

Data preprocessing

Elena Tuzhilina

Stanford University, Department of Statistics

elenatuz@stanford.edu

January 18, 2022

Once upon a time I worked as a data scientist...

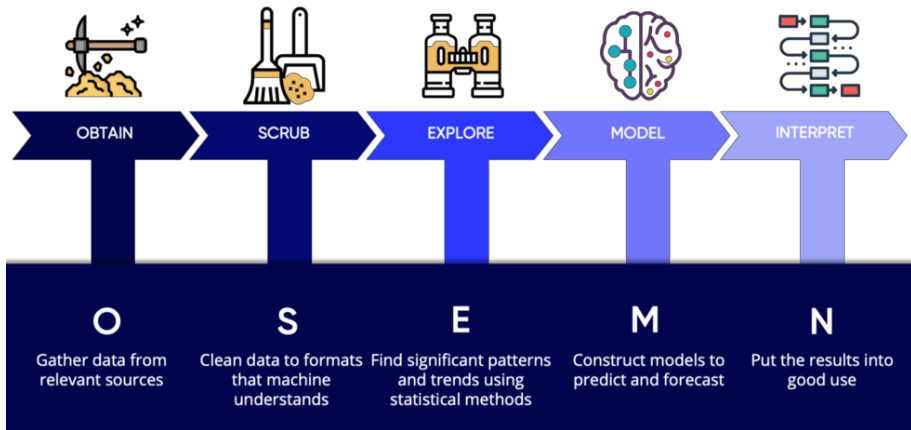


Statistician



Data Scientist

Data science process

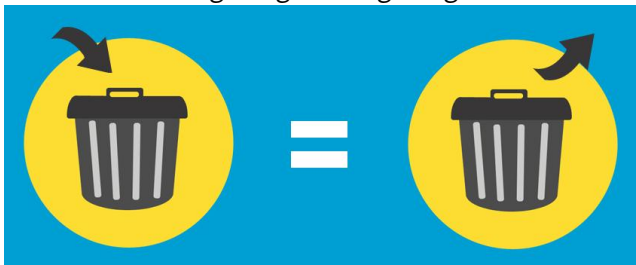


Definition

Data preprocessing is a data mining technique which is used to transform the raw data in a useful and efficient format.

Why is it important? Nonsense input data produces nonsense output.

GIGO: garbage in = garbage out

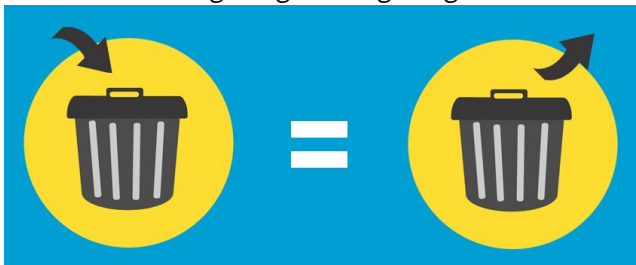


Definition

Data preprocessing is a data mining technique which is used to transform the raw data in a useful and efficient format.

Why is it important? Nonsense input data produces nonsense output.

GIGO: garbage in = garbage out



What is the first step?

Data preprocessing stages

- ① Sanity checks
- ② Data reduction
 - remove uninformative features
 - remove correlated features
- ③ Data cleaning
 - remove duplicates
 - handle missing values
 - handle outliers
- ④ Data transformation
 - do binning
 - handle categorical variables
 - apply normalization/standardization
 - handle skewed distribution

country	year	cases	pop
Afghanistan	1999	745	19997071
Afghanistan	2000	666	20095360
Brazil	1999	37737	17206362
Brazil	2000	40488	17404898
China	1999	212258	1272915272
China	2000	213766	1280428583

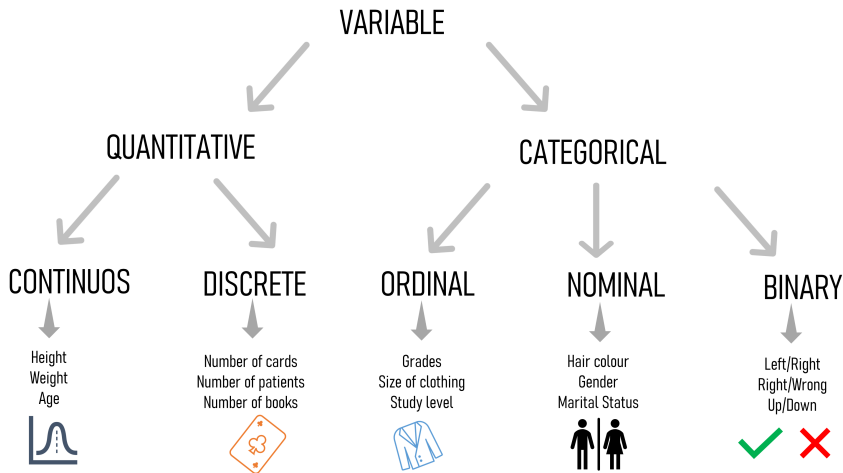
Variables

country	year	cases	pop
Afghanistan	1999	745	19997071
Afghanistan	2000	666	20095360
Brazil	1999	37737	17206362
Brazil	2000	40488	17404898
China	1999	212258	1272915272
China	2000	213766	1280428583

