

Data Cleaning

① Why Data Cleaning is important

- a) See Instances are proper
- b) Labels: Of variables being predicted (If someone labeled it wrong in training it affects algo)
- c) Algo: The algo you use affects results based on available data
- d) Features: Proper feature data is imp for good prediction/classification
- e) Model: Model itself assumes this data is for real world (Not sure)

Garbage -in → Garbage -out

Main Problems

- Lack of Data
- Too much data
- Bad data

② How Can Data be Messy?

- a) Duplicate data → odd weights
- b) Inconsistent text & types
- c) Missing data
- d) Outliers → few values dominate over other values by large amount
- e) Getting data from Multiple Sources like DB, cloud etc and when talking together abt them might create mismatches

→ Its Suggestible to have a look to the data & filter duplicates

→ Remove → you can clean it quick but sometimes you might lose good info and make it biased

Impute → Replace value with mean (or) median

Mask → Create a Category for missing data

→ To find the size plots (Histogram, Density plot, Box plot)
Statistic - (~~Box plots~~ Interquartile Range, Standard deviation)
Residuals (Standardized, Pooled, Standardized)

Data cleaning is very important part before training & analyzing your data. With a wrong data, you get wrong insights and can't fulfill your task although your model is pretty good & robust.

Garbage in → Garbage out

~~How to~~ what to check for cleanliness

- ① Validity
- ② Accuracy
- ③ Completeness
- ④ Consistency
- ⑤ Uniformity

How to check

- ① Inspect data
 - ↳ Using Data Profiling & visualization find if data is proper (or) not. Outliers are there (or) not. Data is skewed (or) not. Data is missing (or) not.

- ② Clean data

↳ Next Important part is cleaning. If for example data is missing, you can either delete the row (or) fill the value using statistical methods. In general, if data is not skewed, "mean" is best to fill. For skewed data, "median" is preferred.

Remove outliers. Sometimes outliers also might carry good information. Need to analyze & take decision.

③ Verify

↳ Verify If ~~data~~ is cleaned (or) not after you clean it

④ Report

↳ find the reason for getting bad data. If you can find the source of problem. You can prevent next observations to be noisy. you can even understand better about Data.