



# Feature Engineering



# Feature Engineering



- Transform Variables
- Extract Features
- Create New Features

# Missing Data

- Scikit-learn and other libraries can't work with missing data



# Missing Data Imputation Techniques

## Numerical Variables



- ☐ Mean / Median Imputation
- ☐ Arbitrary value imputation
- ☐ End of tail imputation

## Categorical Variables



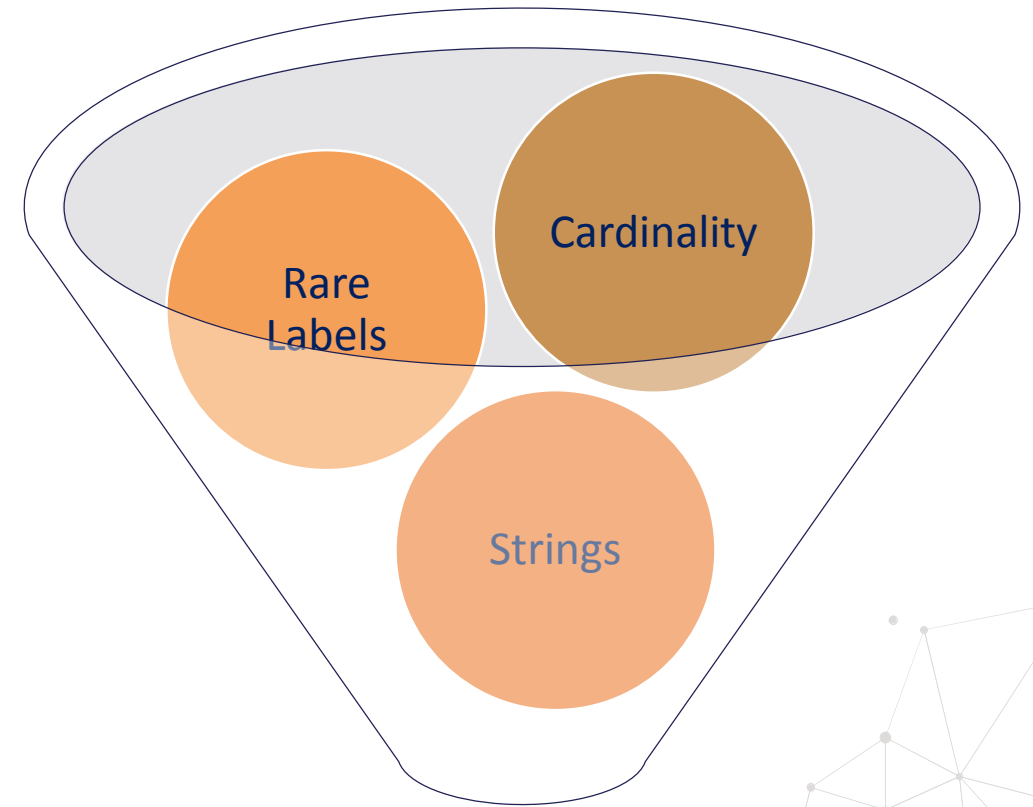
- ☐ Frequent category imputation
- ☐ Adding a “missing” category

