

# Feature Engineering and Selection for Machine Learning

**Soledad Galli, PhD**

CorrelAid meetup

Berlin, 30<sup>th</sup> November 2019

# Intro to the speaker



Soledad Galli, PhD  
Lead Data Scientist

## Finance

- Credit Risk
- Fraud Prevention

## Insurance

- Fraud Prevention
- Motor Claim Liability
- Car repair / scrapping

## Author, instructor

- Online Courses
- Tech Book
- Articles, talks, open-source

# Machine Learning in Finance and Insurance



**Credit Risk**



**Marketing**

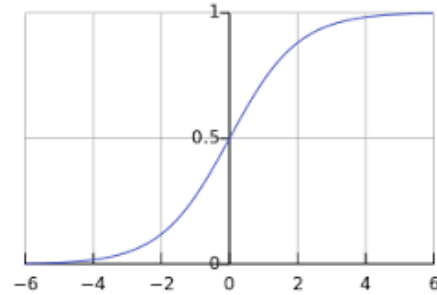


**Pricing**

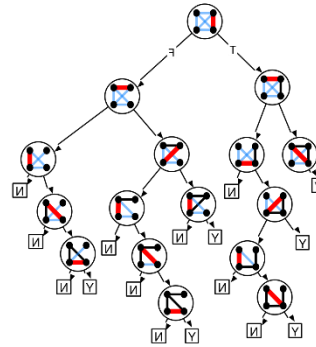
# Machine Learning Pipeline



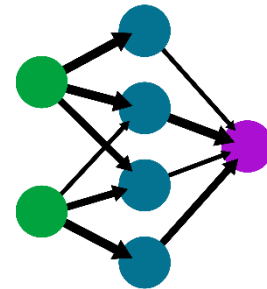
Data Sources



**Linear Models**  
Logistic Regression  
MARS



**Tree Models**  
Random Forests  
Gradient Boosted Trees



**Neural Networks**



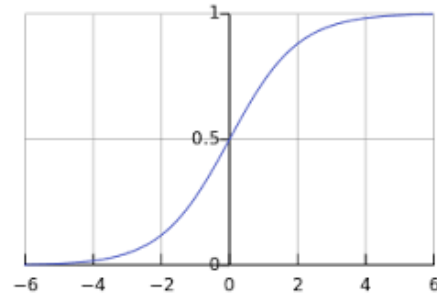
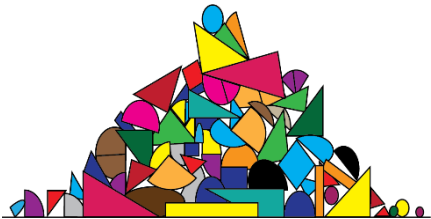
**Average  
Probability**

**Continuous  
Output**

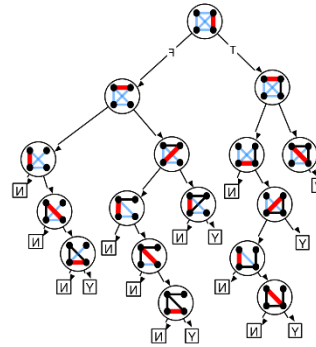
# Machine Learning Pipeline



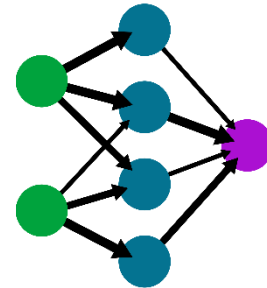
Data Sources



**Linear Models**  
Logistic Regression  
MARS



**Tree Models**  
Random Forests  
Gradient Boosted Trees



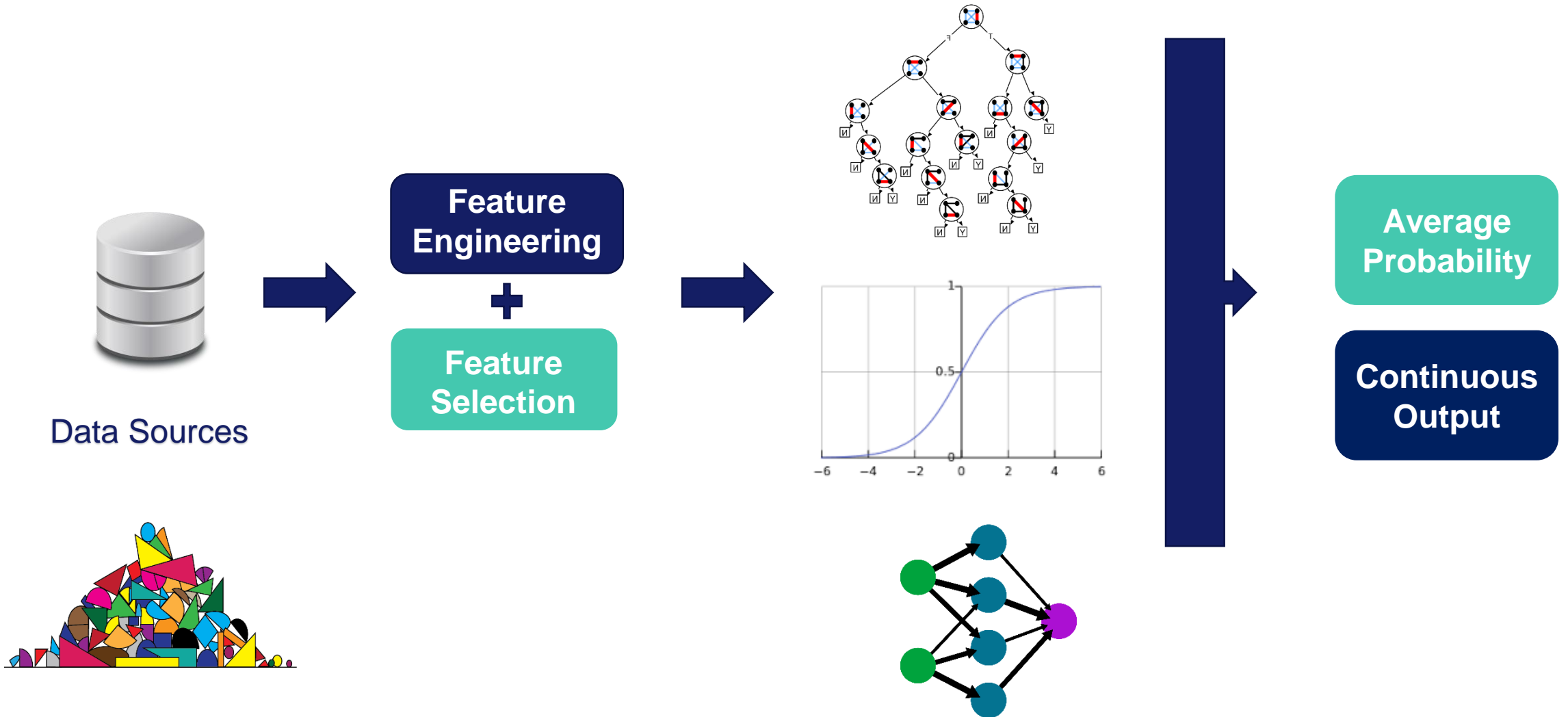
**Neural Networks**



**Average  
Probability**

**Continuous  
Output**

# Machine Learning Pipeline



# Data Preparation Journey

- Common issues found in variables
- Feature / variable engineering: solutions to the data issues
- Feature selection: do we need to select features?
- Feature / variable selection methods
- Overview and knowledge sources

# Data Preparation Journey

- Common issues found in variables
- Feature / variable engineering: solutions to the data issues
- Feature selection: do we need to select features?
- Feature / variable selection methods
- Overview and knowledge sources

