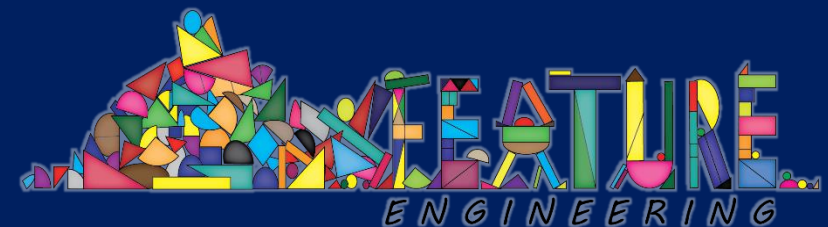


Engineering and selecting features for machine learning

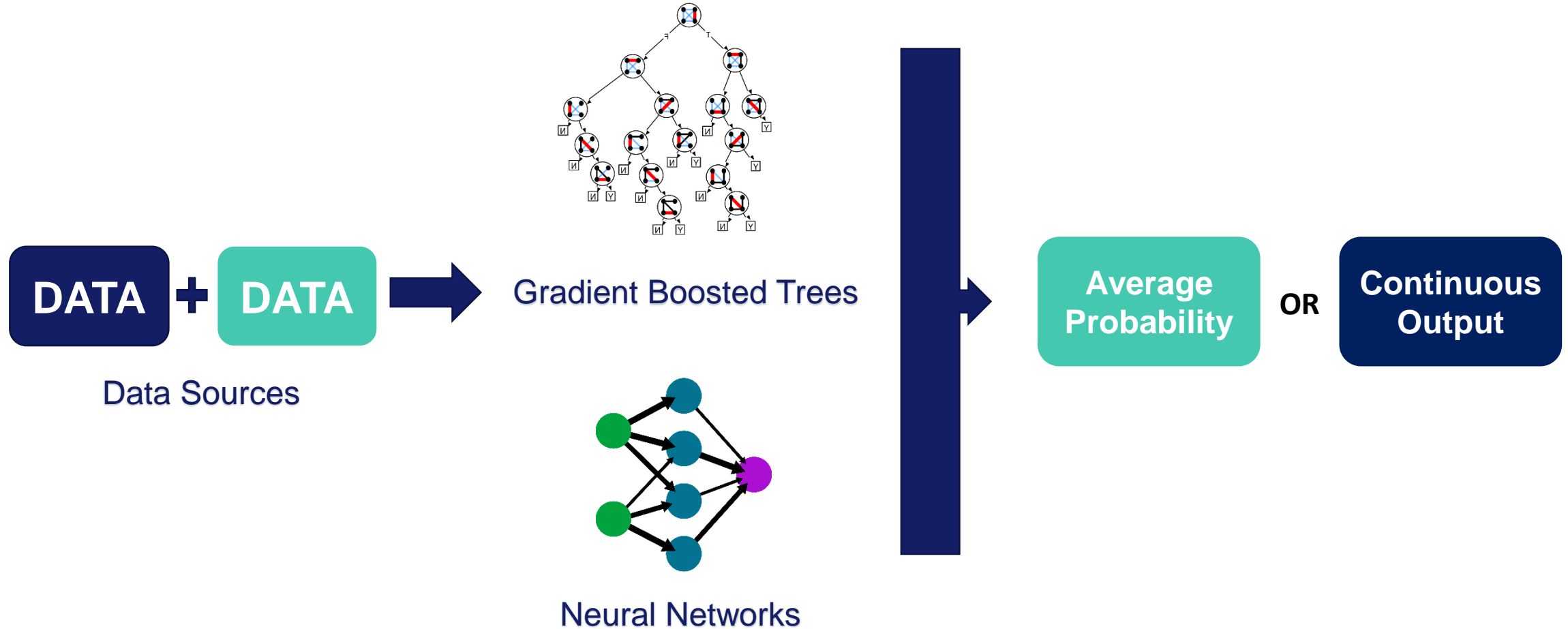
Soledad Galli, PhD

DSF meetup with Busuu

London, 16th October 2018



Machine Learning Pipeline



Machine Learning Finance and Insurance



Claims



Fraud



Credit Risk

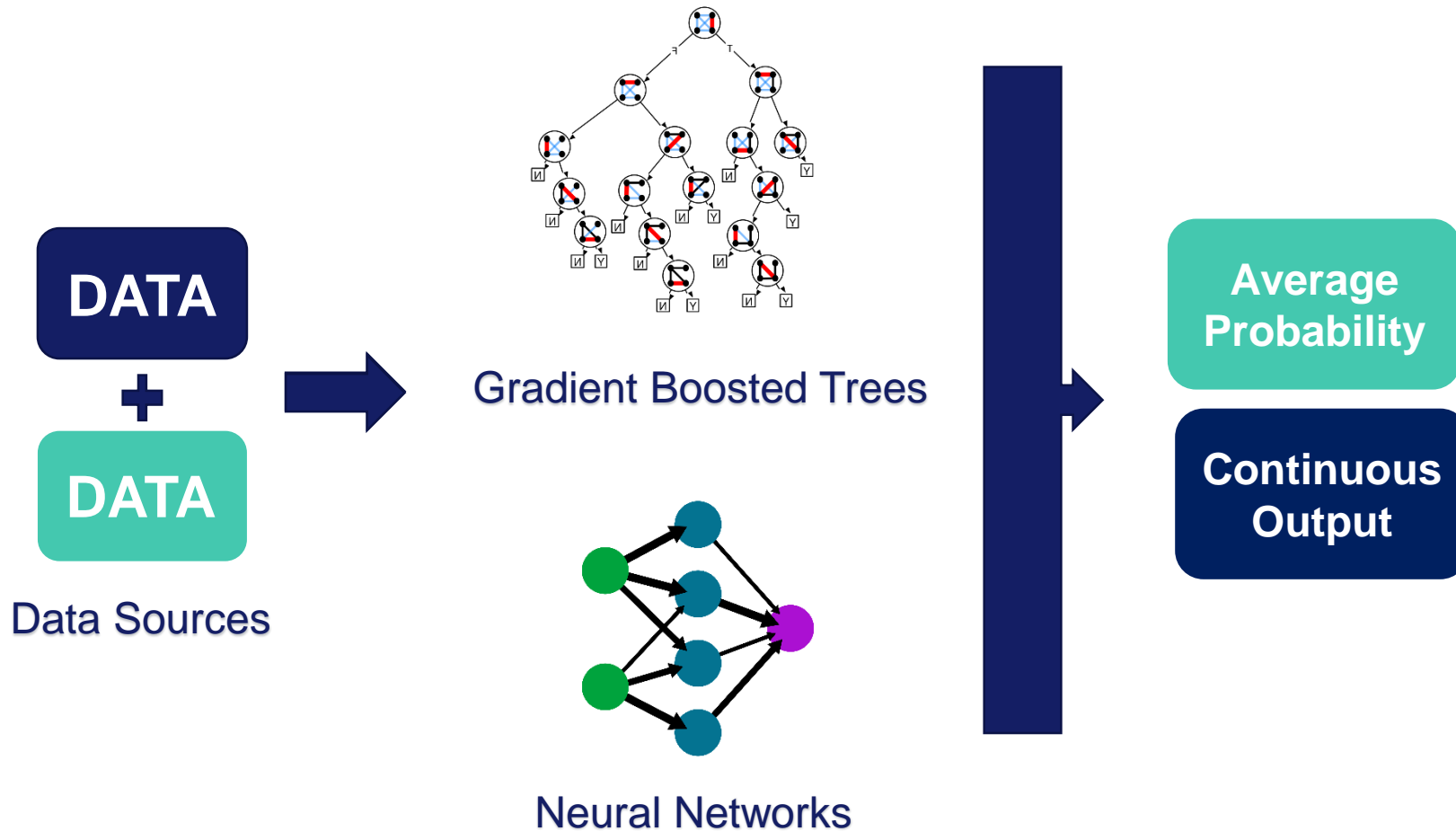


Marketing

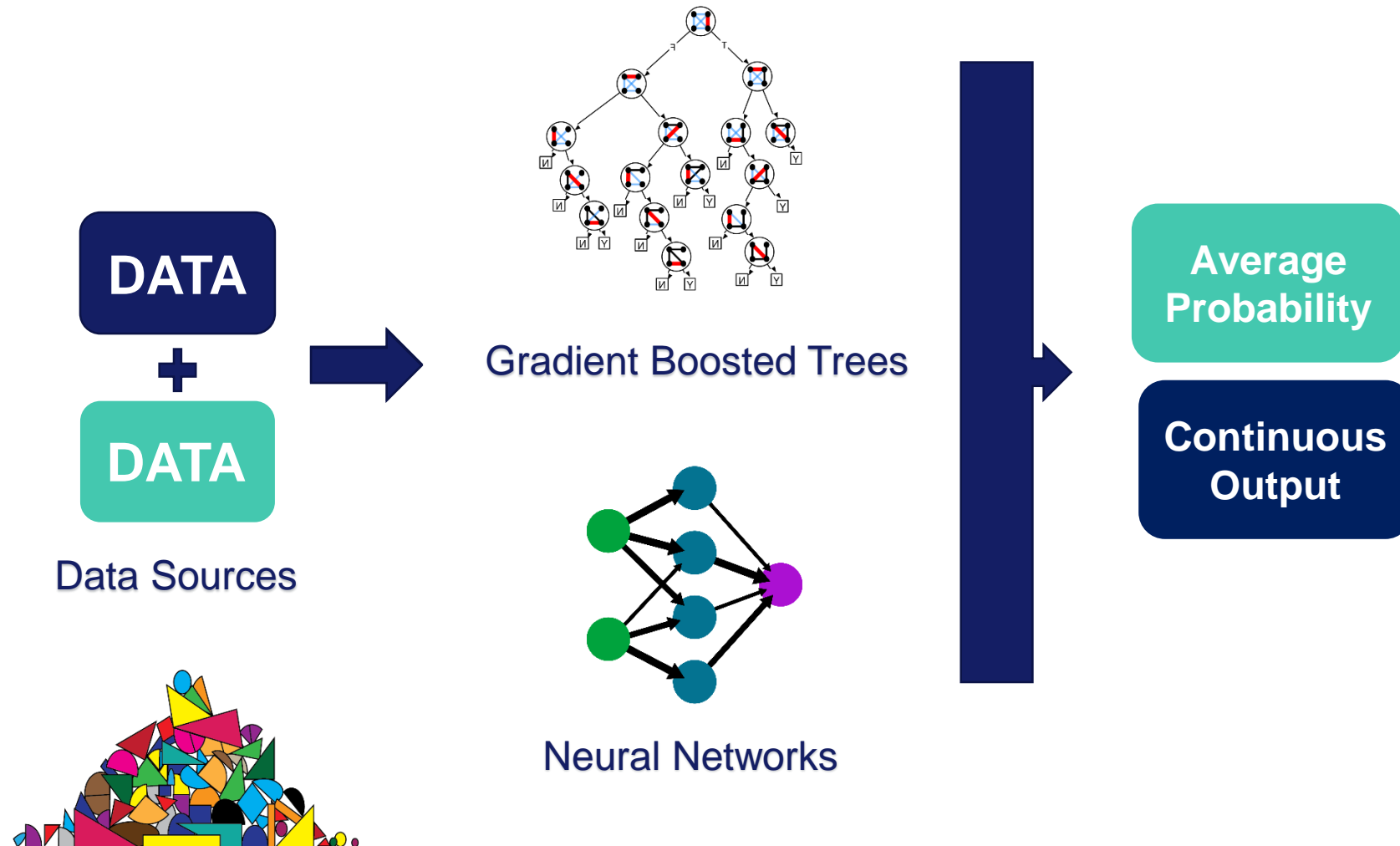


Pricing

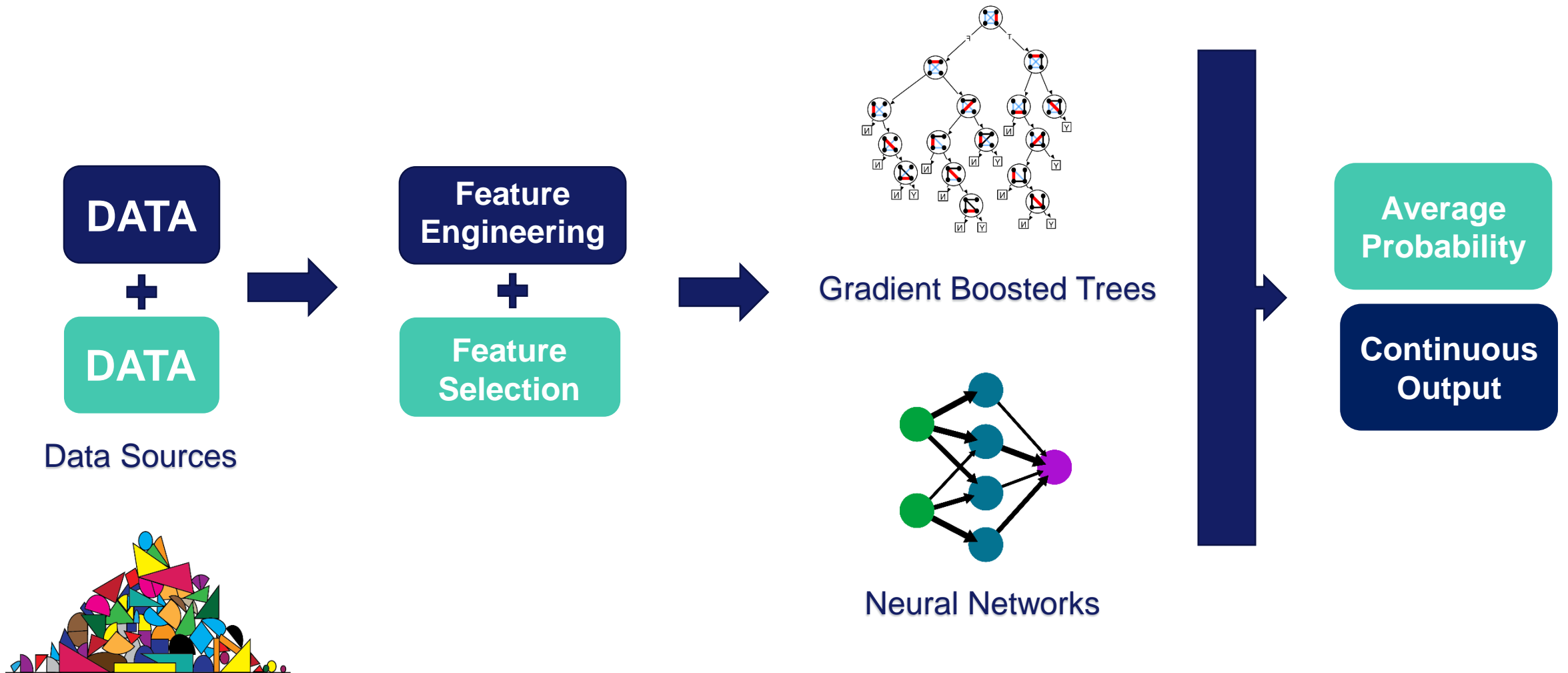
Machine Learning Pipeline



Machine Learning Pipeline



Machine Learning Pipeline



Data Pre-processing 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 Pre-processing Journey

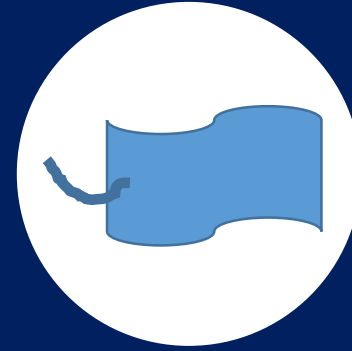
- Common issues found in variables
- Feature / Variable engineering: solutions to the data issues
- Feature selection: do we need to select features?
- Feature / Variable selection methods
- Overview and knowledge sources