## Both



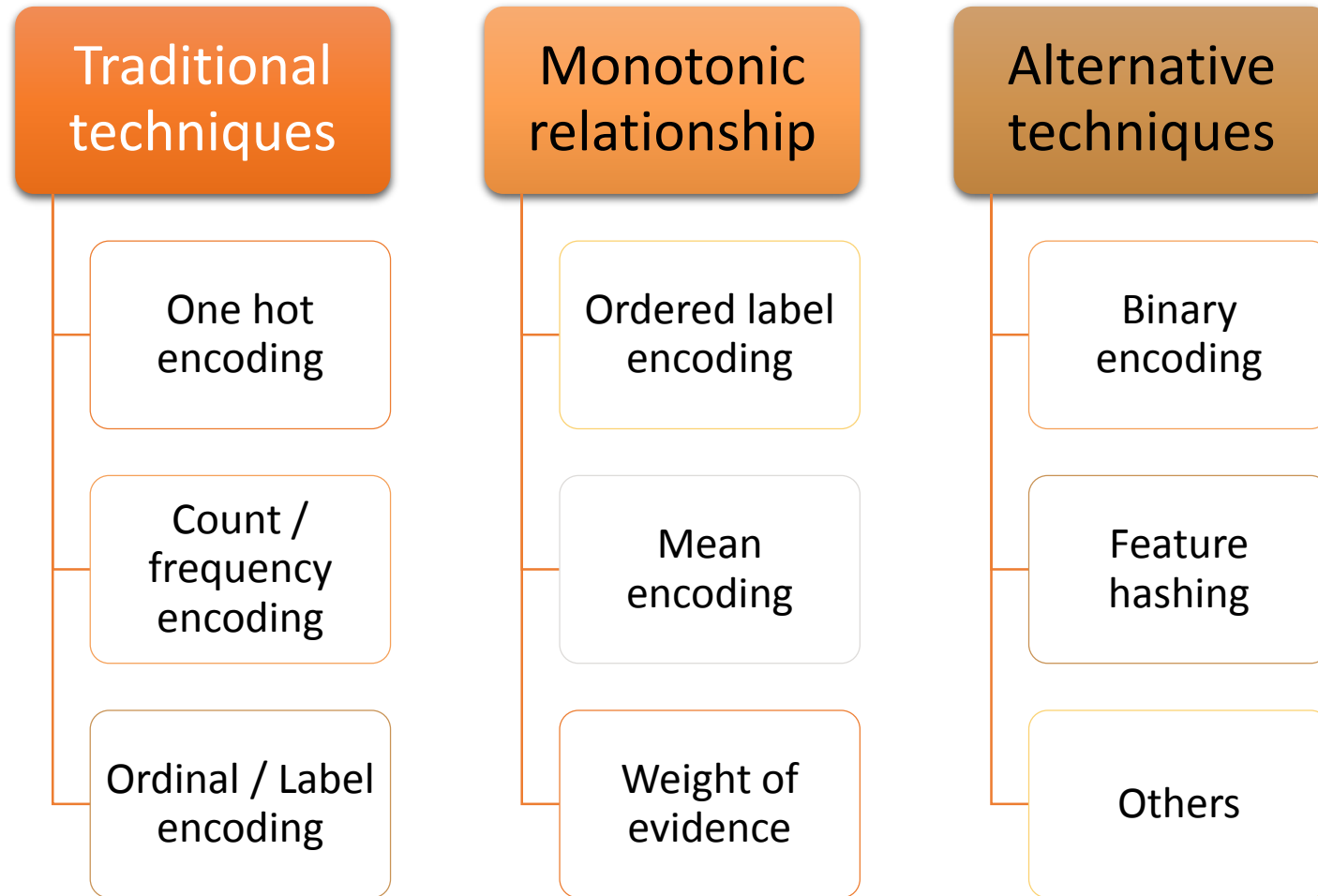
- ☐ Complete Case Analysis
- ☐ Adding a “Missing” indicator
- ☐ Random sample imputation