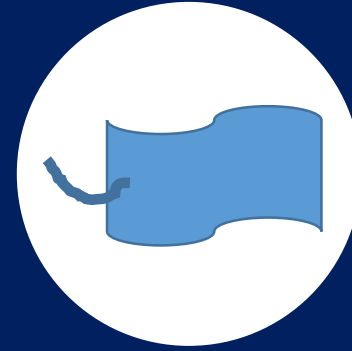


# Problems in Variables



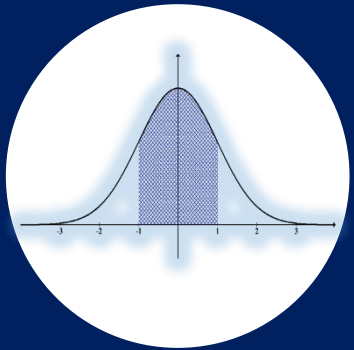
## Missing data

Missing values within a variable



## Labels

Strings in categorical variables



## Distribution

Normal vs skewed  
Scale / magnitude



## Outliers

Unusual or unexpected values

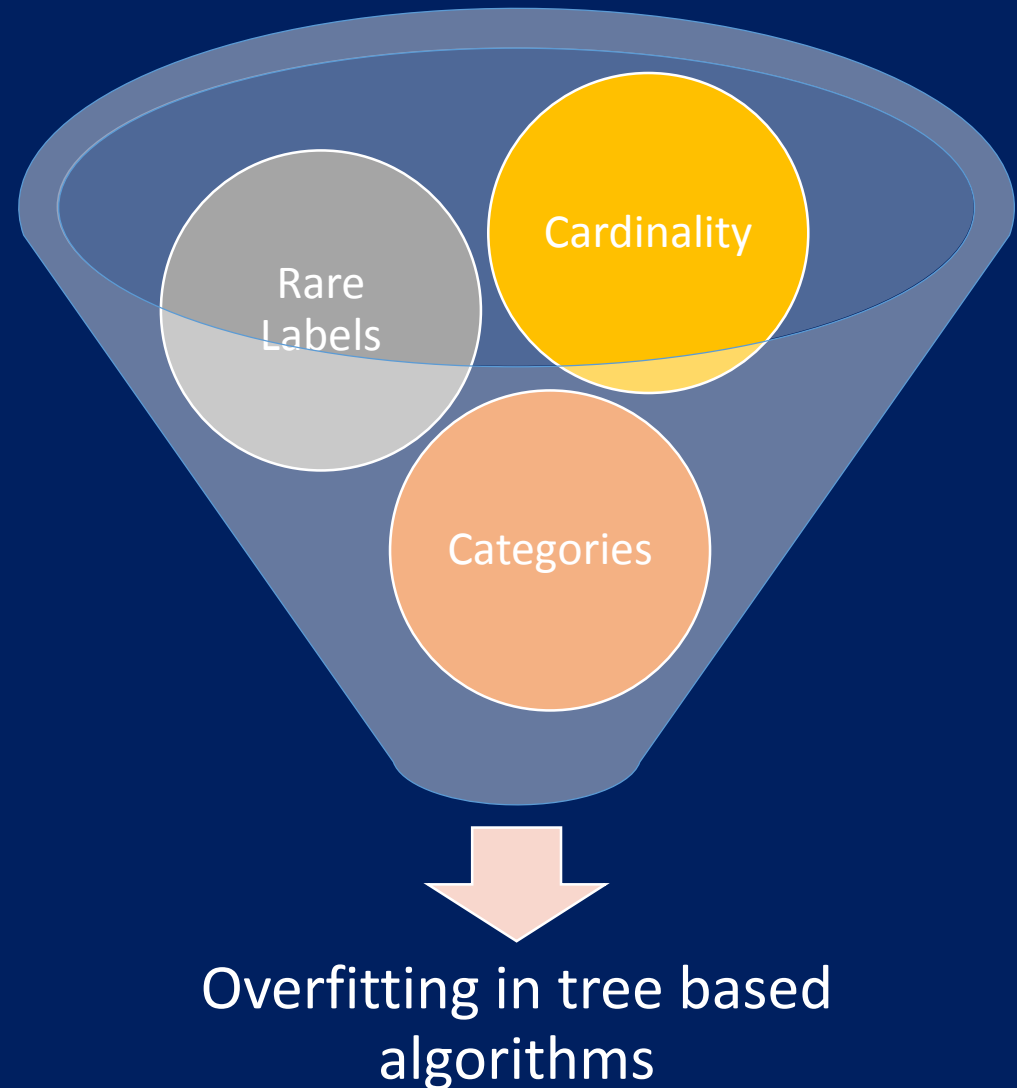
# Missing Data

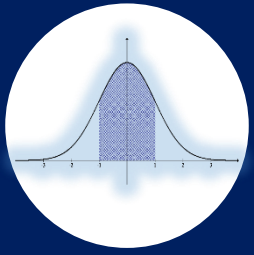
- Missing values for certain observations
- Affects all machine learning models
  - Scikit-learn



# Labels in categorical variables

- Cardinality: high number of labels
- Rare Labels: infrequent categories
- Categories: strings
  - Scikit-learn

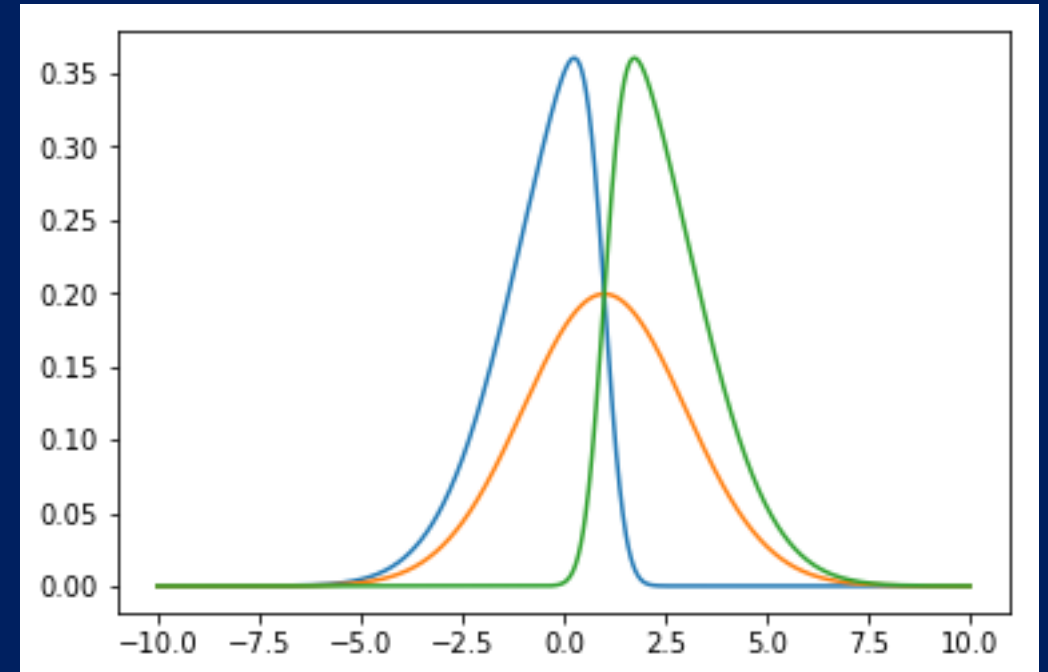




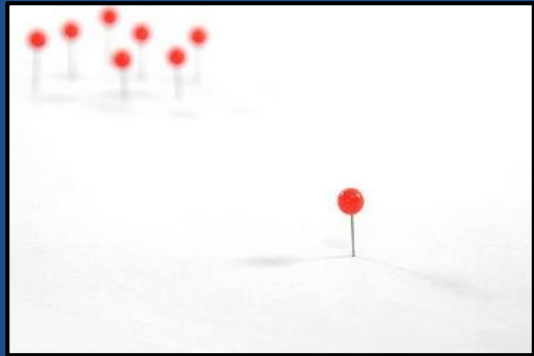
# Distributions

- Linear model assumptions:
  - Variables follow a Gaussian distribution
- Other models: no assumption
  - Better spread of values may benefit performance

