

Data  
/Processing  
/Cleaning  
/Imputation  
Pipeline  
on Real-Time data

Sana Wajid  
12/6/19

# 0P. Step

**P = Processing**  
**C = Cleaning**  
**I = Imputation**

## Assumptions:

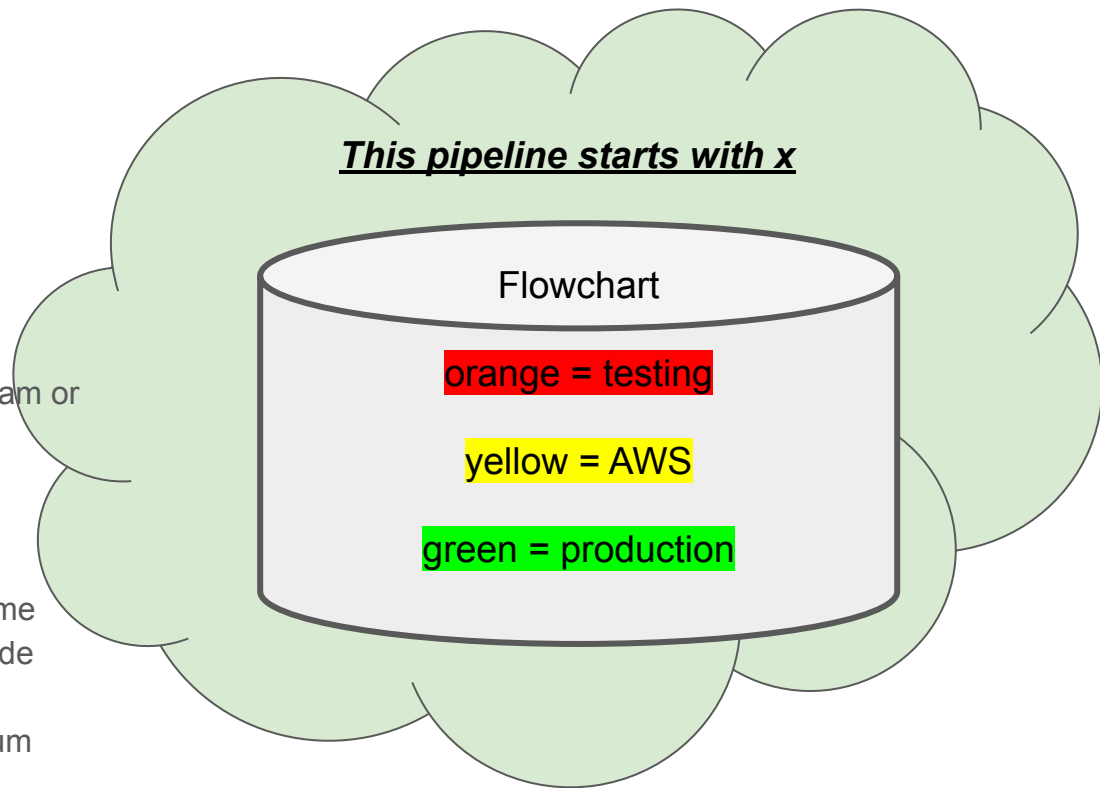
- For testing using csv files
- A pipeline can move it's starting point upstream or downstream.

## Checks:

1. Listed in order of least to most processing time
2. Numbers will refer to function numbers in code comment or headers in Jupyter notebook
  - a. e.g. 1P-2: File contains minimum number of headers