Problems in Variables



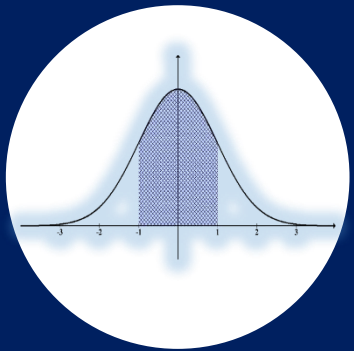
Missing data

Missing values within a variable



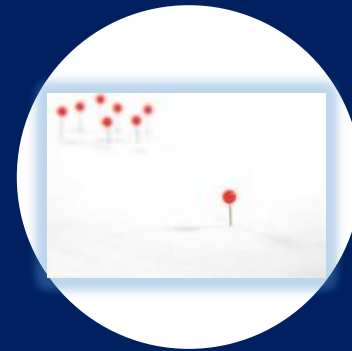
Labels

Strings in categorical variables



Distribution

Normal vs skewed



Outliers

Unusual or unexpected values

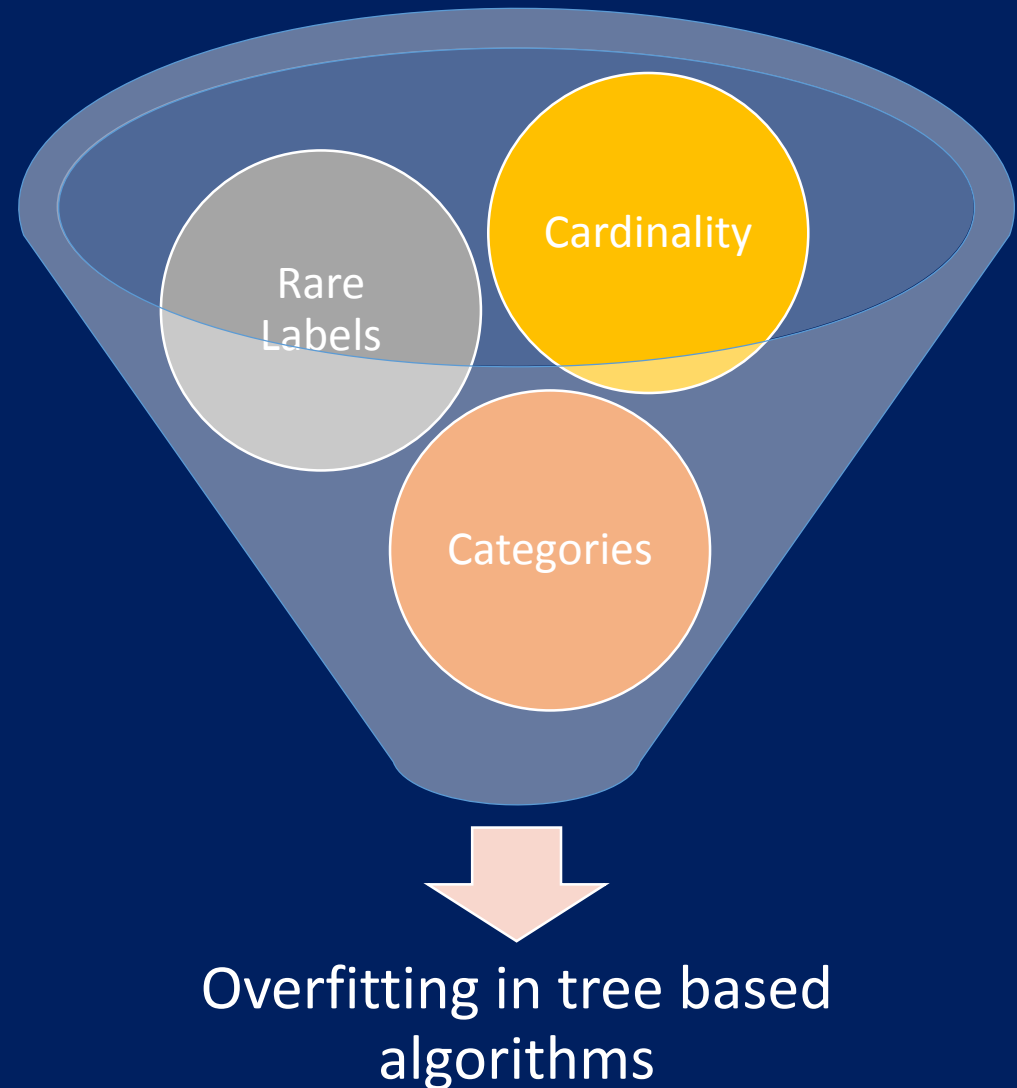
Missing Data

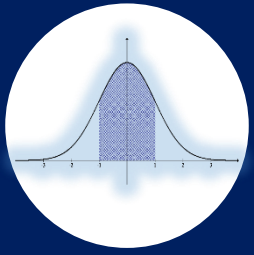
- Missing values for certain observations
- Affects all machine learning models
 - Scikit-learn