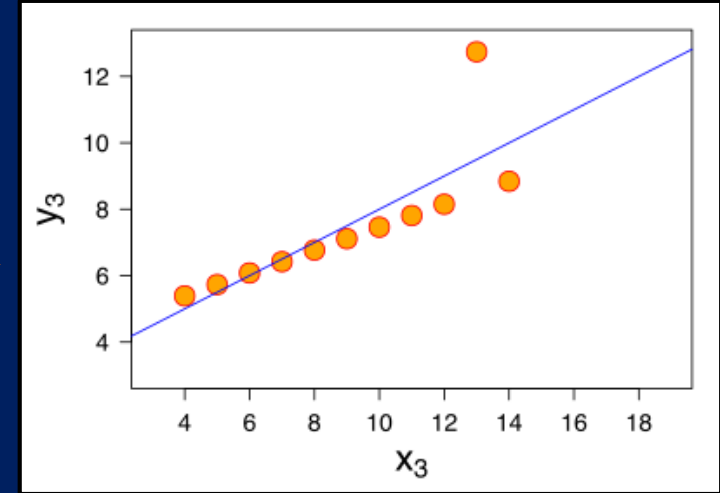
Gaussian vs Skewed



# Outliers



Linear  
models



Adaboost

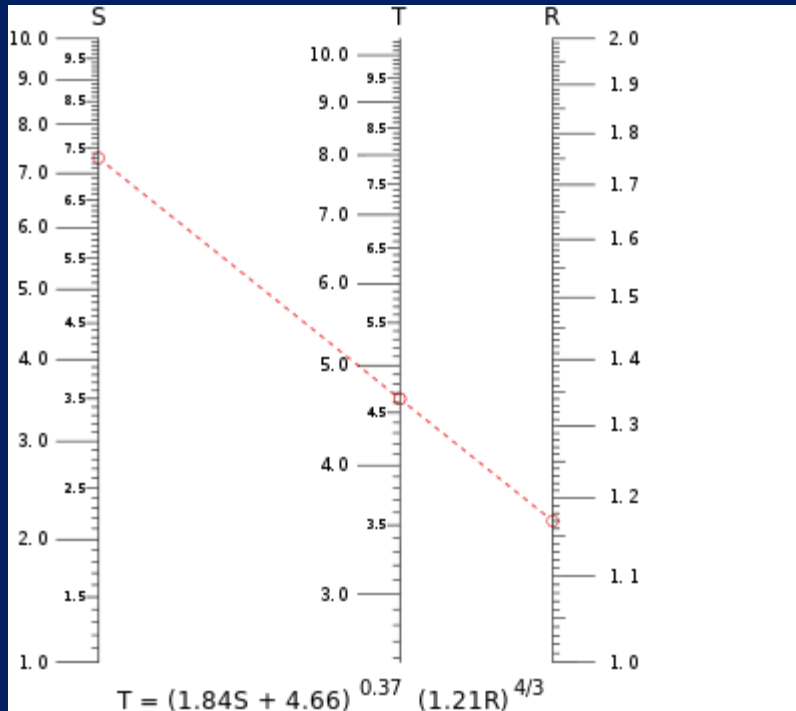


Tremendous  
weights



Bad  
generalisation

# Feature Magnitude - Scale



## Machine learning models sensitive to feature scale:

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

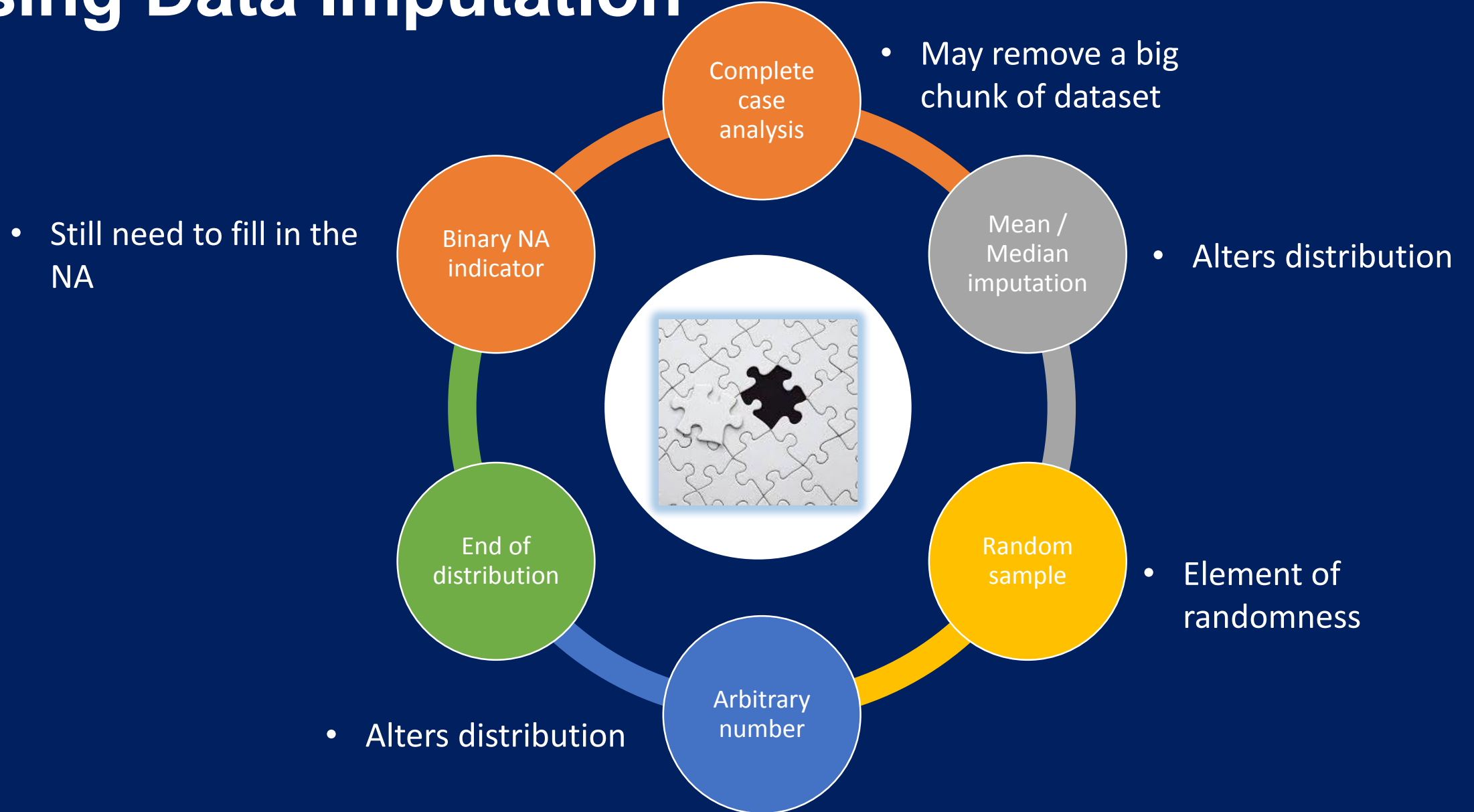
## Tree based ML models insensitive to feature scale:

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

# Data Preparation Journey

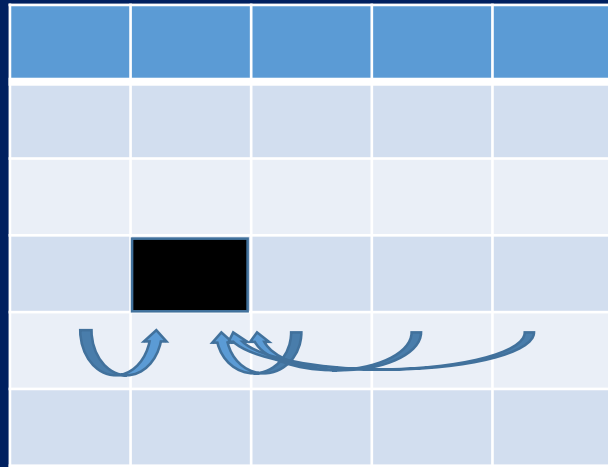
- Common issues found in variables
- Feature / variable engineering: solutions to the data issues
- Feature selection: do we need to select features?
- Feature / variable selection methods
- Overview and knowledge sources