## Libraries:

Libraries used in {pandas, scikit-learn} ∈ Python,  
sed/awk, unix



## 2P. Raw data conversion to pandas

### Assumptions:

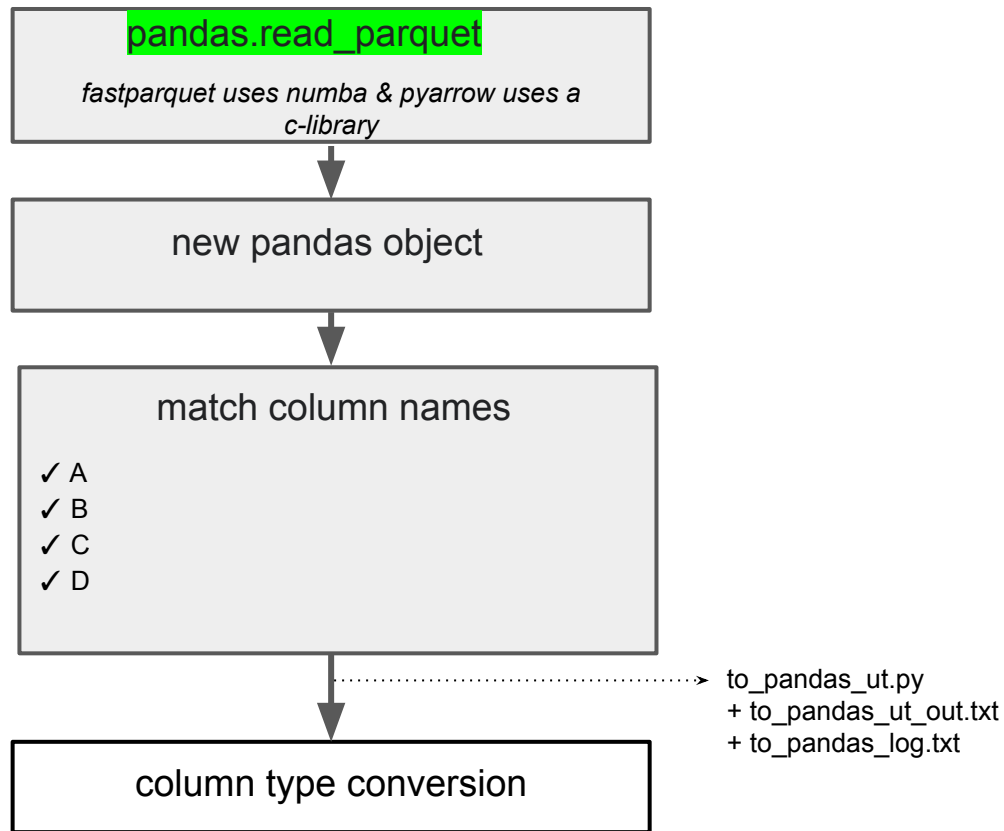
- File contains minimum number of headers
- File contains minimum number of rows
- File contains correct headers
- Column names don't contain spaces

### Functions:

1. `to_pandas`

### Libraries:

`read_parquet`  $\in$  `pandas`  $\in$  `Python`



# 3C. datetime column type conversion

## Assumptions:

- time column is formatted as:

YEAR-MONTH-DAY

space

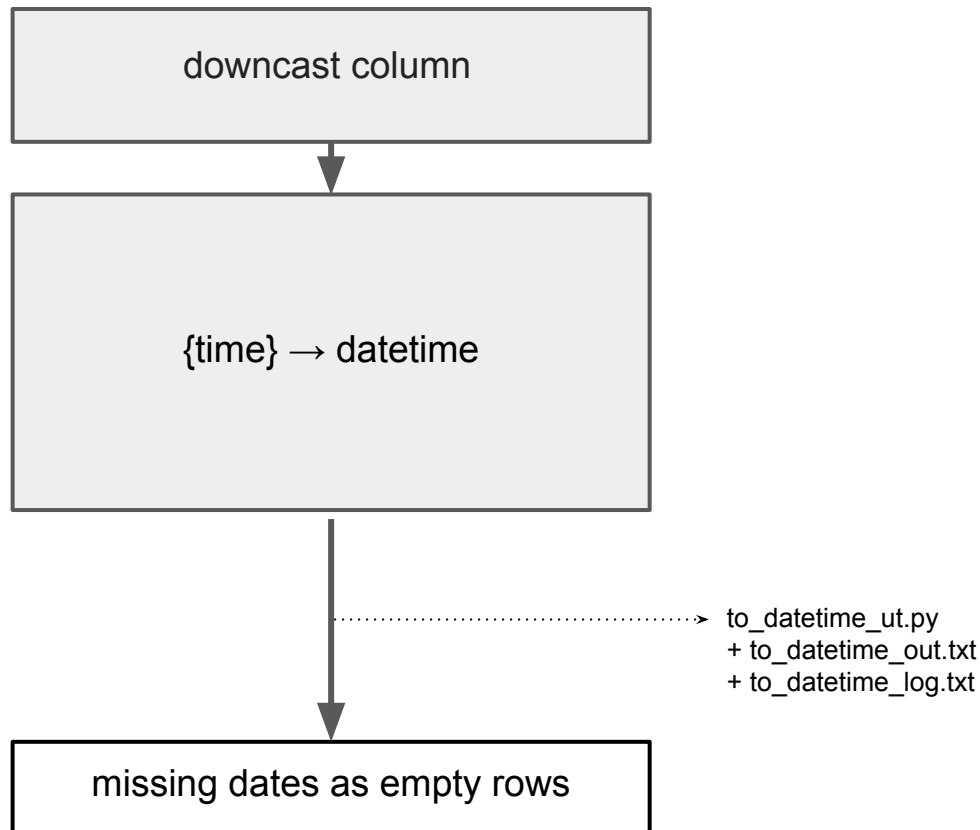
HOUR:MINUTE:SECOND

## Functions:

- column\_to\_datetime

## Libraries:

datetime ∈ Python



# 7I. Impute strategy, simple: scikit-learn

## Assumptions:

- fill value is defined per building and per x
- every  $t_0 : t_{\text{end}}$  has a row

## Functions:

accessor (get) functions return dataframe

1. `get_flanking_cluster` (dataframe obj of pandas type)

## Libraries:

scikit-learn  $\in$  Python

scikit-learn::impute.SimpleImputer

The imputation strategy.