# Categorical Variables

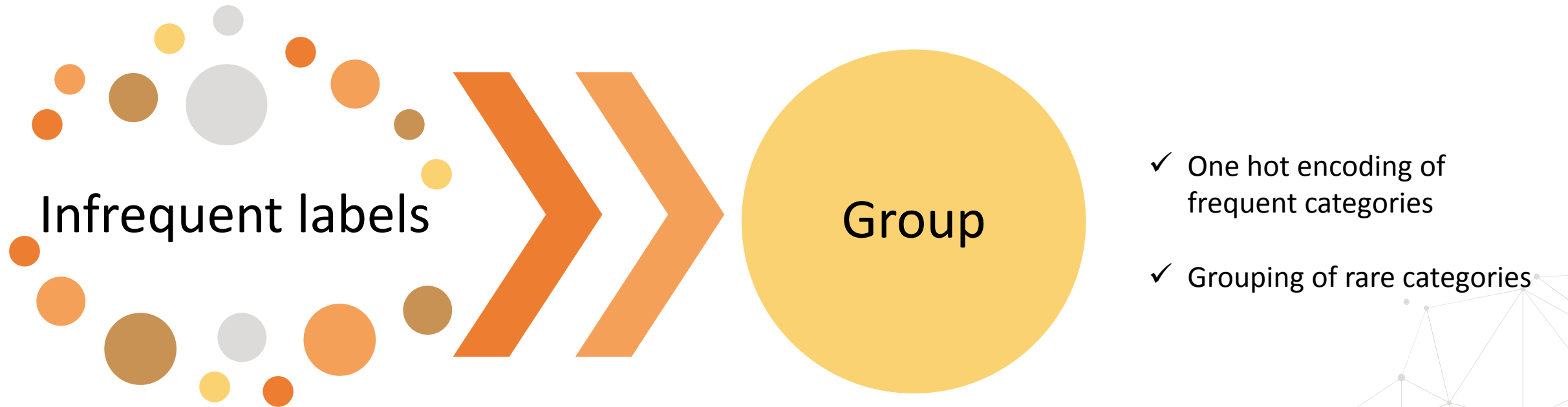


Over-fitting  
Particularly in decision trees

# Categorical Encoding Techniques



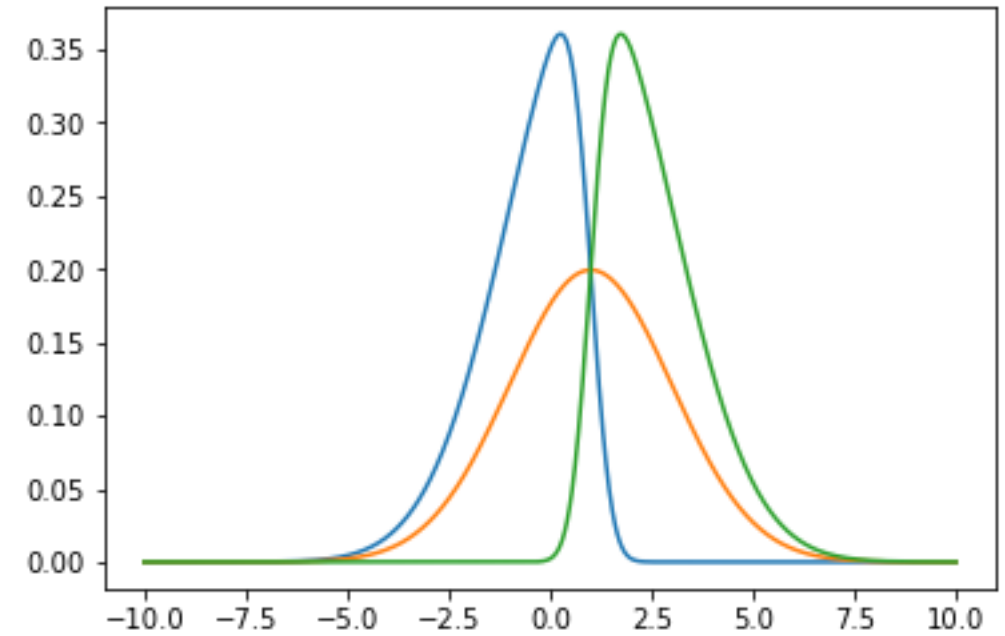
# Encoding Techniques: Rare labels



Particularly important for model deployment

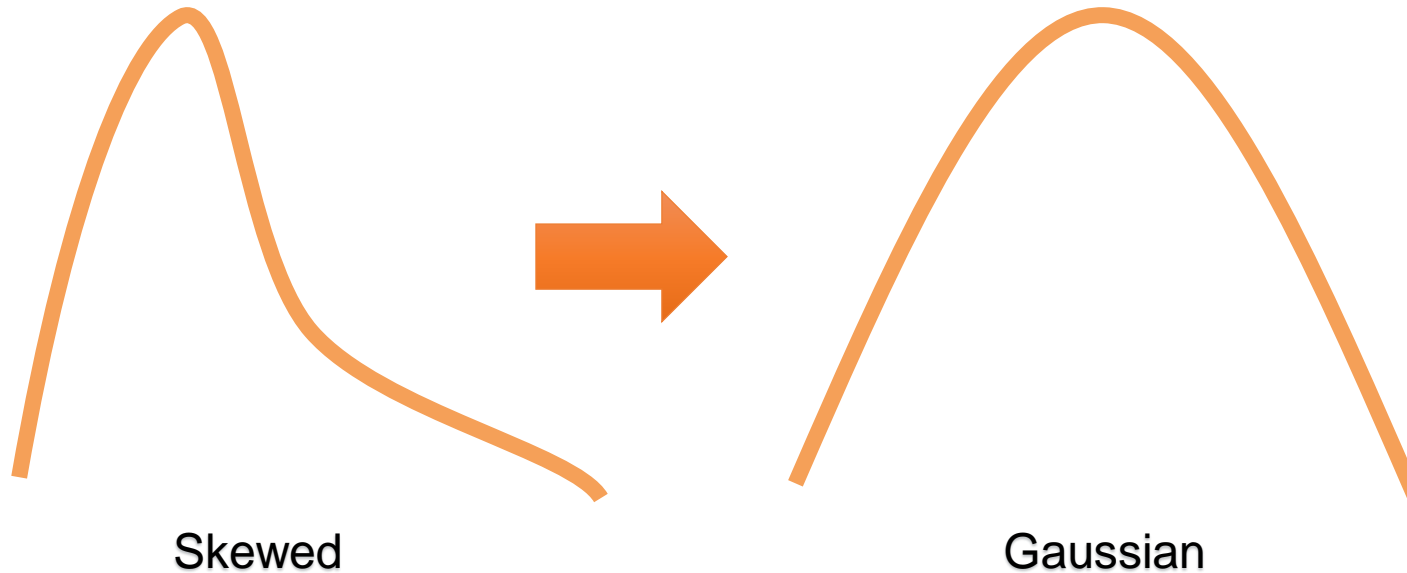
# Distributions

- Some models make assumptions on the variable distributions





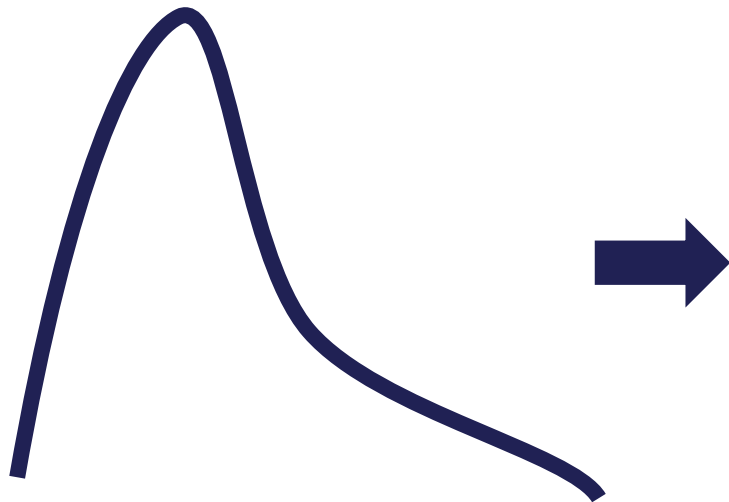
# Mathematical transformations



## Variable transformation

- Logarithmic
- Exponential
- Reciprocal
- Box-Cox
- Yeo-Johnson

# Discretisation



Skewed