# Missing Data Imputation





# More on Missing Data Imputation



Use neighbouring variables to predict the missing value

- KNN
- Regression

AI derived NA imputation

Complex

Insight on real variable value

Useful when only few variables with NA

Complex for production

Prone to errors

Computationally expensive

# Label Encoding

$$\text{WOE} = \ln \left( \frac{\% \text{ of non-events}}{\% \text{ of events}} \right)$$

Weight of evidence

One hot encoding

Count / frequency imputation

Ordinal encoding

Mean encoding

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

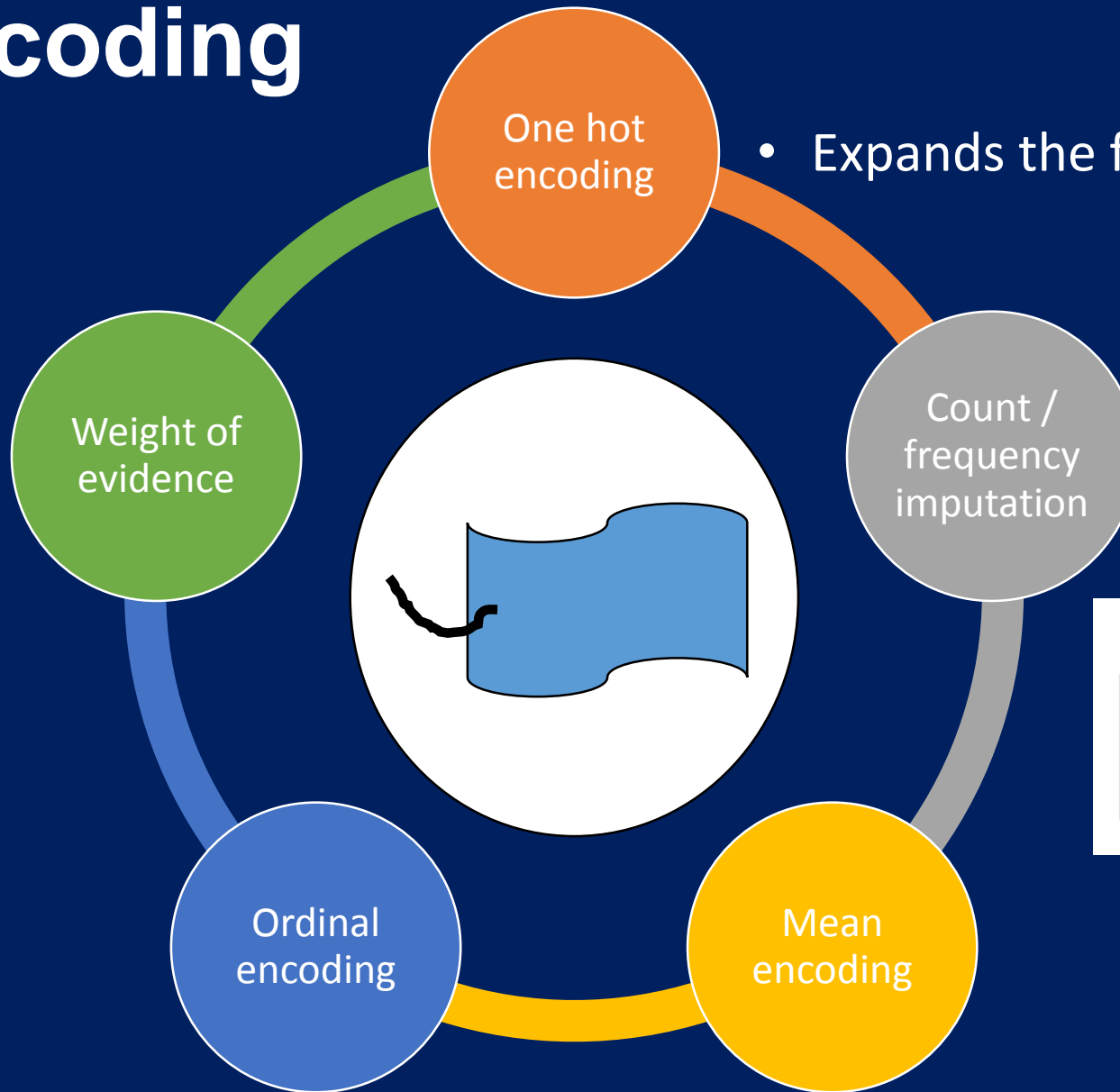
Color	Color
Red	2
Red	2
Yellow	2
Green	1
Yellow	2