Labels in categorical variables

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

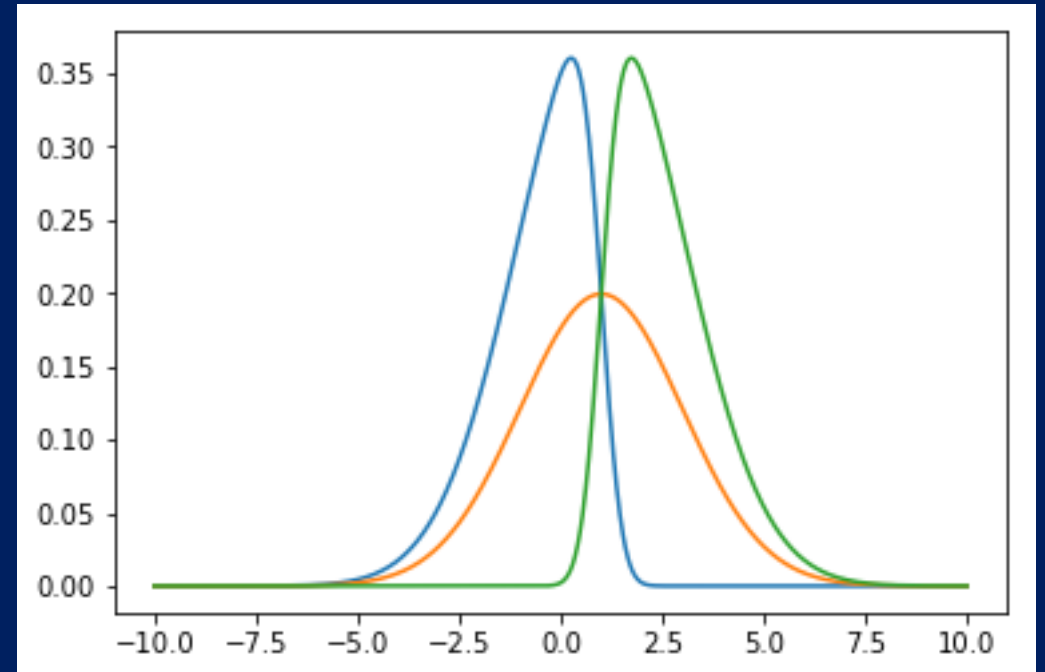




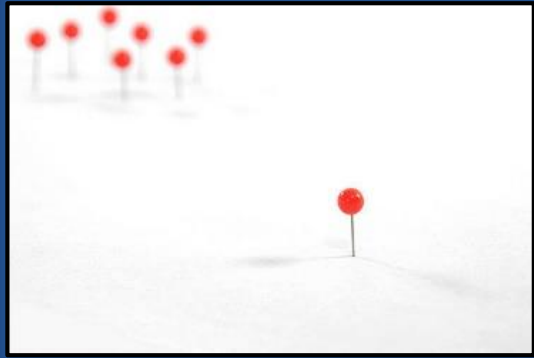
Distributions

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

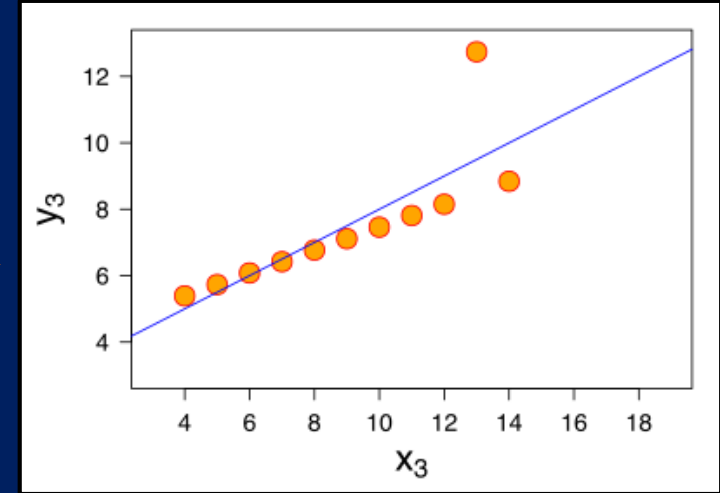
Gaussian vs Skewed



Outliers



Linear
models



Adaboost

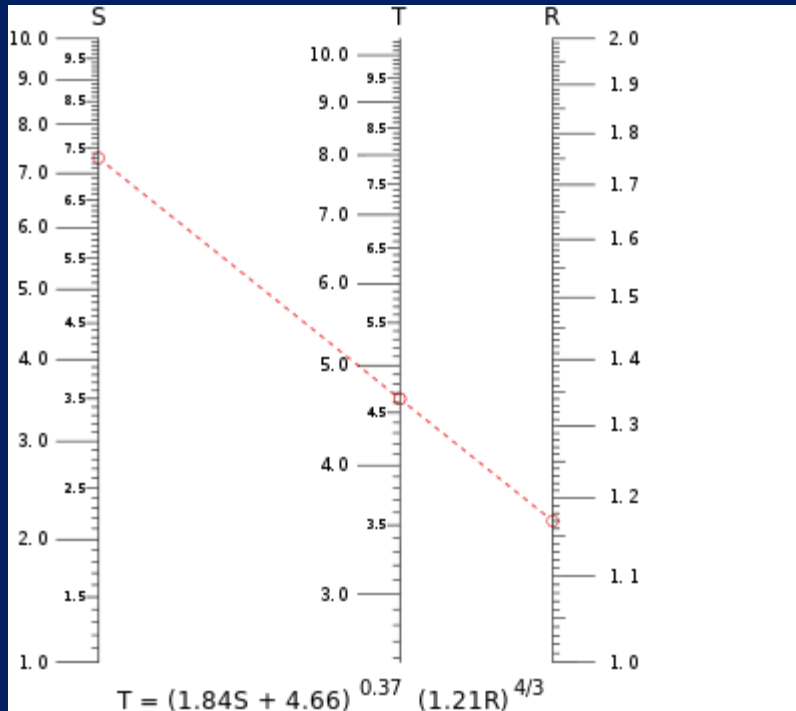


Tremendous
weights



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)

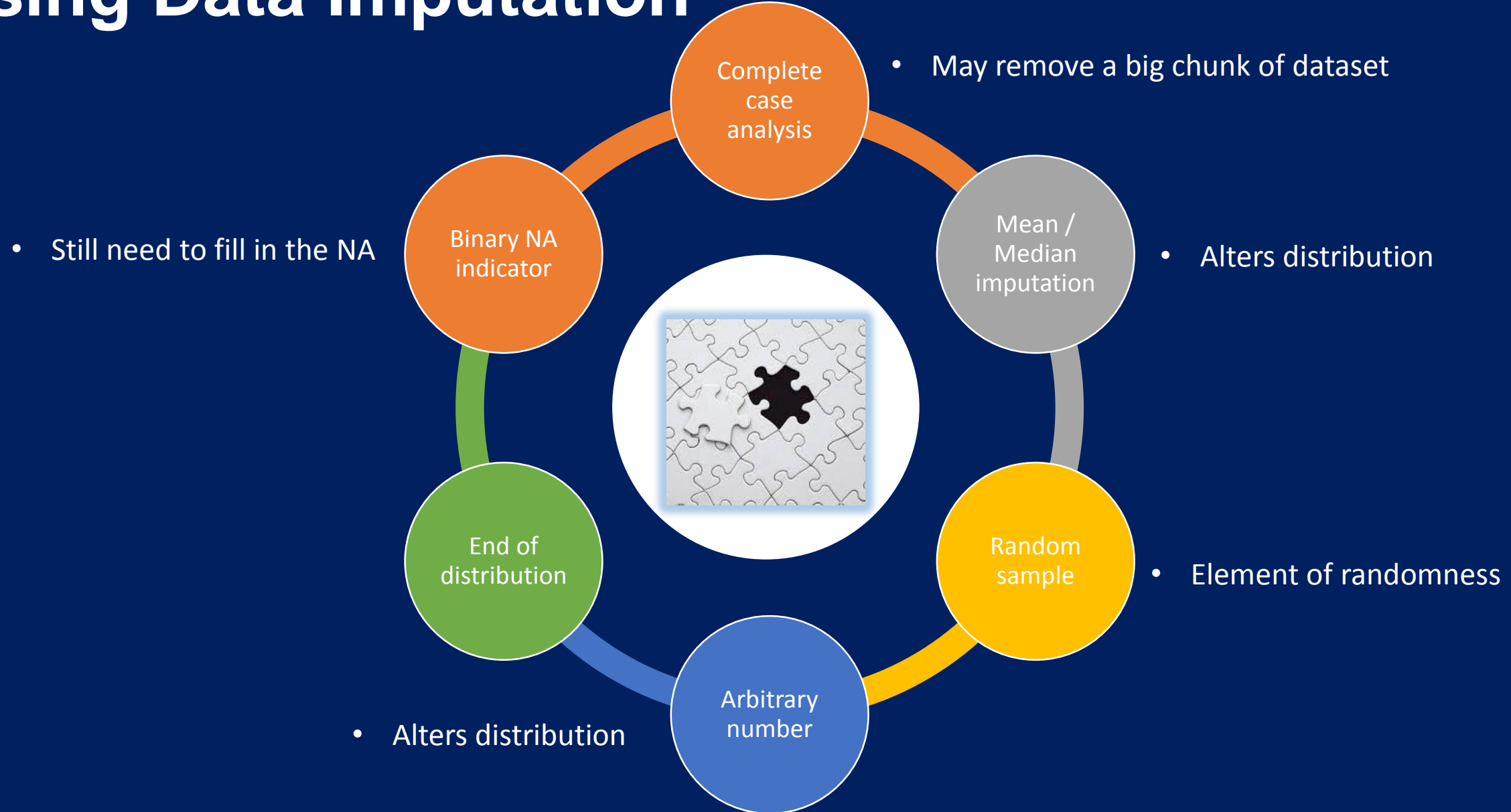
Tree based ML models insensitive to feature scale:

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

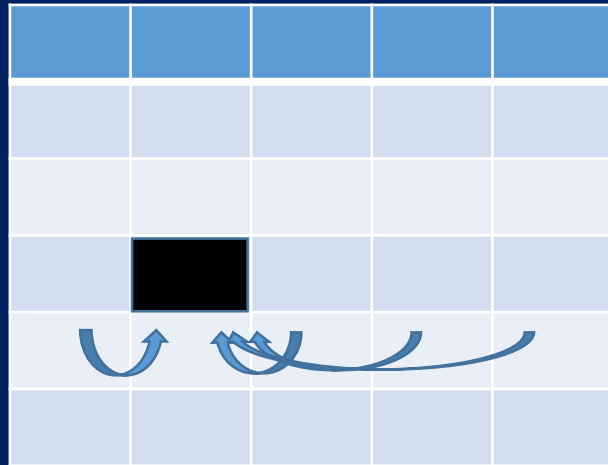
Data Pre-processing Journey

- 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

Improved model performance

Complex for productionisation

Prone to errors

Computationally expensive

Time consuming: 1 model per variable

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

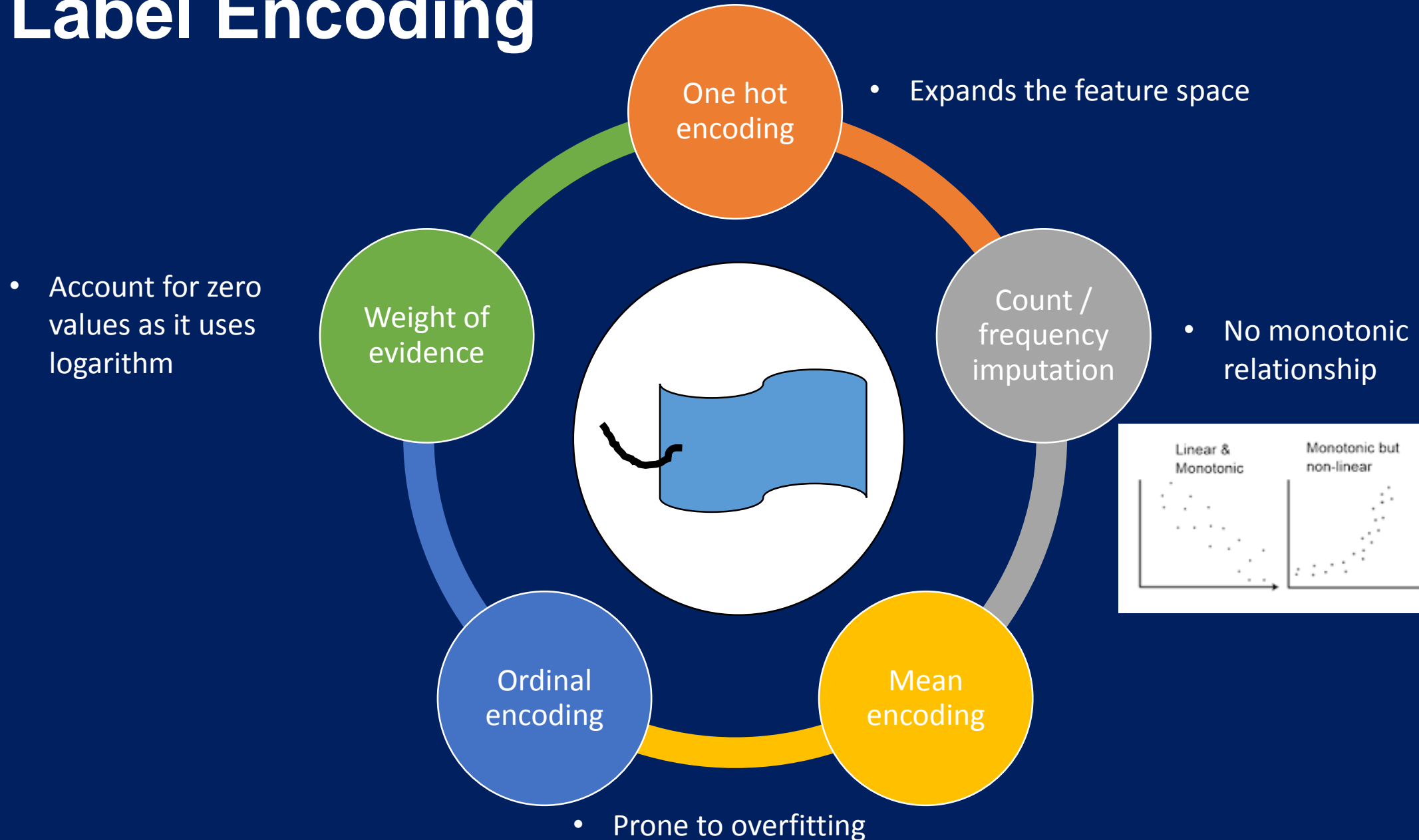
Color	Red	Yellow	Green
Red	1	0	0
Red	1	0	0
Yellow	0	1	0
Green	0	0	1
Yellow	0	0	1

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

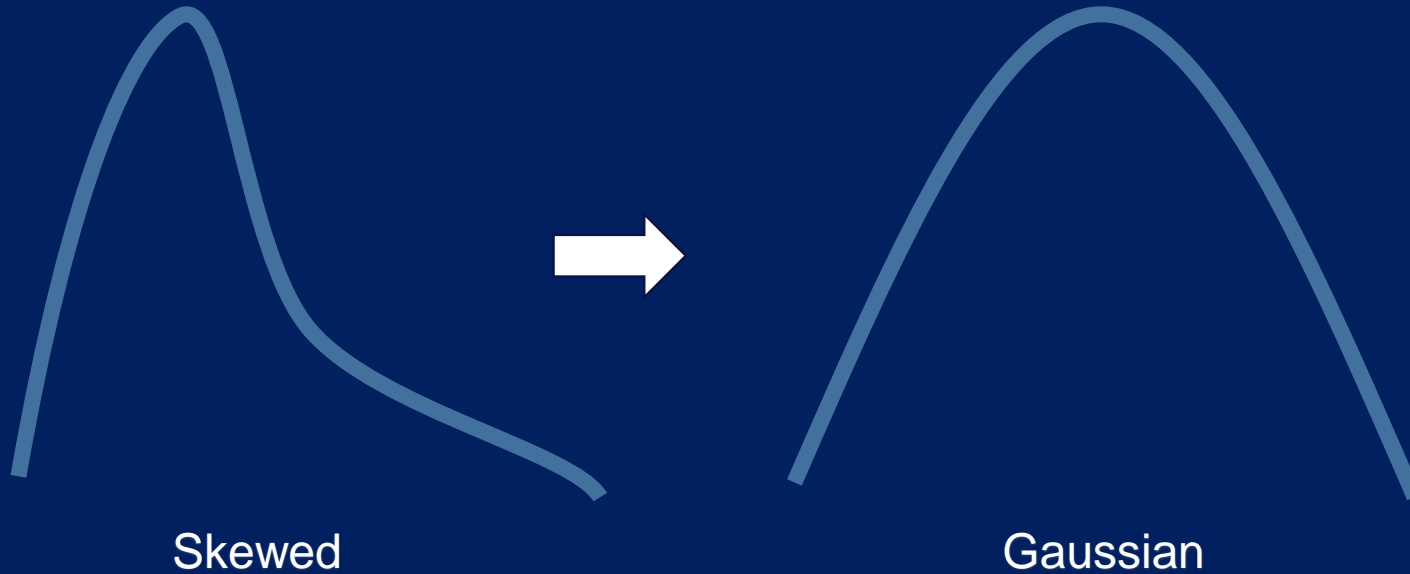
Label Encoding



Rare Labels



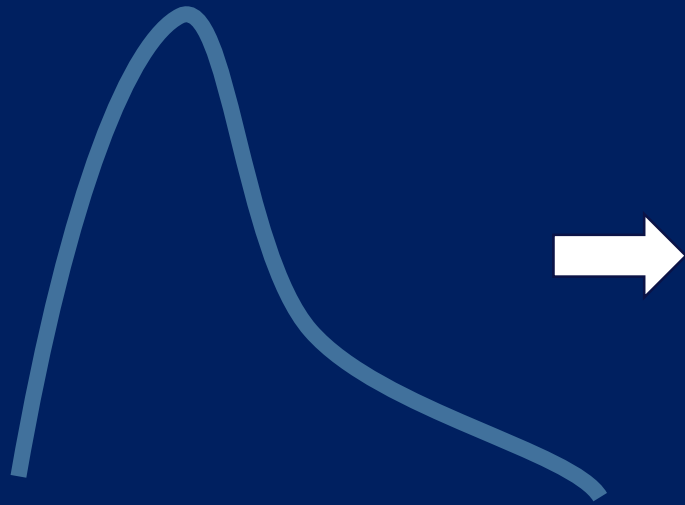
Distribution: Gaussian Transformation



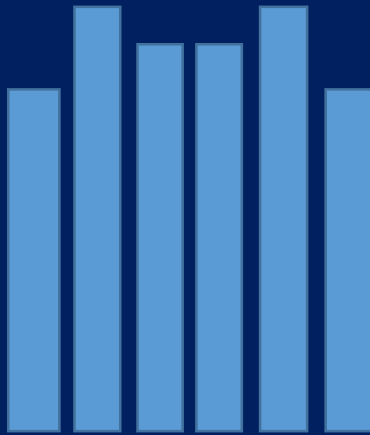
Variable transformation

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

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

Outliers

Trimming



- Remove the observations from dataset

Top | bottom
coding



- Cap top and bottom values

Discretisation



- Equal bin / equal width

Data Pre-processing 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



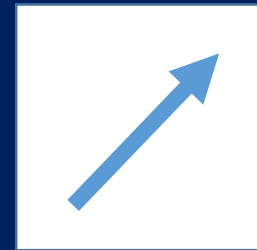
Constant variables
Only 1 value per variable