Improved value spread

Unsupervised

Equal-width

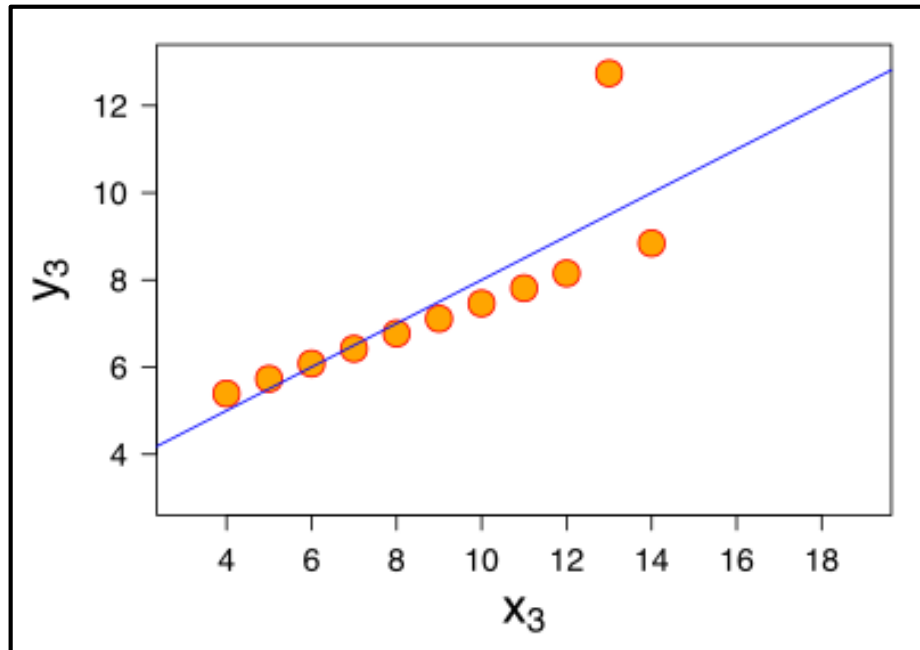
Equal-frequency

K means

Supervised

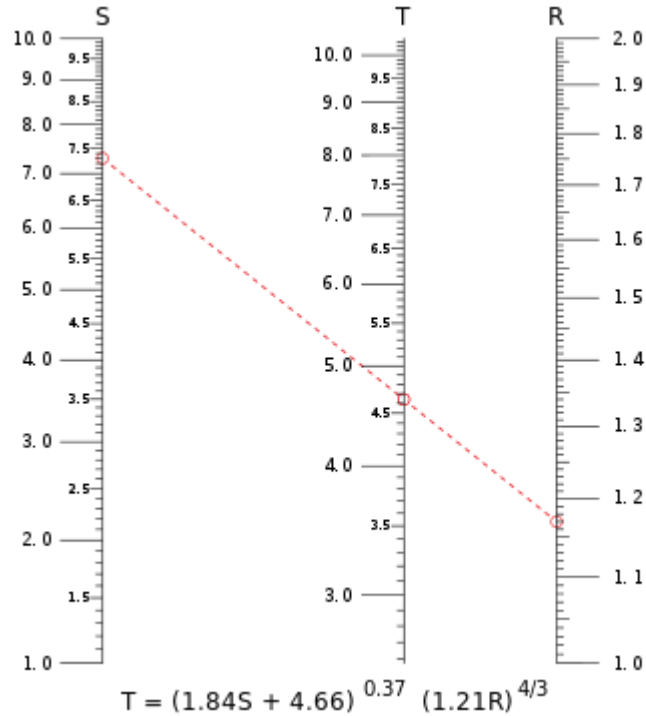
Decision Trees

# Outliers



- Discretisation
- Capping / Censoring
- Truncation

# Variable Magnitude



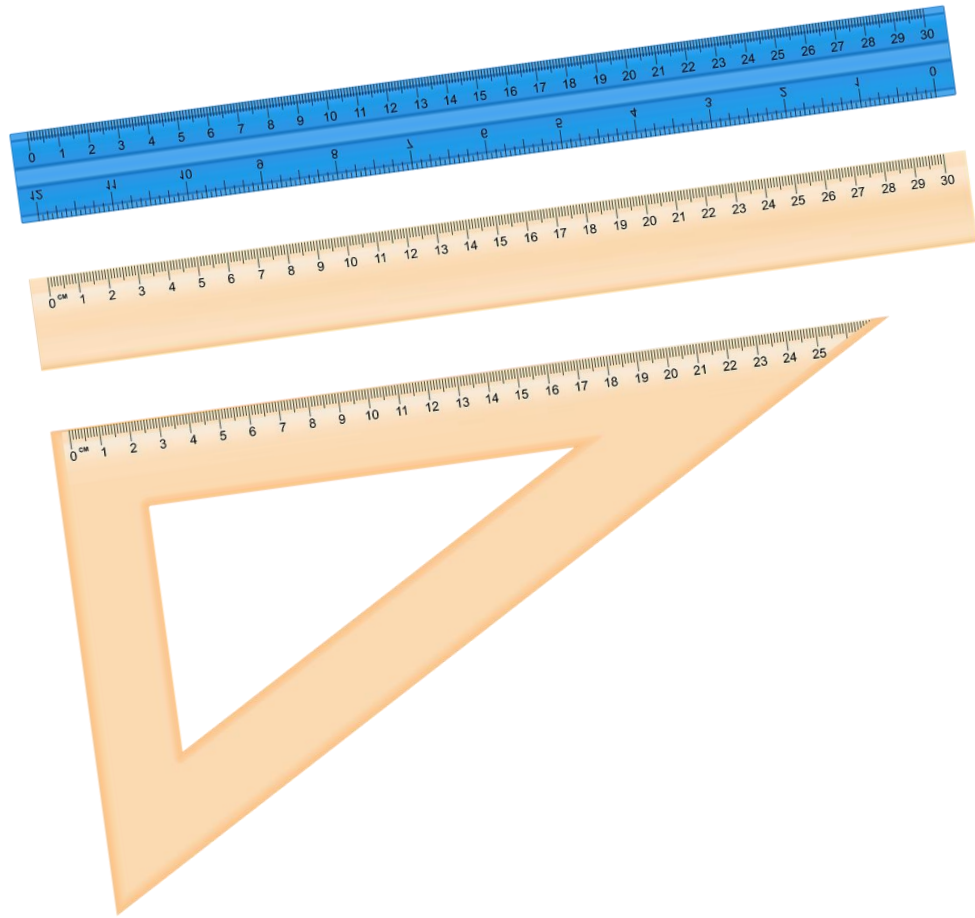
**The machine learning models affected by the magnitude of the feature:**

- Linear and Logistic Regression
- Neural Networks
- Support Vector Machines
- KNN
- K-means clustering
- Linear Discriminant Analysis (LDA)
- Principal Component Analysis (PCA)