Color	Target	Color
Red	0	0.5
Red	1	0.5
Yellow	1	1
Green	0	0
Yellow	1	1

Color	Target	Color
Red	0	2
Red	1	2
Yellow	1	1
Green	0	3
Yellow	1	1

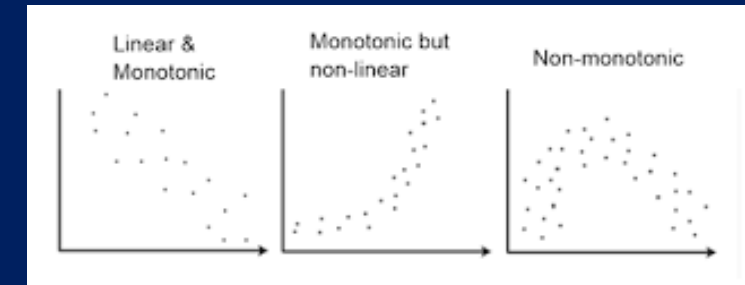
# Label Encoding

- Account for zero values as it uses logarithm



- Expands the feature space

- No monotonic relationship

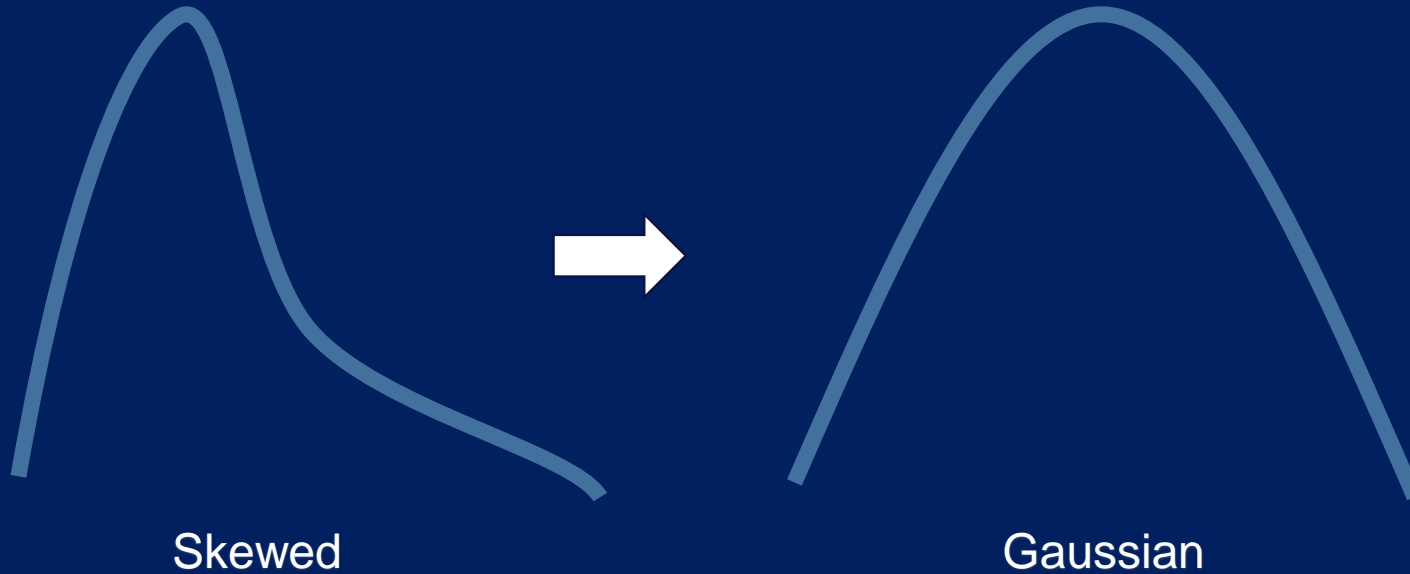


- Prone to overfitting

# Rare Labels



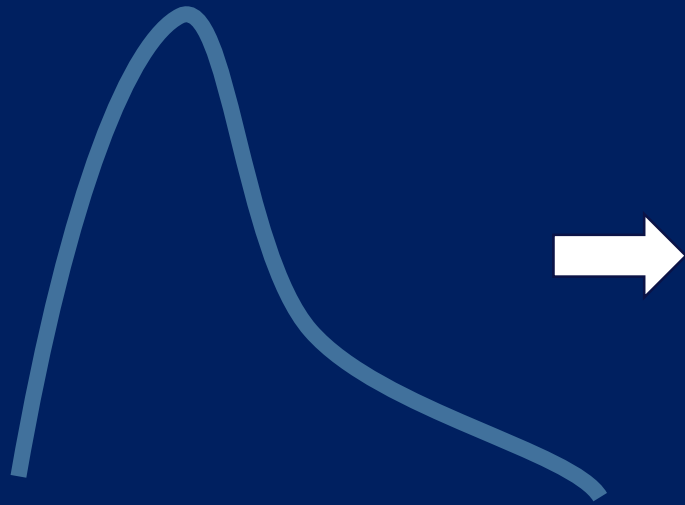
# Distribution: Gaussian Transformation



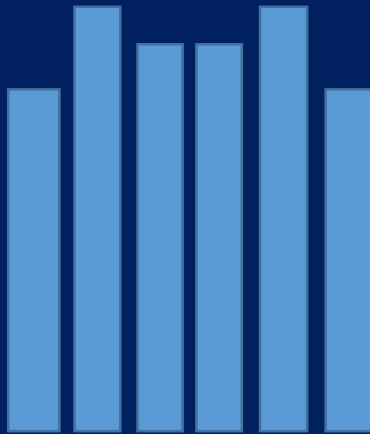
## Variable transformation

- Logarithmic  $\rightarrow \ln(x)$
- Exponential  $\rightarrow x^{\text{Exp}}$  (any power)
- Reciprocal  $\rightarrow (1 / x)$
- Box-Cox  $\rightarrow (x^{\text{Exp}(\lambda)} - 1) / \lambda$ 
  - $\lambda$  varies from -5 to 5
- Yeo-Johnson

# Distribution: Discretisation



Skewed



Improved value spread

## Discretisation

- Equal width bins
  - Bins  $\rightarrow (\text{max} - \text{min}) / n \text{ bins}$
  - Generally does not improve the spread
- Equal frequency bins
  - Bins determined by quantiles
  - Equal number of observations per bin
  - Generally improves spread
- KBins and Decision Trees

# Outliers

Trimming



- Remove the observations from dataset

Top | bottom  
coding



- Cap top and bottom values

Discretisation



- Equal bin / equal width / trees

# Data Preparation Journey

- Common issues found in variables
- Feature / variable engineering: solutions to the data issues
- Feature selection: do we need to select features?
- Feature / variable selection methods
- Overview and knowledge sources



# Why Do We Select Features?

- Simple models are easier to interpret
- Shorter training times
- Enhanced generalisation by reducing overfitting
- Easier to implement by software developers → Model production
- Reduced risk of data errors during model use
- Data redundancy

# Variable Redundancy



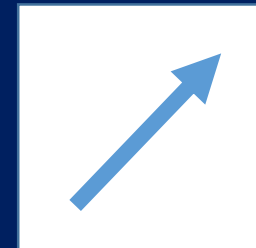
**Constant variables**  
Only 1 value per  
variable



