

DATA PREPROCESSING: HANDLING MISSING VALUES

Something you have to deal with will need to apply human intelligence to solve

pandas - `IsNull().sum()` - for each column and row

- can `dropna()` on the dataframe (rows)

- can drop columns `dropna(axis=1)`

- can apply thresholds, limits on drops, etc

Row: Too much data lost: over fitting, wider confidence intervals, etc.

May bias sample

Columns: May lose information in features (underfitting)

Before dropping or imputing missing values, ask:

- What are the mechanisms causing the missing values?
- Are these missing values missing at random?
- Are there rows or columns missing you are not aware of?

Missing at random? Impute

estimate using mean, median, most frequent
other val from business context (categoricals),

sklearn Imputer

Advanced Methods - Imputing Missing Values

MICE - multiple imputation by chained equations
`sklearn.impute.MICEImputer (v0.20)`

Python fancy impute package

KNN imputer

Soft Impute

MICE

etc