Observations

Preprocessing can be done on observation or variable level

Data types



Preprocessing strongly depends on variable type

1. Sanity checks

Example:

- out-of-range values, e.g. *Income: -100*
- impossible combinations, e.g. *Sex: Male; Pregnant: Yes*



Study the data and run some simple checks

2. Data reduction

Uninformative features

- ① a lot of missing values, e.g. *contains* > 70% of NAs
- ② categorical with too many values, e.g. *UserID: jkaj1daj*
- ③ constant/almost constant, e.g. *Year: 2021*

Correlated features

- ① Check Pearson/Spearman correlation (numerical)
- ② Check Cramer's V statistics (categorical)

*You can remove some redundant features
or
you can do nothing and apply regularization*

2. Data reduction

Uninformative features

- ① a lot of missing values, e.g. *contains* > 70% of NAs
- ② categorical with too many values, e.g. *UserID*: *jkaj1daj*
- ③ constant/almost constant, e.g. *Year*: 2021

Correlated features

- ① Check Pearson/Spearman correlation (numerical)
- ② Check Cramer's V statistics (categorical)

Do you remove features? If yes, what criteria do you apply?

3. Data cleaning: duplicates



You can remove some redundant observations

3. Data cleaning: missing values

- ① Remove observation
- ② Impute numerical
 - use mean/median
 - use k nearest neighbours
 - use low-rank matrix approximation
- ③ Impute categorical
 - add new category "missing"
 - use most frequent values

country	year	cases	population
Afghanistan	1999	725	19087071
Afghanistan	2000	2686	20593360
Brazil	1999	3737	17200362
Brazil	2000	8048	17450898
China	1999	21223	127291272
China	2000	21476	128042583

values

*You can impute missing values
however
some machine learning models can deal with them*

3. Data cleaning: missing values

- ❶ Remove observation
- ❷ Impute numerical
 - use mean/median
 - use k nearest neighbours
 - use low-rank matrix approximation
- ❸ Impute categorical
 - add new category "missing"
 - use most frequent values

country	year	cases	population
Afghanistan	1999	725	19987071
Afghanistan	2000	2676	20593360
Brazil	1999	36797	17200362
Brazil	2000	80483	17450898
China	1999	212293	127291272
China	2000	216706	128042583

values

How do you handle missing values?

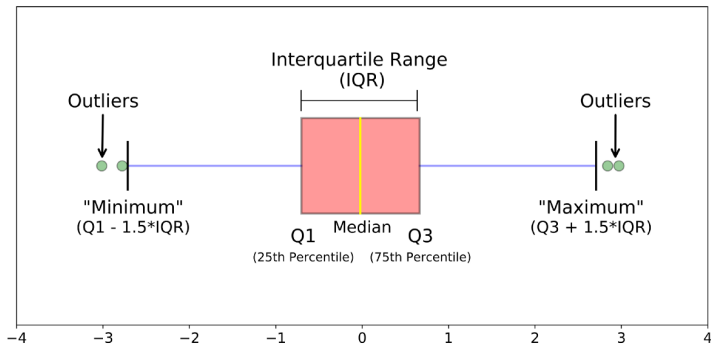
3. Data cleaning: outliers



An outlier is an observation that lies outside the overall pattern of a distribution

3. Data cleaning: outliers

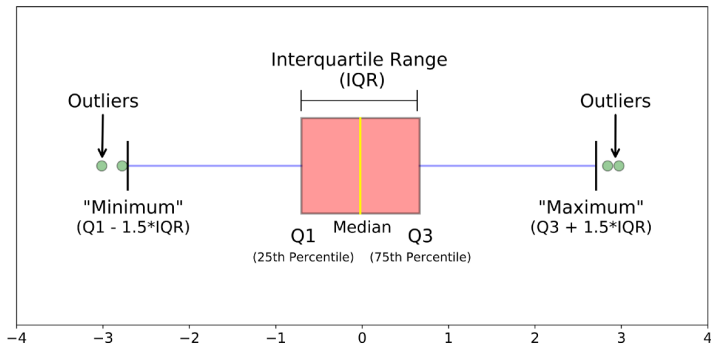
- 1 Remove observation
- 2 Change/trim the value
- 3 Apply transformation, e.g. log-transformation