**Quasi – constant  
Variables**  
> 99% of observations  
show same value



**Duplication**  
Same variable multiple  
times in the dataset

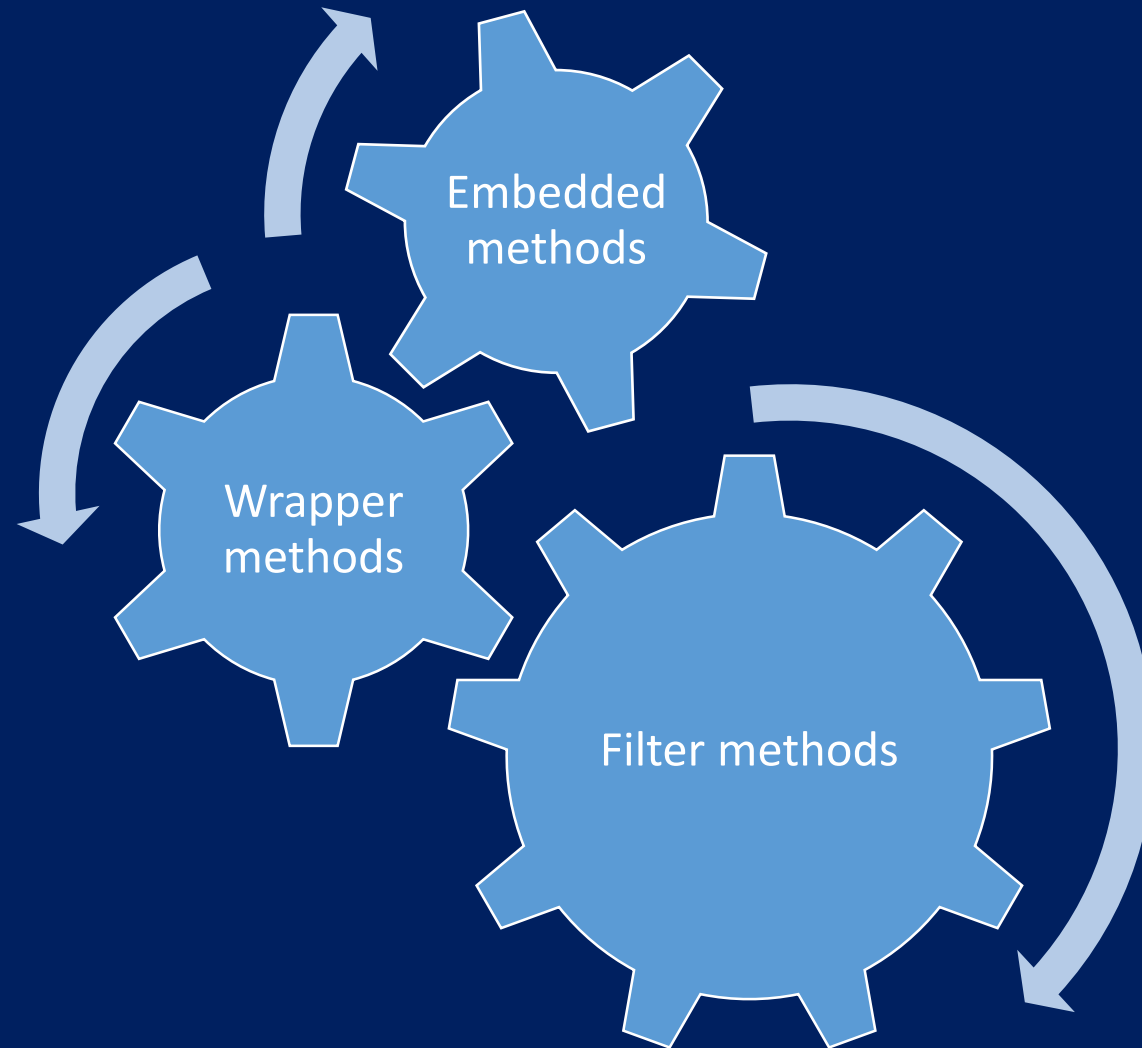


**Correlation**  
Correlated variables  
provide the same  
information

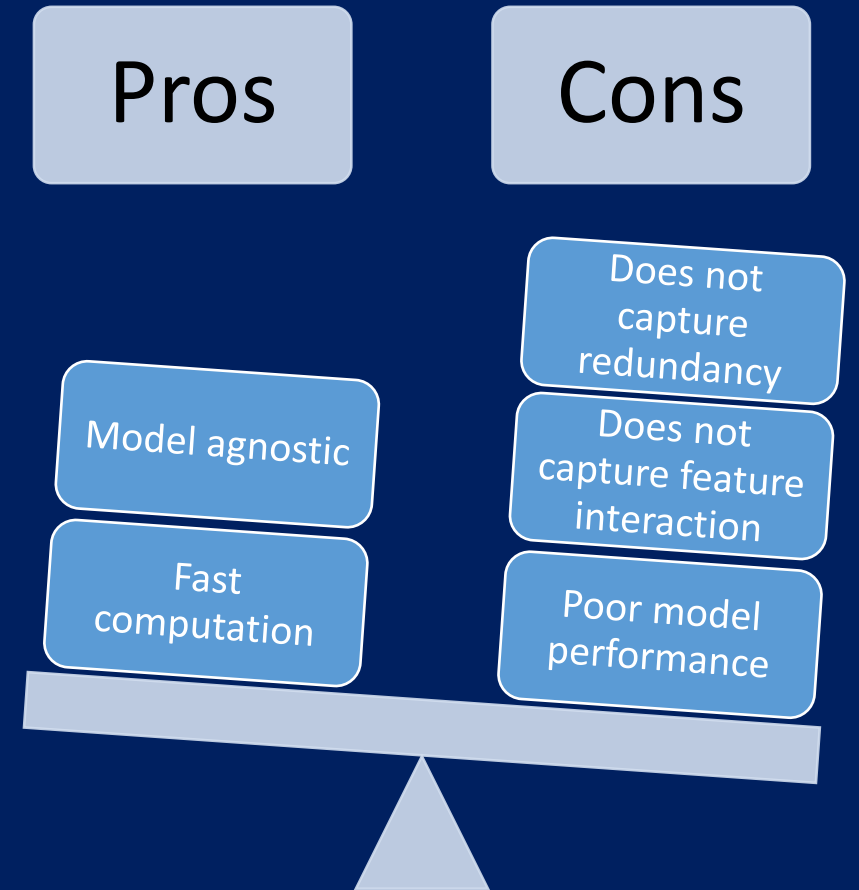
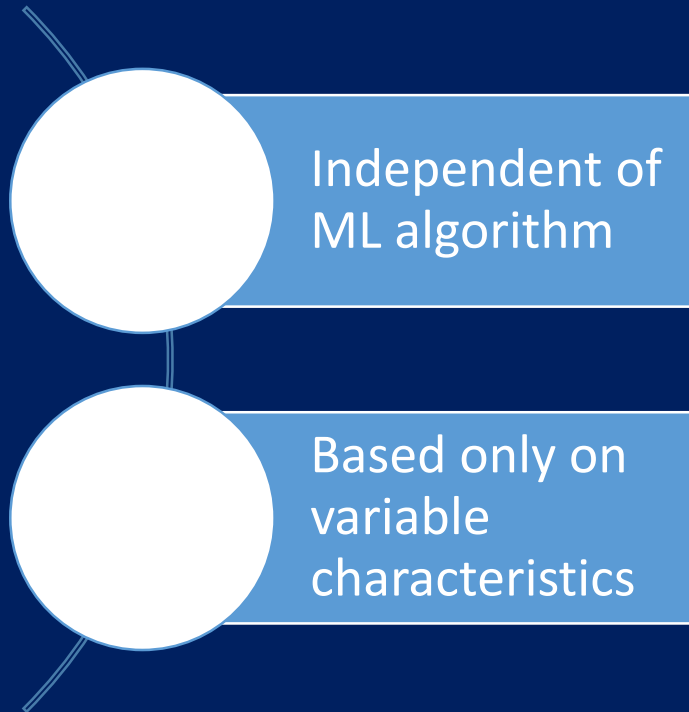
# Data Preparation Journey

- Common issues found in variables
- Feature / variable engineering: solutions to the data issues
- Feature selection: do we need to select features?
- Feature / variable selection methods
- Overview and knowledge sources

