

## Cleaning with PySpark Cheat Sheet

### Defining Schema

```

from pyspark.sql.types import *
Schema = StructType([
    StructField('St ore ', StringType(), nullable=True),
    StructField('St ore Type', StringType(), nullable=True),
    StructField('As sortment', StringType(), nullable=True),
    StructField('Competition Distance', FloatType(), nullable=True),
    StructField('Competition Open_since Month', IntegerType(), nullable=True),
    StructField('Competition Tit ion Ope nSince Year', IntegerType(), nullable=True),
    StructField('Promo 2', IntegerType(), nullable=True),
    StructField('Promo 2Since Week', IntegerType(), nullable=True),
    StructField('Promo 2Since Year', IntegerType(), nullable=True),
    StructField('Promo Int erval', StringType(), nullable=True)
])
df = spark.read.option("header", True).schema(Schema).csv('store.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('store.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 (Digits)	Range
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
Decimal	- 32	0.0	Lots
IType			

### Filtering Data

String Data Types	
StringType	
Varchar	A variant of StringType which has a length limitation. Data writing will fail if the input string exceeds the length limitation
CharType	Reading column of type CharType - pe(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	nts values comprising a sequence of elements
MapType	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.
StructType	Represents values with the structure described by a sequence of Struct Fields

(fields)



```

voter_df.filter(voter_df['name'].isNotNull()).addWithColumn
    ORvoter_df.withColumn("year", voter_df.year)
    voter_df.withColumn("name", voter_df.name)
    voter_df.withColumn("sex", voter_df.sex)
    voter_df.withColumn("age", voter_df.age)

# Multiple Conditions
whereDF = voter_df.where((voter_df['id'] == 1) | (voter_df['id'] == 2))

# Unique Values
voter_df.select("name").distinct().show(10)

User Defined Functions
1. Define a Python method
def reverseString(string):
    return string[::-1]

2. Wrap the function and store as a variable
udfReverseString = udf(reverseString)

3. Use with Spark
user_df = user_df.withColumn("Name", udfReverseString(user_df.Name))

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

```

---

```

.when(<if condition>, <then x>)
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.na.drop = test_df.dropna()
test_df.na.fill = test_df.fillna()
used to replace null value with any other value
df.fillna(value = 99, subset=

```

---

```

functions as F
add_n = udf(lambda x, y: x + y,
IntegerType())
# We register a UDF that adds a column to the DataFrame,
# we cast the id column to an Integer type.
f.registerID =.
offset,
add_n(F.lit(1000),
df.id.cast(IntegerType))

```

---

```

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_id.fromcompany)
# This automatically removes any rows with a company not in the valid_df !

```

---

```

View data/actions:
printSchema(), head(), 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).

```

where as 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.

### Remove duplicate rows & replace values

(cont)

```
> ["Promo2S inc eWeek", "Promo2S inc eYear -  
ar"]].show()  
.withColumn("new_col", when(  
creating a new column, with value equal to  
1 if  
Promo2S <= 2000 otherwise 0  
df.withColumn("new_col", when(promo2S <= 2000,  
when(df.CmpType == "New", 1).otherwise(0))  
alias("new_col").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>