- If “mean”, then replace missing values using the mean along each column. Can only be used with numeric data.
- If “median”, then replace missing values using the median along each column. Can only be used with numeric data.
- If “most\_frequent”, then replace missing using the most frequent value along each column. Can be used with strings or numeric data.
- If “constant”, then replace missing values with `fill_value`. Can be used with strings or numeric data.

dataframe does not contain any missing values

to\_sim\_imputer\_ut.py  
+ to\_sim\_imputer\_out.txt  
+ to\_sim\_imputer\_log.txt

# 71. Impute strategy, simple: scikit-learn

## Assumptions:

- fill value is defined per building and per x
- every  $t_0 : t_{\text{end}}$  has a row

## Functions:

mutators return True or False

2. `impute_cluster_by_mean`
3. `impute_cluster_by_median`
4. `impute_cluster_by_most_freq`
5. `impute_cluster_by_constant`

## Libraries:

scikit-learn  $\in$  Python

scikit-learn::impute.SimpleImputer

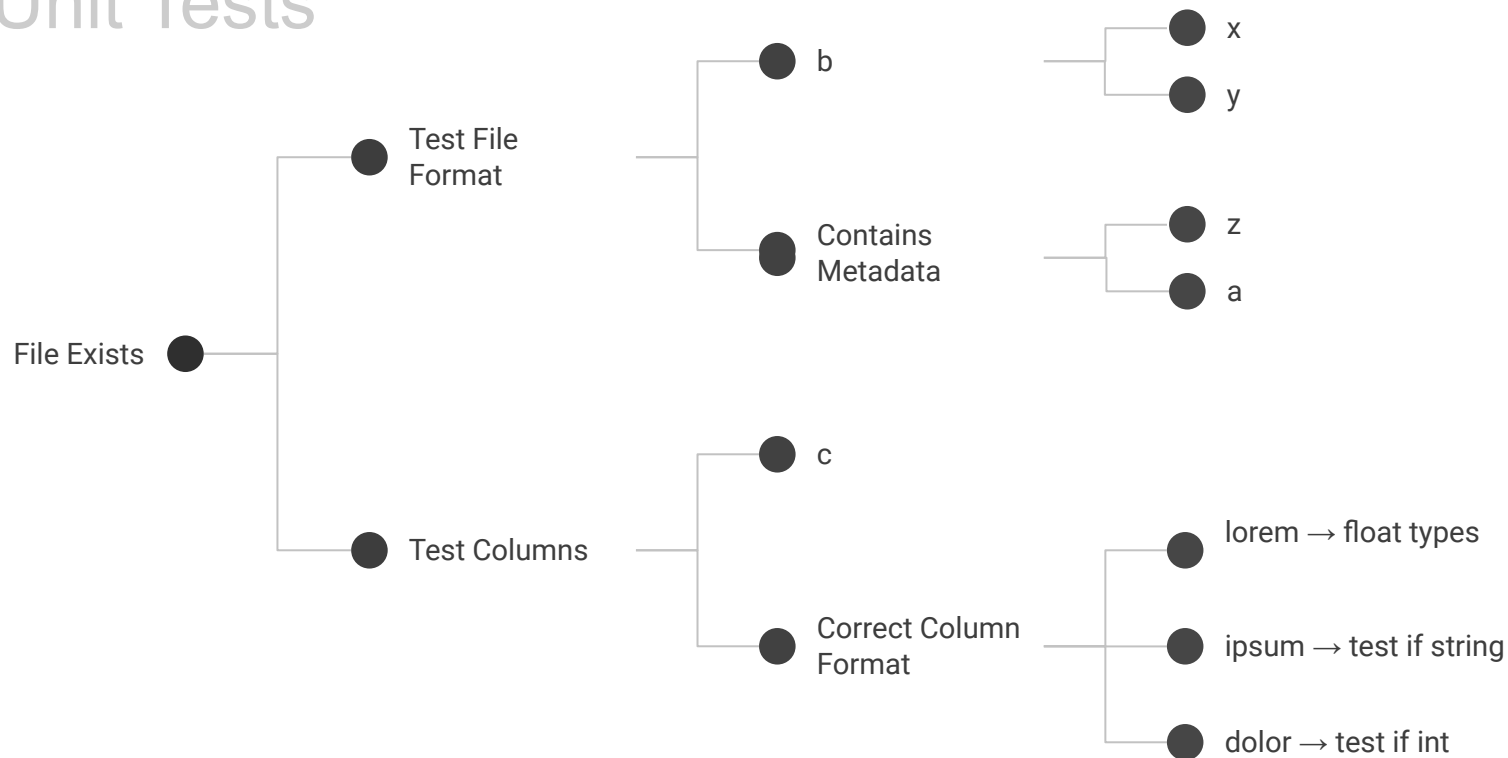
The imputation strategy.

