

Defining Schema

```
from pyspark.sql.types import *
Schema = StructType([
    StructField('St ore ', StringType(), nullable=True),
    StructField('St ore Typ e', StringType(), nullable=True),
    StructField('As sor tme nt', StringType(), nullable=True),
    StructField('Co mpe tit ion Dis tan ce', FloatType(), nullable=True),
    StructField('Co mpe tit ion Ope nSi nce Mon th', IntegerType(), nullable=True),
    StructField('Co mpe tit ion Ope nSi nce Yea r', IntegerType(), nullable=True),
    StructField('Pr omo 2', IntegerType(), nullable=True),
    StructField('Pr omo 2Si nce Wee k', IntegerType(), nullable=True),
    StructField('Pr omo 2Si nce Yea r', IntegerType(), nullable=True),
    StructField('Pr omo Int erv al', StringType(), nullable=True)
])
df = spark.read.option("header", True).schema(Schema).csv('storage.csv')
# We can drop invalid rows while reading the dataset by setting the read mode as "DROPMALFORMED"
df_1 = spark.read.option("header", True).option("mode", "DROPMALFORMED").csv('storage.csv')
df.show()
```

Spark does not detect schema itself properly, so we need to define the schema as well for the data set.

| PySpark DataTypes | | | |
|-------------------|-------------|---------|----------------|
| Type | Size (Byte) | Default | Range (Digits) |
| byte | 1 | 0 | 3 Ints |
| short | 2 | 0 | 5 |
| int | 4 | 0 | 10 |
| long | 8 | 0 | Lots |
| floats | 4 | 0.0f | Lots floats |
| double | 8 | 0.0d | Lots |
| DecimalType | 32 | 0.0 | Lots |

Filtering Data

String Data Types

| | |
|-------------|---|
| StringType | A variant of StringType which has a length limitation. Data writing will fail if the input string exceeds the length limitation |
| CharType(n) | Reading column of type CharType(n) always returns string values of length n. Char type column comparison will pad the short one to the longer length. |

Adding, renaming and removing columns

Complex Data Types

| | |
|--|---|
| ArrayType(elementType, containsNull) | Represents values comprising a sequence of elements |
| MapType(keyType, valueType, valueContainsNull) | Represents values comprising a set of key-value pairs. The data type of keys is described by keyType and the data type of values is described by valueType. For a MapType value, keys are not allowed to have null values. valueContainsNull is used to indicate if values of a MapType value can have null values. |
| StructType(fields) | Represents values with the structure described by a sequence of StructFields (fields) |

If, elif, else equivalent

```

voter_df.filter(voter_df['name'].isNotNull().withColumn
    OR
    voter_df.withColumn(' - '
    voter_df.here(~ voter_df._c1.isNear(''))).filter(e.year)
    voter_df.filter(voter_df.date.year - withColumn(' - ' +
    1800).test_df.sex = test_df.withColumn(' - ' +
    voter_df.where(voter_df['_c0'].contains('Gender'),
    ntains('VOTE'))).drop
#Multiple Conditions
whereDF = flatten.where((col("voter_df.dropped('united_states') -
    state_name") == "xiangrui") | (col("fumang") -
    tName) == "mic haeli").sort(asc�mapspark.sql import
state_name"))
whereDF.show(truncate=False)
#Unique Values
voter_df = df.select(df["VOTER_NAME"]).distinct()
# Show the rows with 10 highest IDs
set
voter_df.orderBy(voter_df.id.desc().offset(1000),
    df.id.cast(IntegerType())))

```

```

.when(<if condition>, <then xx>)
df.select(df.Name, df.Age,
.when(df.Age >= 18, "Adult")
.when(df.Age < 18, "Minor"))
.otherwise() is like else
df.select(df.Name, df.Age,
.when(df.Age >= 18, "Adult")
.otherwise("Minor"))

```

Remove duplicate rows & replace values

```

dropDuplicates()
test_df.dropDuplicates()
test_df.fillna('Unknown', 'Age').dropDuplicates()
fillna()
used to replace null value with
any other value
df.fillna(value=-99, subset=

```

User Defined Functions

- Define a Python method
- ```
def reverseString(string):
 return string[::-1]
```
- Wrap the function and store as a variable
- ```
udfReverseString = udf(reverseString, StringType())
```
- Use with Spark
- ```
user_df = user_df.withColumn('Reverse Name',
 udfReverseString(u.user - df.Name))
```

### Using SQL to clean script

```
df.createOrReplaceTempView("table1")
df2 = spark.sql("SELECT field1, field2
 FROM table1")
```

### Validating with Joins

```

parsed_df =
spark.read.parquet('parsed_data.parquet')
company_df = spark.read.parquet('company.parquet')
verified_df = parsed_df.join(company_df,
 parsed_df.company == company_df.company)
This automatically removes any rows
with a company not in the valid_df !

```

### View data/actions:

`printSchema(), head(), show(), count(), columns and describe()`

`show()` - Displays/Prints a number of rows in a tabular format. By default it displays 20 rows and to change the default number, you can pass a value to `show(n)`.

`take(n)` returns first n rows as Array of row objects. It is an alias for `first()`.

`count()` - total rows

Published 3rd September, 2022.

Last updated 12th September, 2022.

Page 2 of 3.

Sponsored by [Readable.com](#)

Measure your website readability!

<https://readable.com>



By [datamansam](#)

[cheatography.com/datamansam/](https://cheatography.com/datamansam/)

### Remove duplicate rows & replace values

(cont)

```
> ["Promo2SinceWeek","Promo2SinceYear"]).show()
.withColumn() ,when()
creating a new column, with value equal to
1 if
Promo2SinceYear > 2000 otherwise 0
df.withColumn("greater_than_2000",
when(df.CompetitionDistance==2000,1).otherwise(0)
.alias('value_desc')).show()
```



By **datamansam**

[cheatography.com/datamansam/](https://cheatography.com/datamansam/)

Published 3rd September, 2022.

Last updated 12th September, 2022.

Page 3 of 3.

Sponsored by **Readable.com**

Measure your website readability!

<https://readable.com>