**Machine learning models insensitive to feature magnitude are the ones based on Trees:**

- Classification and Regression Trees
- Random Forests
- Gradient Boosted Trees

# Feature scaling methods



## Scaling methods

- **Standardisation**
- Mean normalisation
- **Scaling to maximum and minimum**
- Scaling to absolute maximum
- Scaling to median and quantiles
- Scaling to unit norm

# Datetime Variables



- Day, Month, semester, year
- Hour, min, sec
- Elapsed Time
  - Time between transactions
  - Age

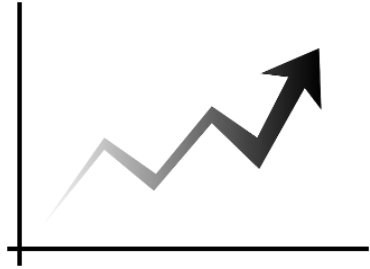
# Text



An insurance claim  
A formal request to an ins  
payment based on the te  
Insurance claims are re  
out to the insured

- Characters, words, unique words
- Lexical diversity
- Sentences, paragraphs
- Bag of Words
- TFIDF

# Transactions and Time Series



Aggregate data

- Number of payments in last 3, 6, 12 months
- Time since last transaction
- Total spending in last month



# Geo Data



- Distances

# Feature Combination



- **Ratio:** Total debt with income → Debt to income ratio
- **Sum:** Debt in different credit cards → total debt
- **Subtraction:** Income without expenses → disposable income

# Thank you

[www.trainindata.com](http://www.trainindata.com)