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



Duplication
Same variable multiple times in the dataset

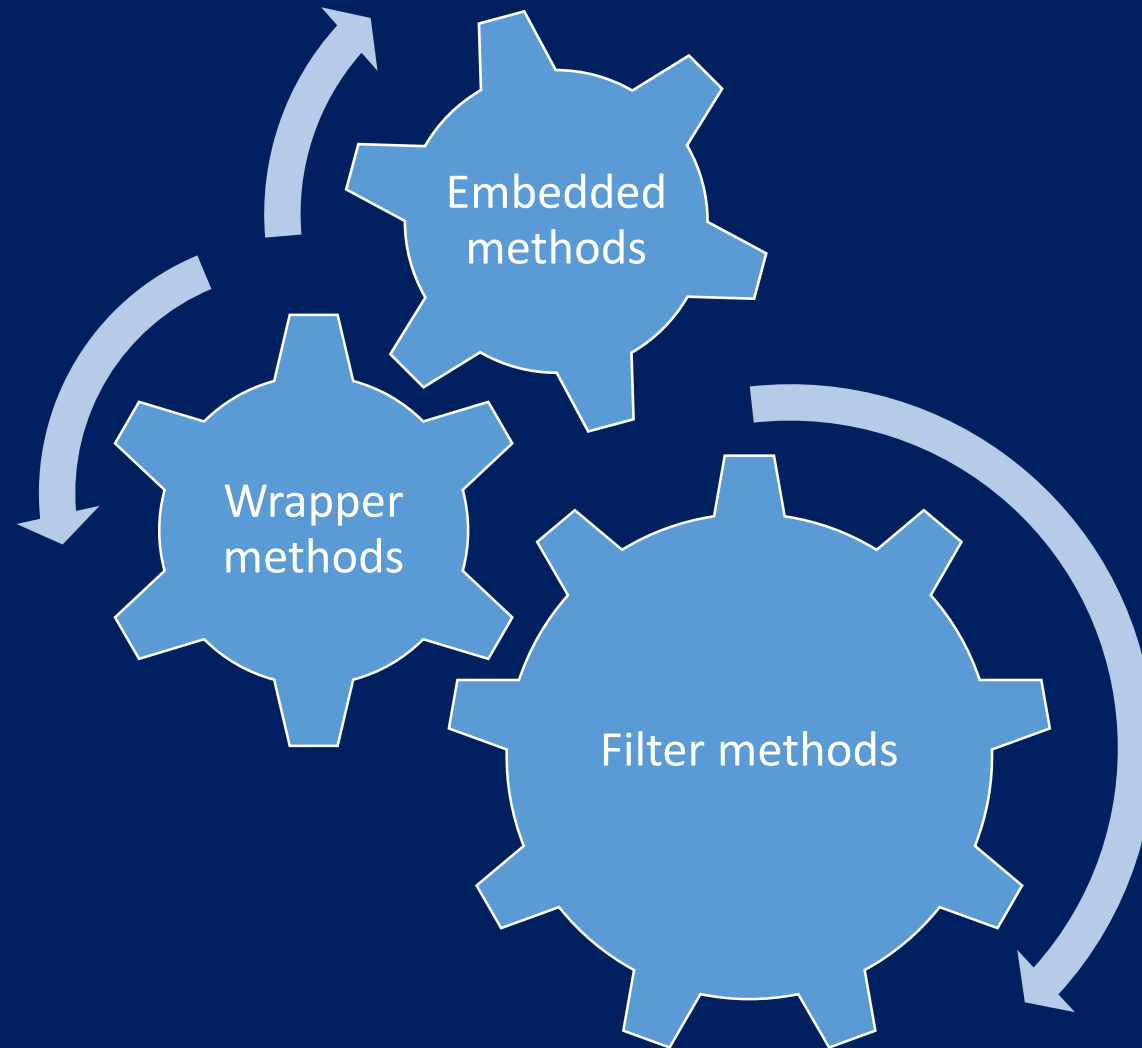


Correlation
Correlated variables provide the same information

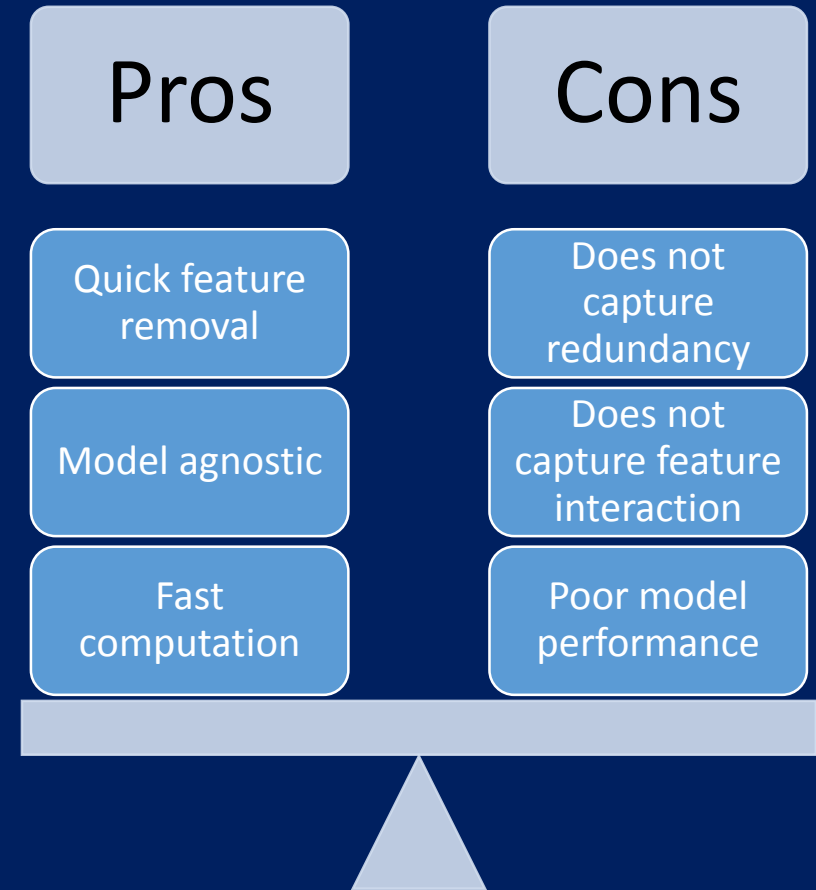
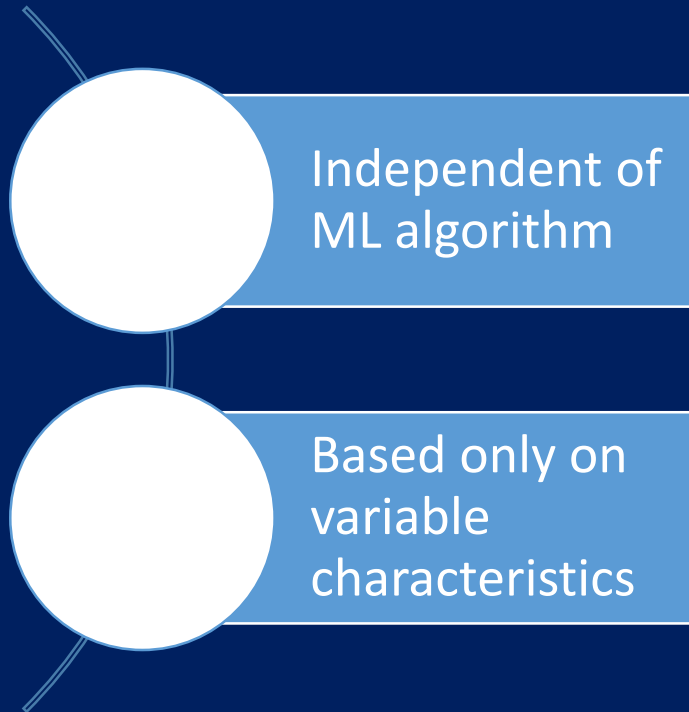
Data Pre-processing Journey

