  null| null|
+----+------+-----+
```

```
df4.na.replace(10, 20).show()
+----+------+-----+
| age|height| name|
+----+------+-----+
|  20|    80|Alice|
|   5|  null|  Bob|
|null|  null|  Tom|
|null|  null| null|
+----+------+-----+
```

# REPLACE (CONT)

```
1  df4.na.replace('Alice', None).show()
2  +----+------+----+
3  | age|height|name|
4  +----+------+----+
5  |  10|    80|null|
6  |   5|  null| Bob|
7  |null|  null| Tom|
8  |null|  null|null|
9  +----+------+----+
```
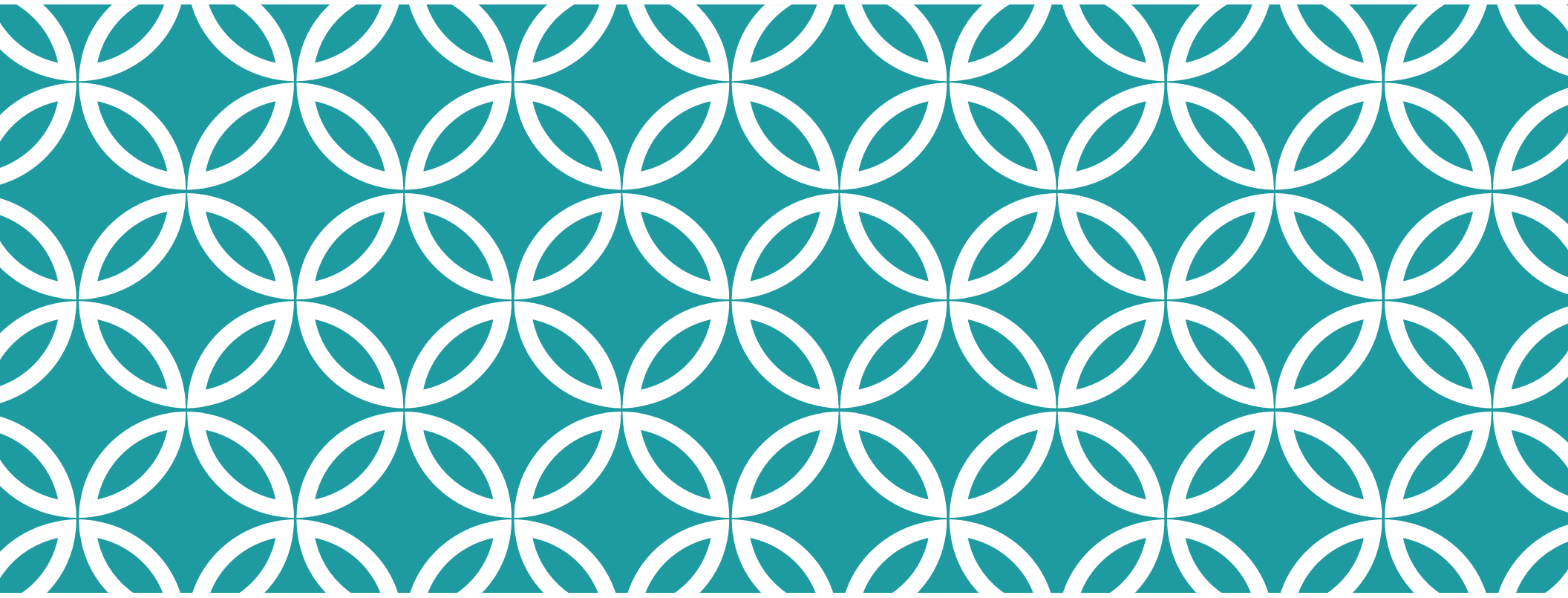
```
1  df4.na.replace({10 : 100, 5 : 50, 80 : 10},subset=['age']).show()
2  +----+------+-----+
3  | age|height| name|
4  +----+------+-----+
5  | 100|    80|Alice|
6  |  50|  null|  Bob|
7  |null|  null|  Tom|
8  |null|  null| null|
9  +----+------+-----+
```

COMPLEX TYPES

# STRUCTS