- If “mean”, then replace missing values using the mean along each column. Can only be used with numeric data.
- If “median”, then replace missing values using the median along each column. Can only be used with numeric data.
- If “most\_frequent”, then replace missing using the most frequent value along each column. Can be used with strings or numeric data.
- If “constant”, then replace missing values with `fill_value`. Can be used with strings or numeric data.

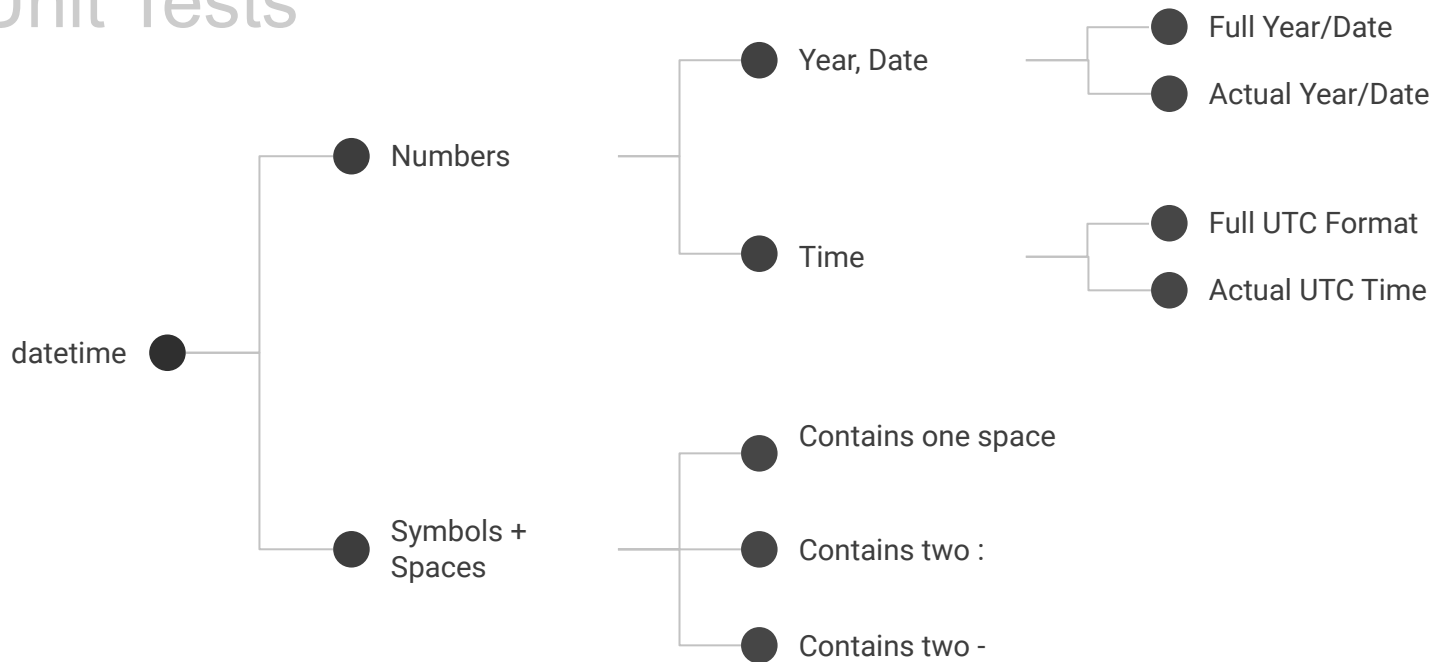
dataframe does not contain any missing values

to\_sim\_imputer\_ut.py  
+ to\_sim\_imputer\_out.txt  
+ to\_sim\_imputer\_log.txt

# Unit Tests



# Unit Tests





# References

etc, more will be added