

Feature Engineering with PySpark • (divyanshu)

① Outlier Removal - Quantitative Filtering

- from pyspark.sql.functions import mean, stddev
- $\mu = df.agg(\{ 'x' : 'mean' \}).collect()[0][0]$
- $\sigma = df.agg(\{ 'x' : 'stddev' \}).collect()[0][0]$
- low = $\mu - 3 * \sigma$
- High = $\mu + 3 * \sigma$
- $df = df.filter((df['x'] \leq \text{low}) \& (df['x'] \geq \text{high}))$
↓
df.where also works.

② Text based filtering : where + like

- $df = df.where(df['gender'].like('MALE'))$
↓
NOT

③ Dropping columns : -df.drop(*['x1', 'x2']) ↓ unpacker.

④ Dropping NULL Records

- - df.dropna()
- - df.dropna(how='all', subset=['x1', 'x2'])
- - df.dropna(thresh=2)

⑤ Drop Duplicates w.r.t a col - df.dropDuplicates(['zipcode'])

⑥ FUNCTION for Minmax scaling (UDF) →

- def min_max_scaled(df, cols_to_scale):

for col in cols_to_scale:

max = df.agg({col: 'max'}).collect()[0][0]

min = df.agg({col: 'min'}).collect()[0][0]

new_col_name = 'scaled_' + col

df = df.withColumn(new_col_name, ((df[col] - min) / (max - min)))

return df

⑦ Create a UDF that drops columns with null > Threshold

def column_dropper (df, threshold):

```
total_row = df.count()
```

for col in df.columns:

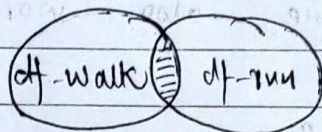
$$\text{missing_rows} = df.filter(df[col].isnull()) \cdot \text{count()}$$
$$\text{missing \%} = \text{missing_vows} / \text{total_vows}$$

if missing% > threshold:

$$df = df + \text{drop}(\text{COL})$$

return of

⑧ LEFT JOIN Example →



```
- df_wake = df_wake.withColumn('LAT', df_wake['LAT'].cast('double'))
```

- df-walk = df-walk with column ('lon', df-walk['lon'], cast[longitude])

- df-run = df-run.withColumn('LAT', round('LAT', 5))

$(\text{const } \text{round}(\text{val}))$

- # join condn

$$- \text{cond}^n = [\text{df-walk}['LAT'] \equiv \text{df-run}['LAT'], \text{df-walk}['LOW'] \equiv \text{df-run}['LOW']]$$

- Jom

$$- df = df_wait \cdot join(df_run, on=cond, how='left')$$

⑤ Creating New Features off already existing.

- $df = df_{\text{with Column ('Area', } df['L'] * df['B'])}$

⑥ Working with DateTime Features.

- from pyspark.sql.functions import to_date
- df = df.withColumn('DATE', to_date('mydate'))

→ if you wish to keep the time-component, use → to_timestamp

→ other options : from pyspark.sql.functions import year
: " " " " import month.
dayofmonth
weekofyear

to take delta b/w 2 dates.

datediff

⑦ Creating LAGGED Features →

- window()
- lag(col, count=1)

- from pyspark.sql.functions import lag ~~withColumn~~
- from pyspark.sql.window import window

- w = Window().orderBy(df['Date'])

- df = df.withColumn('Xlagby1wk', lag('x', count=1).over(w))

↓

DATE	X	Xlagby1wk
2013-01-01	2.3	null
2013-01-08	3.4	2.3
2013-01-15	4.2	3.4

Note: PySpark week starts - Sunday (=1) → Saturday (=7)

⑧ Joining on DateTime Features

→ B = B.withColumn('Year', year('Date'))

→ condn = [A['id'] == B['id'],

A['Year'] == B['Year']

→ A = A.join(B, on=condn, how='left')

A → has 'Year' column

B → has 'Date' column.