```
1   from pyspark.sql.functions import struct
2   complexDF = df.select(struct("Description", "InvoiceNo").alias("complex"))
3   complexDF.createOrReplaceTempView("complexDF")
4
5   complexDF.select("complex.Description")
6   complexDF.select(col("complex").getField("Description"))
7   complexDF.select("complex.*")
```
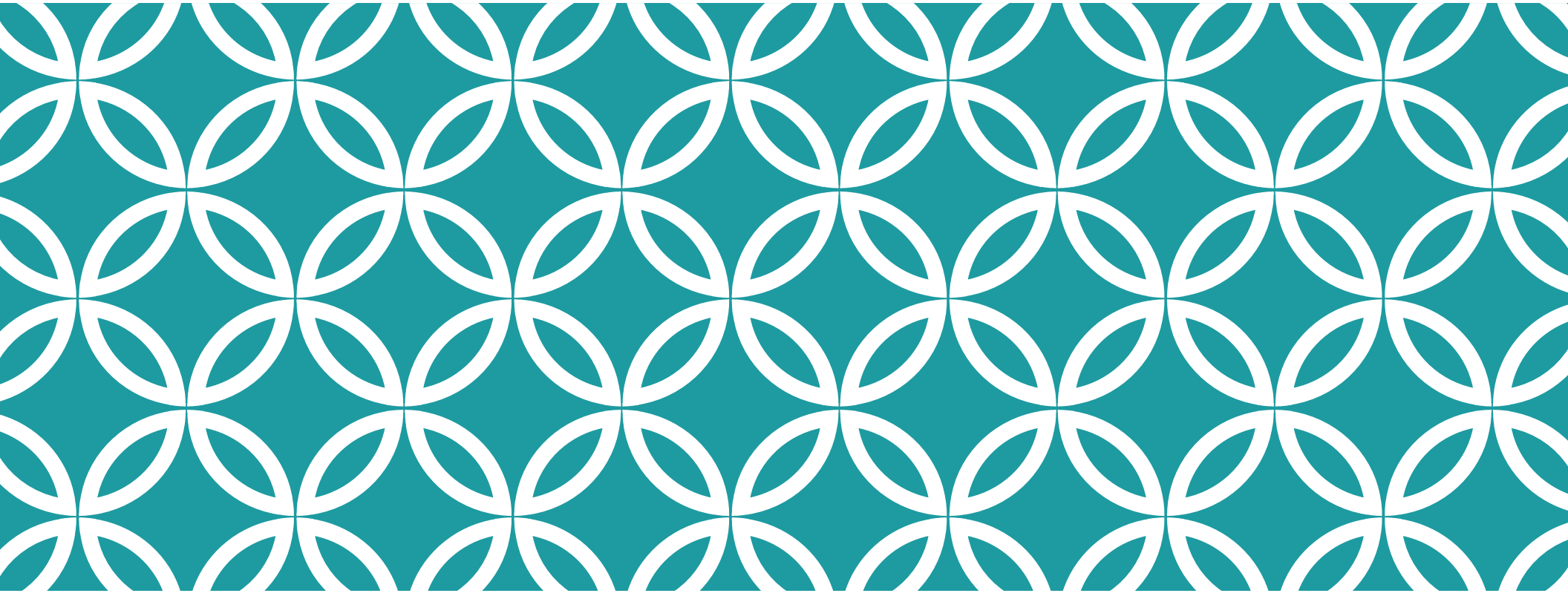
# ARRAYS

```
1  from pyspark.sql.functions import split, size, array_contains
2
3  df.select('Description',
4     split(col("Description"), " ").alias("asArray"),
5     array_contains(split(col("Description"), " "), "WHITE"),
6     size(split(col("Description"), " ")).alias('wordsCounter')
7  ).show(2, truncate=False)
8  +--------------------------------+-----------------------------------------+------
9  |Description                     |asArray                                  |array_
10 +--------------------------------+-----------------------------------------+------
11 |WHITE HANGING HEART T-LIGHT HOLDER|[WHITE, HANGING, HEART, T-LIGHT, HOLDER]|true
12 |WHITE METAL LANTERN             |[WHITE, METAL, LANTERN]                  |true
13 +--------------------------------+-----------------------------------------+------
```

# ARRAYS

We can explode an array column to generate a line per element:

```
from pyspark.sql.functions import split, explode

df.withColumn("splitted", split(col("Description"), " "))\
    .withColumn("exploded", explode(col("splitted"))) \
    .withColumn('size',size("splitted")) \
    .select("Description", "InvoiceNo", "exploded").show(5)
+--------------------+---------+--------+
|         Description|InvoiceNo|exploded|
+--------------------+---------+--------+
|WHITE HANGING HEA...|   536365|   WHITE|
|WHITE HANGING HEA...|   536365| HANGING|
|WHITE HANGING HEA...|   536365|   HEART|
|WHITE HANGING HEA...|   536365| T-LIGHT|
|WHITE HANGING HEA...|   536365|  HOLDER|
+--------------------+---------+--------+
```

# AGGREGATION

# COUNT

```
from pyspark.sql.functions import count

df4.show()
+----+------+-----+
| age|height| name|
+----+------+-----+
|  10|    80|Alice|
|   5|  null|  Bob|
|null|  null|  Tom|
|null|  null| null|
+----+------+-----+
df4.select(count('*'),
           count("age"),
           count("height"),
           count('name')).show()
+--------+----------+-------------+-----------+
|count(1)|count(age)|count(height)|count(name)|
+--------+----------+-------------+-----------+
|       4|         2|            1|          3|
+--------+----------+-------------+-----------+
```

# COUNT

```
1   from pyspark.sql.functions import countDistinct
2
3   dfx = spark.sparkContext.parallelize([(2, 'Alice'), (2, 'Bob')])\
4          .toDF(StructType([StructField('age', IntegerType()),
5                            StructField('name', StringType())]))
6
7   dfx.show()
8   +---+-----+
9   |age| name|
10  +---+-----+
11  |  2|Alice|
12  |  2|  Bob|
13  +---+-----+
14  dfx.agg(countDistinct(dfx.age).alias('c')).show()
15  +---+
16  |  c|
17  +---+
18  |  1|
19  +---+
20  dfx.agg(countDistinct(dfx.name).alias('c')).show()
21  +---+
22  |  c|
23  +---+
24  |  2|
25  +---+
26  dfx.agg(countDistinct(dfx.age, dfx.name).alias('c')).show()
27  +---+
28  |  c|
29  +---+
30  |  2|
31  +---+
```

# FIRST AND LAST

```
: df.show(truncate=False)

+---+----------------------------+---+
|uid|name                        |age|
+---+----------------------------+---+
|1  |'Fabio Nogueira'            |47 |
|2  |'Andrea Vasconcelos'        |47 |
|3  |'Thiago Vasconcelos Nogueira'|25 |
+---+----------------------------+---+
```

```
: from pyspark.sql.functions import first, last
  df.select(first('uid'), last('age')).show()

+----------------+----------------+
|first(uid, false)|last(age, false)|
+----------------+----------------+
|               1|              25|
+----------------+----------------+
```

# SUM, MIN, AND MAX