- 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

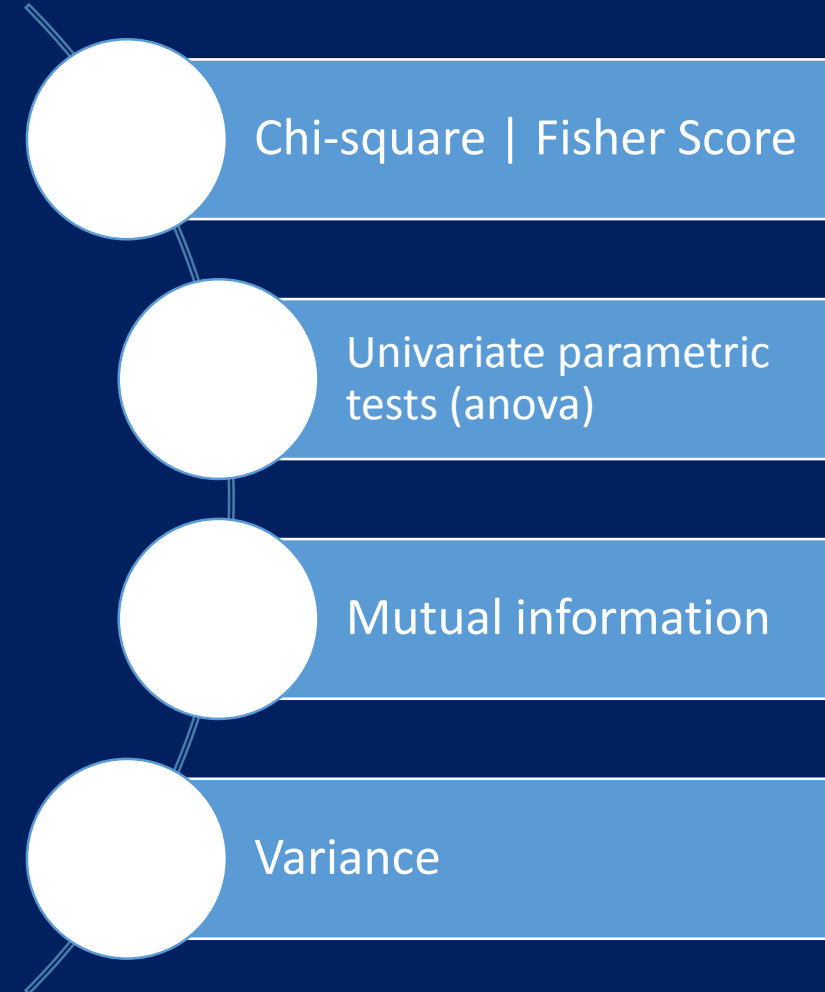
Feature Selection Methods



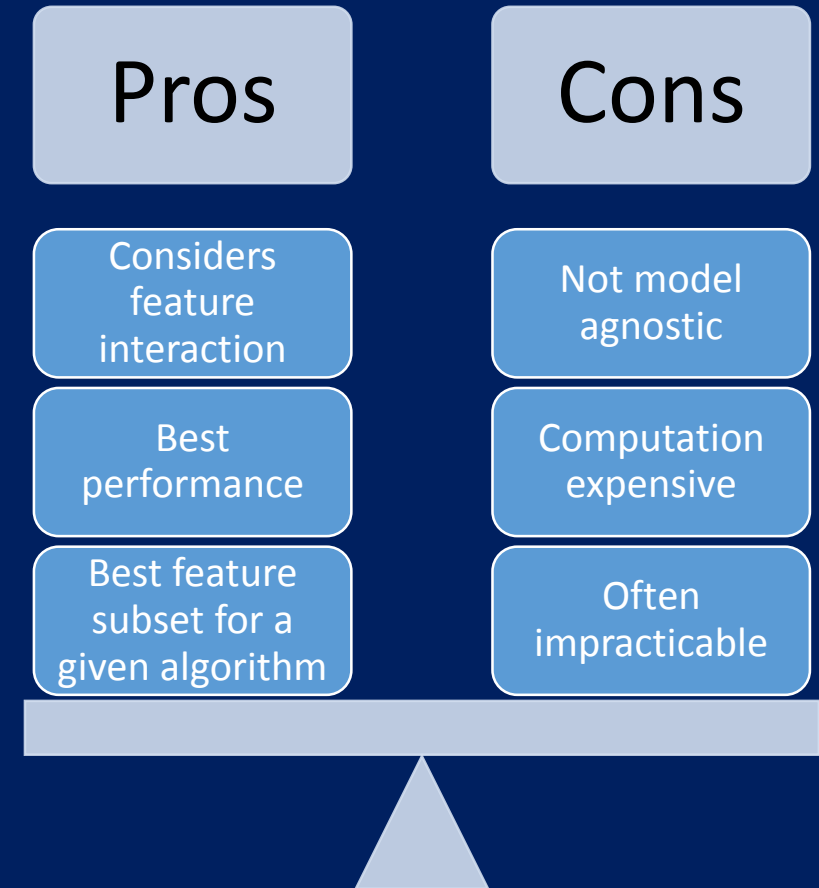
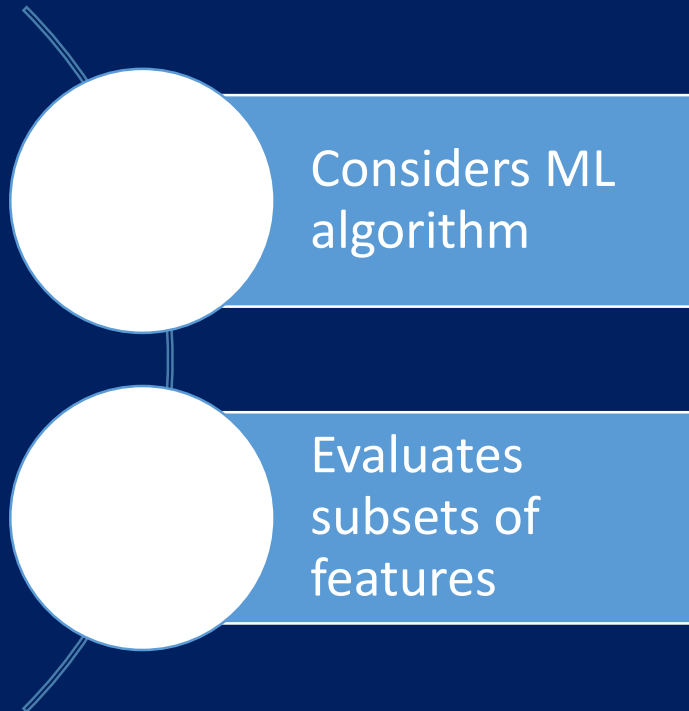
Filter methods



