

# ⇒ Cleaning Data :-

⇒ To clean data, we need column expressions.  
Column expressions are constructed from  
existing columns, operators, and built-in func.

Standard SQL query commands to express  
transformations :- DISTINCT, WHERE, GROUP BY

## ⇒ Counting NOT NULL values in a tabl :-

Suppose, we have this  
table named users-dirty

User-id	timestamp	email	updated

⇒ To get count of all the records in a tabl,  
we run the command :-

[SELECT \* FROM users-dirty]

✓ It will also count the rows, that has all the  
attributes as NULL.

⇒ Now, to get number of NOT NULL values in a specific column, we need to specify that column name in COUNT() :-

[SELECT COUNT(email) FROM user-dirify]

it will give the no. of records in which the email column contains some value (and is not null)

⇒ To get the no. of NULL email records :-

[SELECT count(\*) - count(email) FROM users-dirify]

✓ Some result can also be obtained as:-

[SELECT COUNT\_IF email IS (l) FROM users-dirify  
gives no. of NULL email records.

⇒ COUNT\_IF() function :-

⇒ Takes a condition as parameter, and counts the no. of records that matches the passed condition.

Ex:-

[SELECT count\_if (len(email)>20) FROM users-dirify]

⇒ Python code to count no. of

## NULL values in a column :-

```
%python  
from pyspark.sql.functions import col  
usersDF = spark.read.table("users_dirty")
```

```
usersDF.selectExpr("count_if(email IS NULL)")  
usersDF.where(col("email").isNull()).count()
```

Can also be written as:-

→ usersDF.email

⇒ usersDF["email"]

gives same

result

i.e., 848 in  
our case

## Getting all the Distinct Rows :-

[sql

SELECT DISTINCT (\*) FROM users\_dirty

[python

usersDF.distinct().count()

distinct() dataframe method returns all the rows that has combination of values of all records unique (if no params passed).

⇒ Casting a column to diff. datatype

in PySpark :-

⇒ for it, we use the .cast() method:-

Suppose, in a table (test), we have a column → containing seconds from epoch.

sec-from-epoch
10901010

This column can be converted to timestamp as follows:-

In Python

```
from pyspark.sql.functions import col  
myDF = spark.table("test")
```

```
myDF2 = myDF.select(col("sec_from_epoch").CAST("timestamp"))
```

Casting a column to any other datatype

## Another way of casting in SQL :-

Suppose, user\_id in a table is of INT datatype.  
It can be casted to string as:-

[SELECT string (user\_id) FROM <table\_name>

If any column data is Binary Encoded, casting it to string datatype with above method will successfully find the string representation of the binary.

Kafka data is generally binary encoded, and can be converted to human readable string by the above method.

python

```
reqDF = (spark.table ("<table-name>")
         .select (col ("user_id").cast ("string"))
         )
```

## ⇒ date\_format( ) :-

⇒ Can be used in Spark-SQL Ansi SQL, as well as in pyspark.

[`date_format (<time-stamp column>, " <format-string>")`]

"mmm d,yyyy" ⇒ Jul 4, 2020

"HH:mm:ss" ⇒ 15:38:40

- - . and so on.

⇒ Used to get a date in any custom-format, from timestamp type / datetime type.

① [SELECT date\_format(created, "mmm d,yyyy")  
FROM tbl]

② from pyspark.sql import functions as F

[df.withColumn("ab", F.date\_format("created", "HH:mm:ss")).display())

## ⇒ regexp-extract()

⇒ Uses regular expressions to search for and perform string manipulation.