

# PySpark 3.1.1 Cheat Sheet

A curated list of DataFrame transformations to reduce your StackOverflow time and get you familiarized with the Apache Spark DataFrame Python API.

## Legend:

**D** = prefix for DataFrame functions  
**F** = prefix for pyspark.sql.functions  
**C** = equivalent to F.col(a)  
**T** = prefix for pyspark.sql.types

**Σ** = groupby function  
**λ** = lambda function  
**Γ** = set of options

**a, b, c** = column names

**s** = string parameter  
**i** = integer parameter  
**d** = datetime parameter  
**p** = any parameter

**r** = regular expression  
**\*** = collection | array

## Types | Casting

**C**.astype(**T**)  
**C**.cast(**T**)  
**F**.rint  
**F**.from\_csv  
**F**.from\_json  
**F**.from\_timestamp  
**F**.from\_utc\_timestamp  
**F**.to\_csv  
**F**.to\_date  
**F**.to\_json  
**F**.to\_timestamp  
**F**.to\_utc\_timestamp

## Nulls

**C**.isNull  
**C**.isNotNull  
**C**.eqNullSafe(**b**)  
**F**.isnan  
**F**.isnull  
**F**.nanvl(**i**, **i**)  
**F**.coalesce(**a**, **b**)  
**D**.dropna

## Strings

**C**.startswith  
**C**.endswith  
**C**.contains  
**C**.isin(**\***)  
**C**.substr(**a**, **i**)  
**F**.lower  
**F**.upper  
**F**.initcap  
**F**.length  
**F**.trim  
**F**.ltrim  
**F**.rtrim  
**F**.split(**a**, **s**, **i**)  
**F**.substring(**a**, **s**)  
**F**.substring\_index(**a**, **s**, **i**)  
**F**.instr(**a**, **s**)  
**F**.locate(**a**, **s**, **i**)  
**F**.concat  
**F**.concat\_ws(**s**, **\***)  
**F**.repeat  
**F**.reverse  
**F**.translate(**a**, **s**, **s**)  
**F**.ascii  
**F**.encode(**a**, **Γ**)  
**F**.decode(**a**, **Γ**)  
**F**.overlay(**a**, **b**, **i**)  
**F**.levenshtein(**a**, **b**)  
**F**.soundex

## Dates | Time | Timezones

**C**.between(**d**, **d**)  
**F**.from\_unixtime  
**F**.from\_utc\_timestamp(**a**, **Γ**)  
**F**.timestamp\_seconds  
**F**.unix\_timestamp(**a**, **Γ**)  
**F**.current\_date  
**F**.current\_timestamp  
**F**.next\_day(**a**, **Γ**)  
**F**.date\_add(**d**, **i**)  
**F**.date\_format(**a**, **Γ**)  
**F**.date\_sub(**d**, **i**)  
**F**.date\_trunc(**Γ**, **Γ**)  
**F**.dayofweek  
**F**.dayofmonth  
**F**.dayofyear  
**F**.second  
**F**.minute  
**F**.hour  
**F**.days  
**F**.last\_day  
**F**.week\_of\_year  
**F**.month  
**F**.quarter  
**F**.year  
**F**.add\_months(**a**, **i**)  
**F**.months\_between(**a**, **b**)  
**F**.datediff(**a**, **b**)

## Arrays | Lists

**C**.getItem(**i**)  
**F**.size  
**F**.array  
**F**.array\_contains(**a**, **p**)  
**F**.array\_distinct  
**F**.array\_except(**a**, **b**)  
**F**.array\_intersect(**q**, **b**)  
**F**.array\_join(**a**, **s**)  
**F**.array\_union(**a**, **b**)  
**F**.slice(**a**, **i**, **i**)  
**F**.array\_sort  
**F**.sort\_array  
**F**.array\_position(**a**, **p**)  
**F**.array\_remove(**a**, **p**)  
**F**.array\_repeat(**a**, **i**)  
**F**.array\_overlap(**a**, **b**)  
**F**.array\_zip(**\*a**)  
**F**.zip\_with(**a**, **b**, **λ**)  
**F**.array\_max  
**F**.array\_min  
**F**.element\_at(**a**, **i**)  
**F**.explode  
**F**.explode\_outer  
**F**.posexplode  
**F**.posexplode\_outer  
**F**.flatten  
**F**.forall(**a**, **λ**)

## Maps | Dicts

**C**.getField(**s**)  
**C**.dropFields(**\*a**)  
**C**.withField(**s**, **a**)  
**F**.size  
**F**.struct(**\*a**)  
**F**.create\_map(**\*a**)  
**F**.json\_tuple(**a**, **\*s**)  
**F**.map\_concat(**\*a**)  
**F**.map\_entries  
**F**.map\_filter(**a**, **λ**)  
**F**.map\_from\_arrays(**a**, **b**)  
**F**.map\_from\_entries(**a**)  
**F**.map\_keys  
**F**.map\_values  
**F**.map\_zip\_with(**a**, **b**, **λ**)  
**F**.filter(**a**, **λ**)  
**F**.exists(**a**, **λ**)  
**F**.aggregate(**a**, **p**, **λ**)  
**F**.transform(**a**, **λ**)  
**F**.transform\_keys(**a**, **λ**)  
**F**.transform\_values(**a**, **λ**)

## High Order Functions on Arrays | Maps

## Aggregations

**Σ**.count()  
**Σ**.sum()  
**Σ**.max()  
**Σ**.min()  
**Σ**.mean()  
**Σ**.avg()  
**Σ**.pivot(**a**)  
**Σ**.agg(**\*λ**)  
**D**.crossJoin(**D**)  
**D**.crosstab(**a**, **b**)  
**D**.cube

## Regular Expressions

**C**.rlike(**r**)  
**F**.regexp\_extract(**a**, **r**, **i**)  
**F**.regexp\_replace(**a**, **r**, **s**)  
**D**.colRegex(**r**)