To detect outliers you can use IQR

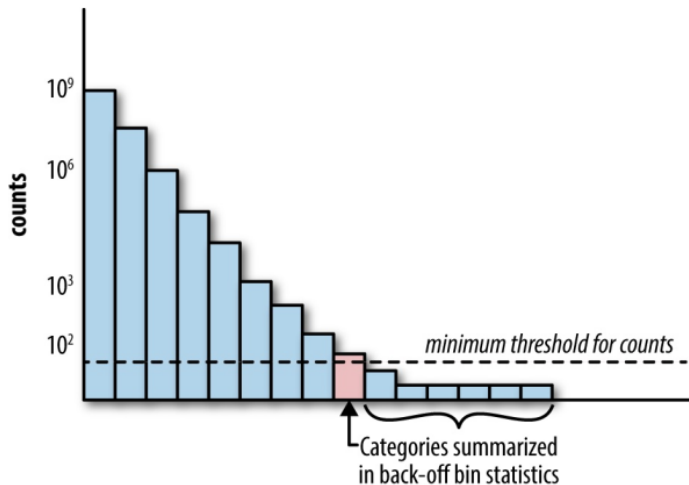
3. Data cleaning: outliers

- 1 Remove observation
- 2 Change/trim the value
- 3 Apply transformation, e.g. log-transformation



How do you deal with outliers?


4. Data transformation: binning



You can apply binning to reduce number of different values of a categorical feature

4. Data transformation: categorical variable

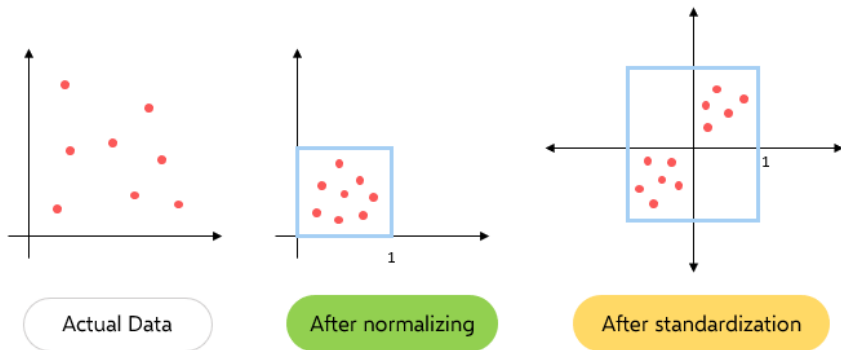
Color	
Red	
Red	
Yellow	
Green	
Yellow	



Red	Yellow	Green
1	0	0
1	0	0
0	1	0
0	0	1
0	1	0

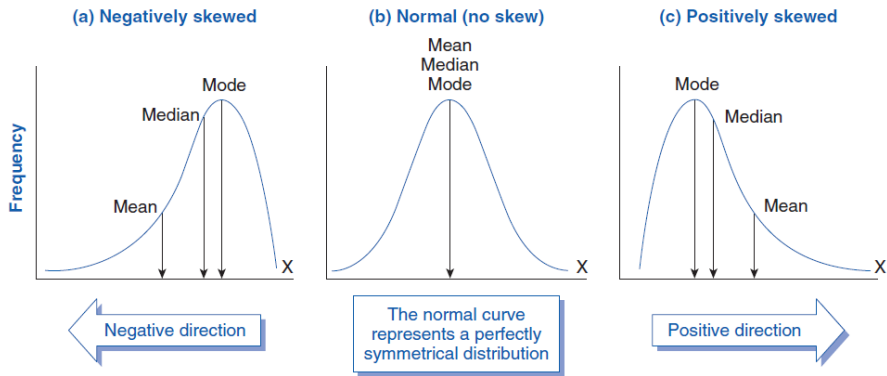
You need to convert all categorical variables to numeric format

4. Data transformation: normalization/standardization



Some ML techniques requires data scaling and centering

4. Data transformation: skewed distribution



*Some ML techniques does not work well for very skewed distributions
You can apply log-transformation to these features*

- ① **Do you do any data preprocessing? How complex is it?**
- ② **From your point of view, how important is data preprocessing step?**
- ③ **What tools do you use for data preprocessing?**

Images used

- 1 Data Scientist
- 2 OSEMN
- 3 GIGO
- 4 Variables, Observations, Values
- 5 Not sure I can trust this data
- 6 Duplicates
- 7 Outliers
- 8 Boxplot
- 9 Binning
- 10 Feature encoding
- 11 Data standardization
- 12 Skewed distribution