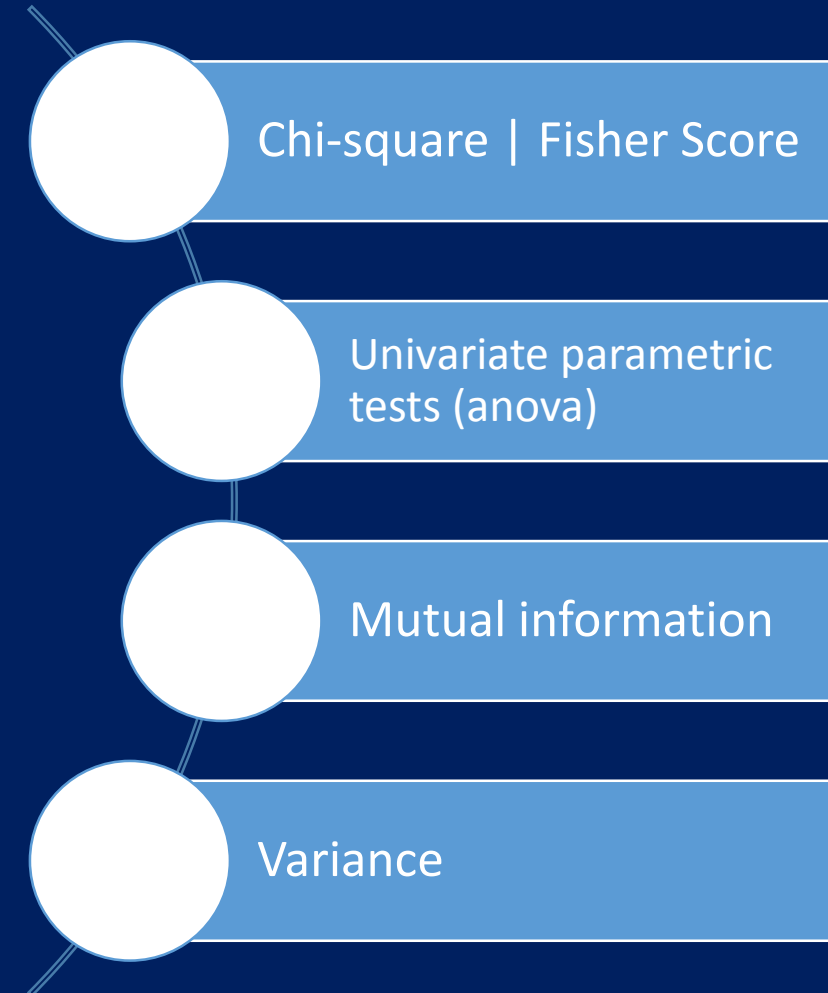
# Feature Selection Methods



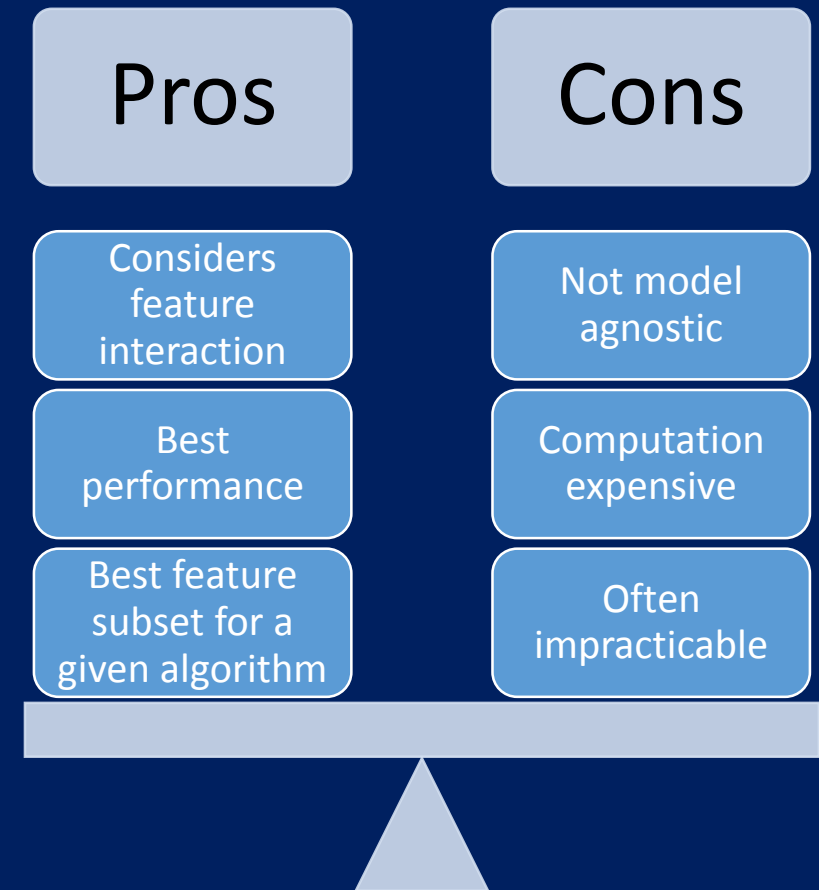
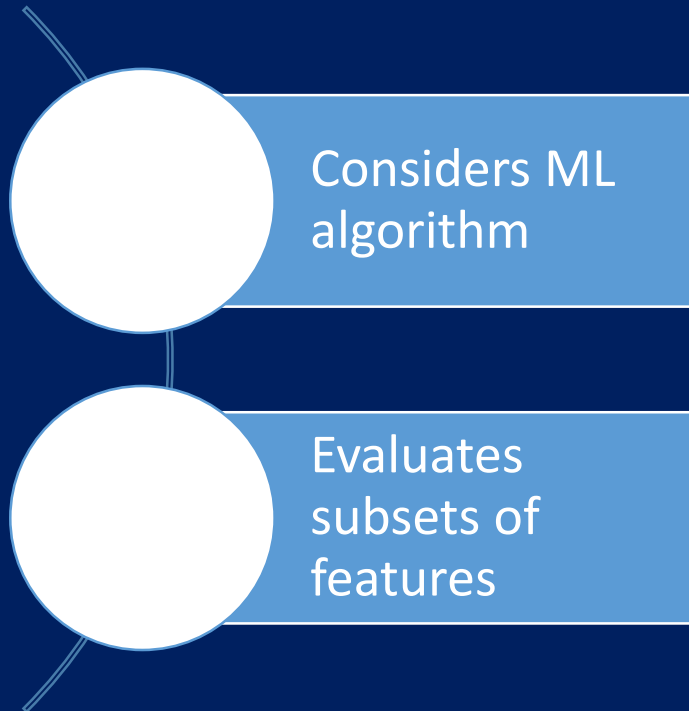
# Filter methods



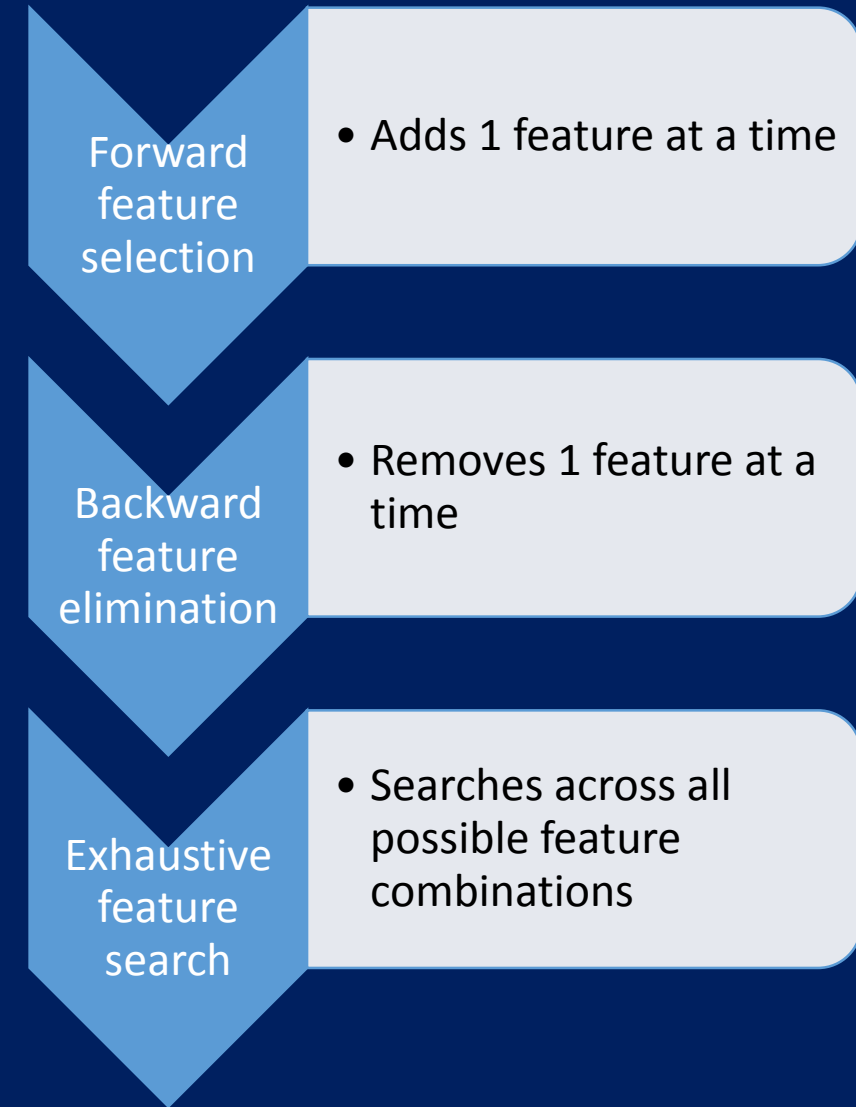
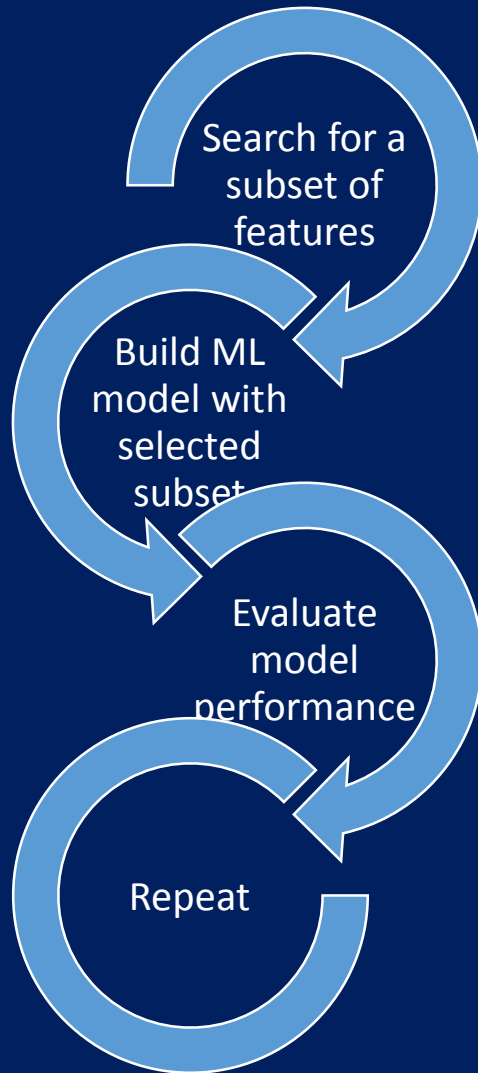
# Filter methods



