

# PySpark Fingertip Commands Cheat Sheet

## Reading data from a file

```
df = spark.read.csv("file.csv", header=True)
df =
spark.read.option("header", True)
```

## Casting a column to a different data type

```
df.withColumn("col1", col("col1").cast("double"))
```

## Displaying the schema of a Dataframe

```
df.printSchema()
```

## Get distinct count of columns

```
df.select("col").distinct().count()
```

## Filtering rows based on a condition

```
# Filter entries of age, only keep those records
# of
# which the values are >24
df.filter(df["age"] > 24).show()
```

## Renaming columns of DataFrame

```
#Syntax
df.withColumnRenamed("old_name", "new_name")
#Example
df =
df.withColumnRenamed("CallNumber", "Phone Number")
```

## Inspect Data

```
# Return first n rows
df.head()
# Return first row
df.first()
# Return the first n rows
df.take(2)
# Print the schema of df
df.printSchema()
# Print the (logical and physical) plans
df.explain()
#Get All column names from DataFrame
df.columns
```

## Get count

```
# Get_row count
rows = empDF.count()
# Get_columns count
cols = len(empDF.columns)
```

## Selecting specific columns of a Dataframe

```
df.select("col1", "col2").show()
# Select All columns
df.select("*").show()
```

## Full content of the columns without truncation

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

## Handling missing or null values

```
# Fill all null values with 0
df.fillna(0)
#Fill specific columns with specified values
df.fillna({col1: 0,
"col2": "missing"})

```

## Joining two dataframes

```
#syntax:
join_new_df = df1.join(df2,
on="key_column",
how="inner")
#example
join_new_df =
empDF.join(deptDF, empDF.emp_id == deptDF.dept_id, "inner")
\ .show(truncate=False)
```

## Adding a new column to a DataFrame

```
df.withColumn("new_col", col("col1") + col("col2"))
```

## Dropping columns from a Dataframe

```
df.drop("col1")
```

# PySpark Fingertip Commands Cheat Sheet

## Grouping data by a colm and agg. with a function

```
#syntax
df.gro upB y("c ol1 " ).a gg( {"co 12": " mea n"})
#examples
df.gro upB y("d epa rtm ent ") \
    .ag g(s um( " sal ary " ).a lia s("s um_ - 
sal ary "), \
        avg ("sa lar y").a li as( " - 
avg _sa lar y"), \
        sum ("bo nus " ).a lia s("s um_ - 
bon us"), \
        max ("bo nus " ).a lia s("m ax_ - 
bon us") \
    ) \
    .sh ow( tru nca te= False)
```

## Stopping the SparkS ession

```
spark.stop()
```