

# The Six Steps to build a preprocessor

Before starting!  
**\*\*Know your columns\*\***

What are the dtypes? -> df.dtypes

What are the values? -> df[col].value\_counts()

01

Are you missing  
any values?

Imports

```
from sklearn.impute import SimpleImputer
from sklearn.preprocessing import OneHotEncoder
from sklearn.pipeline import make_pipeline
from sklearn.compose import
make_column_transformer
```

For EACH column where you're missing data, but NOT the ones that aren't, instantiate an imputer

- imputer = SimpleImputer(strategy = ' ')
- strategies: 'mean', 'median', 'most\_frequent', 'constant' (fill\_value = ' ')
- use argument add\_indicator = True to negate effect of imputation

- Nominal categorical columns should be one-hot encoded
  - ohe = OneHotEncoder()
- If the column has too many unique values relative to the length of the dataset, consider dropping the column instead of one-hot encoding

02

Does it need to be  
one-hot encoded?

03

Does it need to  
be scaled?

- Numeric columns whose values do NOT fall between 0 and 1 should be scaled for MOST models
  - scaler = StandardScaler()
- Numeric columns whose values fall between 0 and 1 are already scaled

- You should have one pipeline for each group of columns that are receiving different treatments
- ex) group1\_pipe = make\_pipeline(mean\_imputer, scaler)
- ex) group2\_pipe = make\_pipeline(constant\_imputer, ohe)

04

Collect steps into  
pipelines

05

Create  
tuples

Combine pipelines and columns together in tuples

- ex) group1\_tuple = (group1\_pipe, group1\_columns)
- ex) group2\_tuple = (group2\_pipe, group2\_columns)

- preprocessor = make\_column\_transformer(group1\_tuple, group2\_tuple, remainder = ' ')
- Use remainder parameter option 'drop' or 'passthrough'

Create the  
preprocessor

06

You have a preprocessor!  
Congratulations!