# Wrapper methods

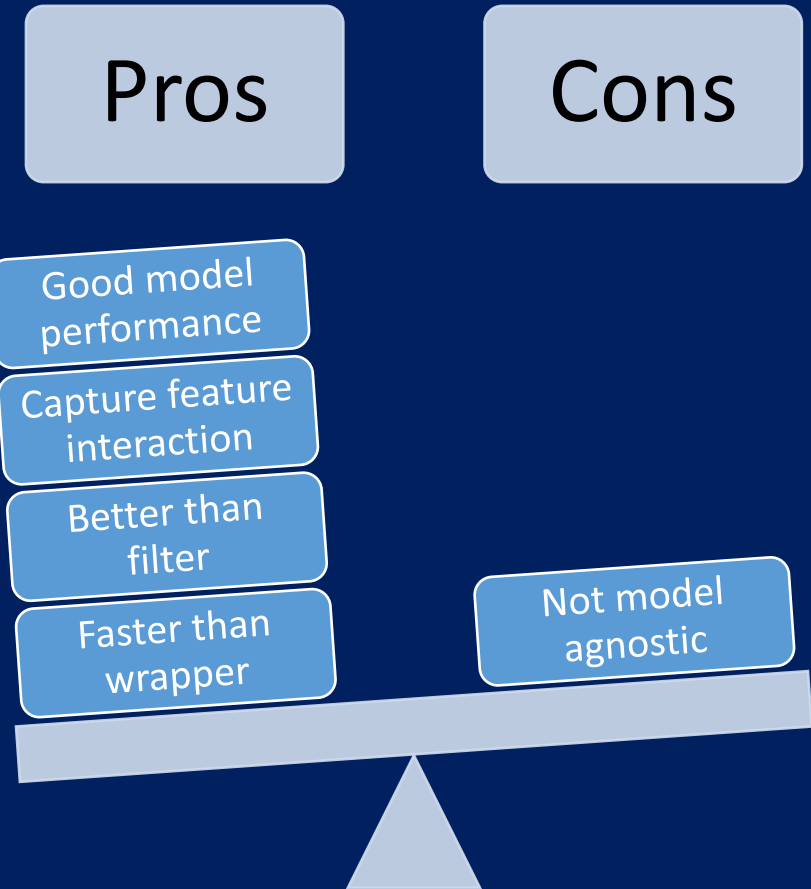
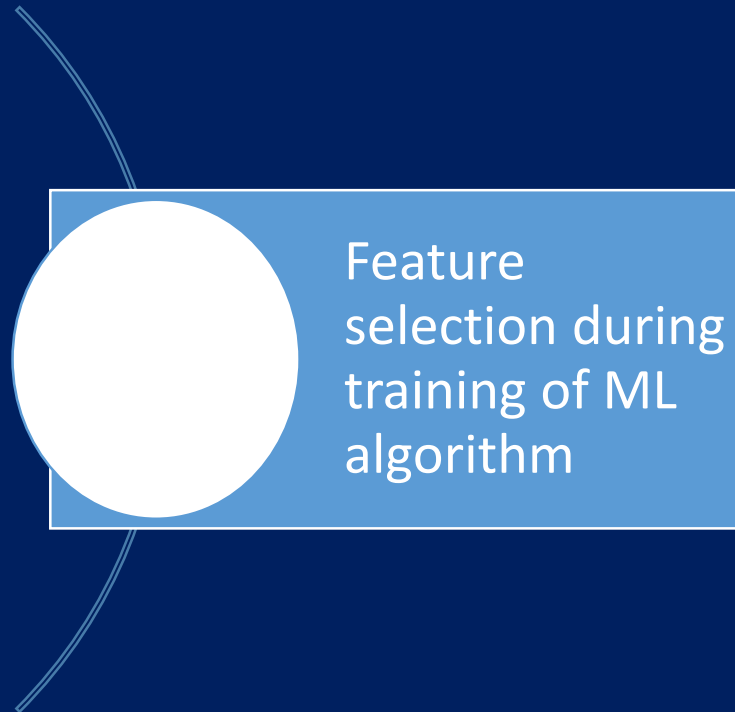


# Wrapper methods

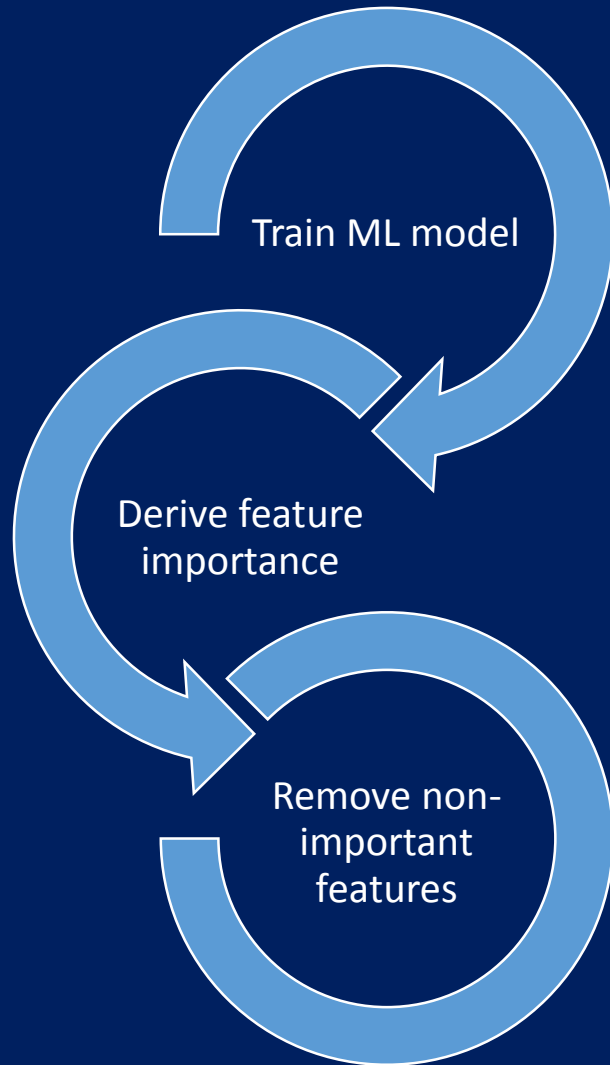




# Embedded methods



# Embedded methods



# Data Preparation Journey

- Common issues found in variables
- Feature / variable engineering: solutions to the data issues
- Feature selection: do we need to select features?
- Feature / variable selection methods
- Overview and knowledge sources

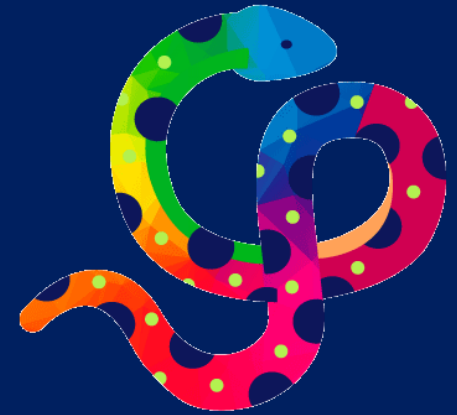
# How can we engineer features?



- + Fast computation
- + Cross-validation
- Little versatility to select features

## 🏠 Category Encoders

- + Alternative encoding procedures
- Bad for interpretability



## Feature-Engine

- + More engineering steps
- + Can apply to subset of features
- Need to decide step a priori

# How can we select features?



+ Filter and embedded methods



+ Wrapper methods  
- Slow

# How can we learn more?

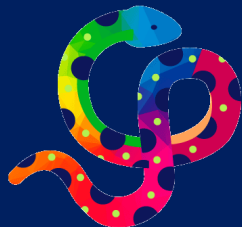
## The 2009 Knowledge Discovery in Data Competition (KDD Cup 2009) Challenges in Machine Learning, Volume 3

Gideon Dror, Marc Boullé, Isabelle Guyon,  
Vincent Lemaire, and David Vogel, editors

Summary of learnings from the winners



Documentation



How to Win a Data  
Science Competition:  
Learn from Top Kagglers



Udemy.com, includes code