```
1   from pyspark.sql.functions import sum, min, max
2   dfx.show()
3   +----+-----+
4   | age| name|
5   +----+-----+
6   |   2|Alice|
7   |   5|  Bob|
8   |null| Fool|
9   |null| null|
10  +----+-----+
11  dfx.select(sum("age"),
12      min("age"),
13      max("age")).show()
14  +--------+--------+--------+
15  |sum(age)|min(age)|max(age)|
16  +--------+--------+--------+
17  |       7|       2|       5|
18  +--------+--------+--------+
```
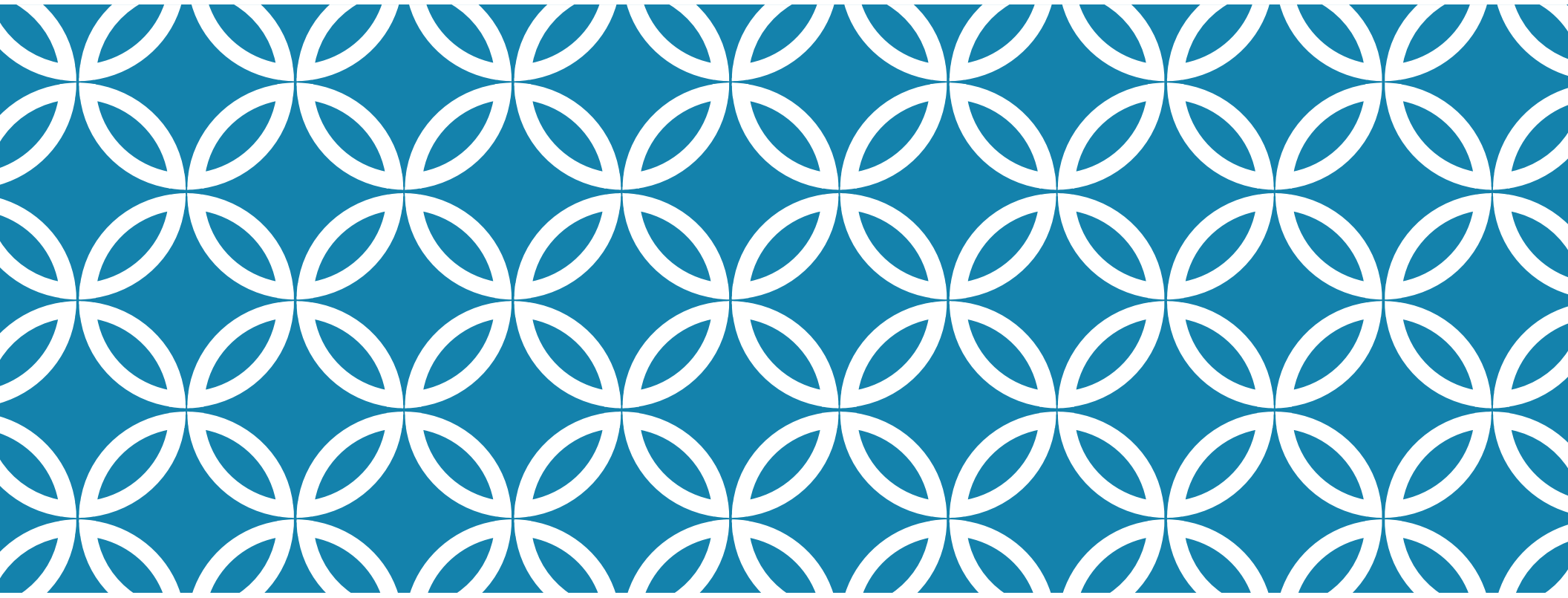
# AVG, VAR, AND STDEV

```
1  dfx.select(
2      count("age"),
3      sum("age"),
4      sum('age') / count('age') ,
5      avg("age")).show()
6  +----------+--------+---------------------+--------+
7  |count(age)|sum(age)|(sum(age) / count(age))|avg(age)|
8  +----------+--------+---------------------+--------+
9  |         2|       7|                  3.5|     3.5|
10 +----------+--------+---------------------+--------+
```

```
1  from pyspark.sql.functions import var_pop, stddev_pop
2  from pyspark.sql.functions import var_samp, stddev_samp
3
4  df.select(var_pop("Quantity"), var_samp("Quantity"),
5    stddev_pop("Quantity"), stddev_samp("Quantity")).show()
6  +----------------+-----------------+--------------------+--------------------+
7  |var_pop(Quantity)|var_samp(Quantity)|stddev_pop(Quantity)|stddev_samp(Quantity)|
8  +----------------+-----------------+--------------------+--------------------+
9  |695.2492099104054| 695.4729785650273|  26.367578764657278|  26.371821677029203|
10 +----------------+-----------------+--------------------+--------------------+
```

# GROUP BY

```
1  from pyspark.sql.functions import *
2  usersDf.groupBy('genre').count().\
3      na.replace({'F':'Female','M':'Male'}).show()
4  +------+-----+
5  | genre|count|
6  +------+-----+
7  |Female|  273|
8  |  Male|  670|
9  +------+-----+
```

# QUESTIONS???