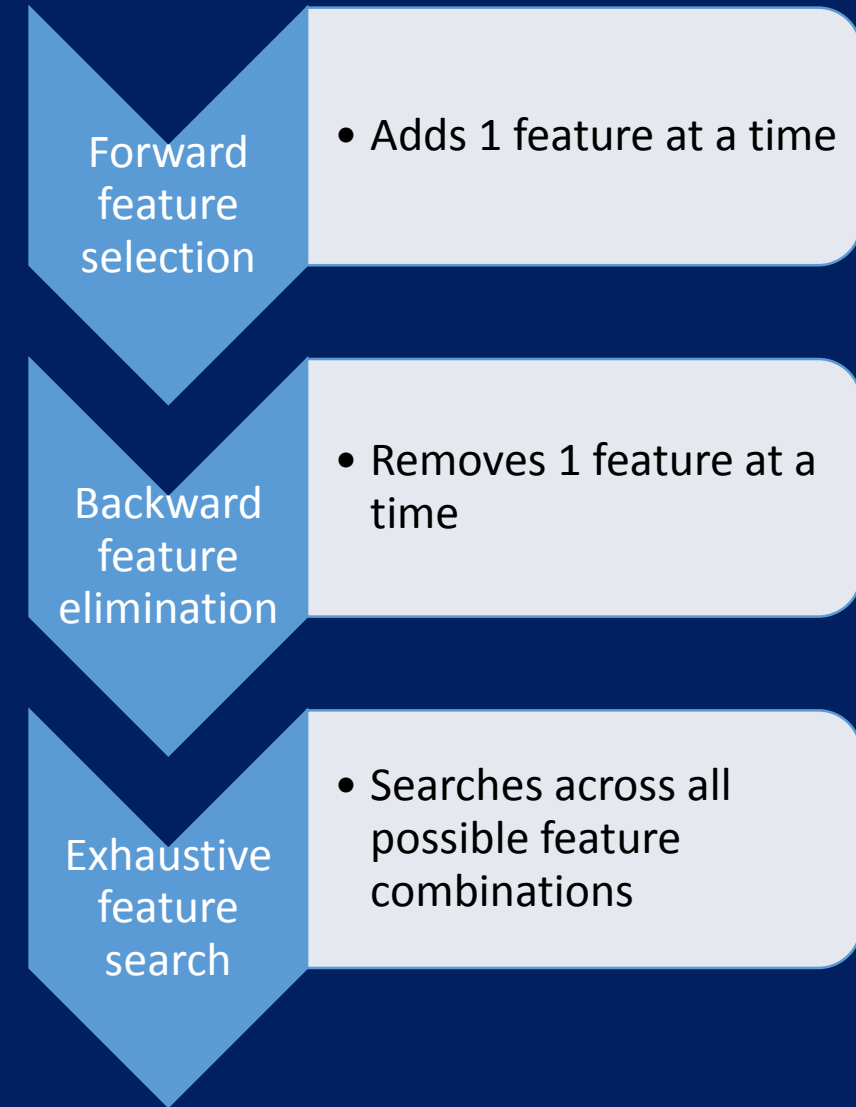
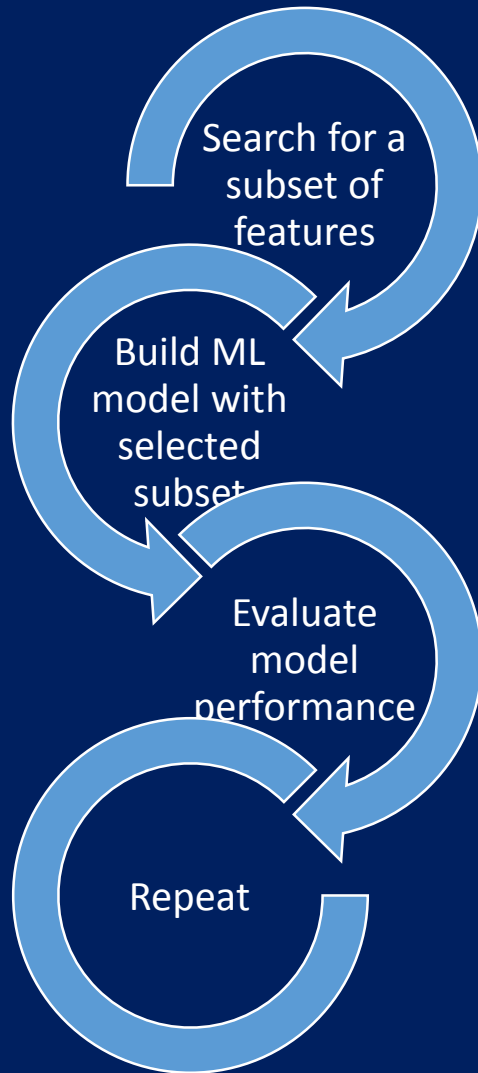
Filter methods



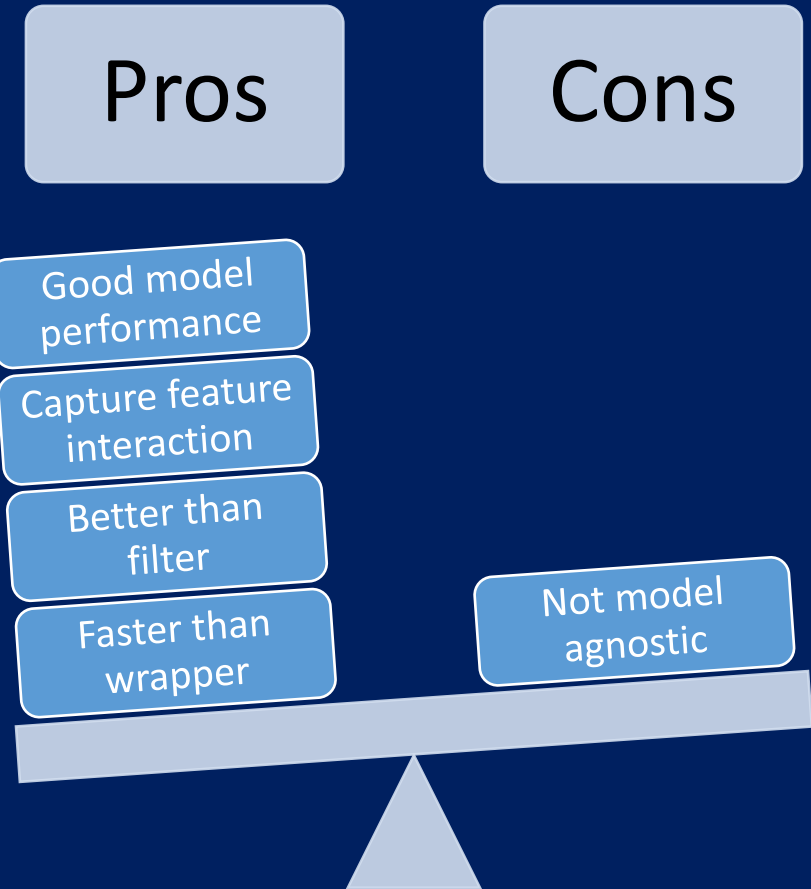
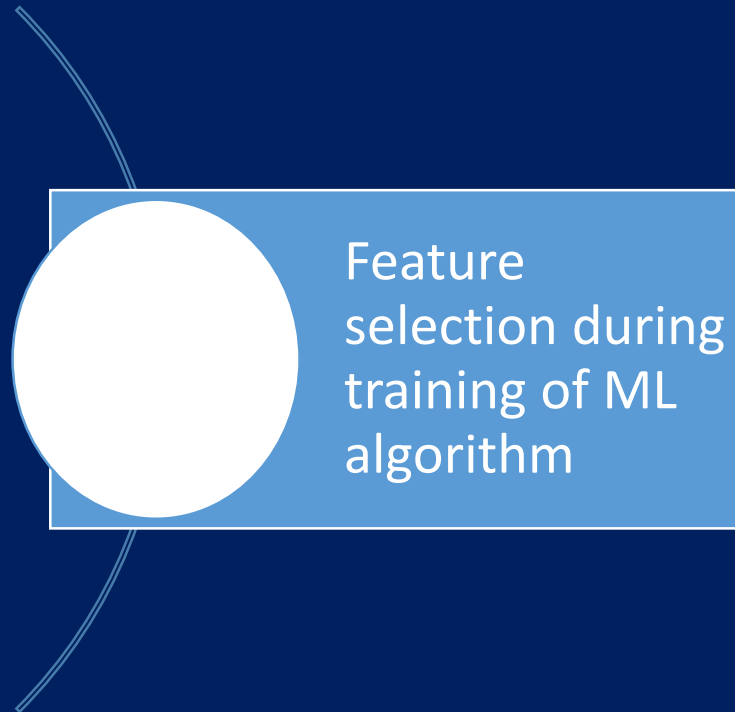
Wrapper methods



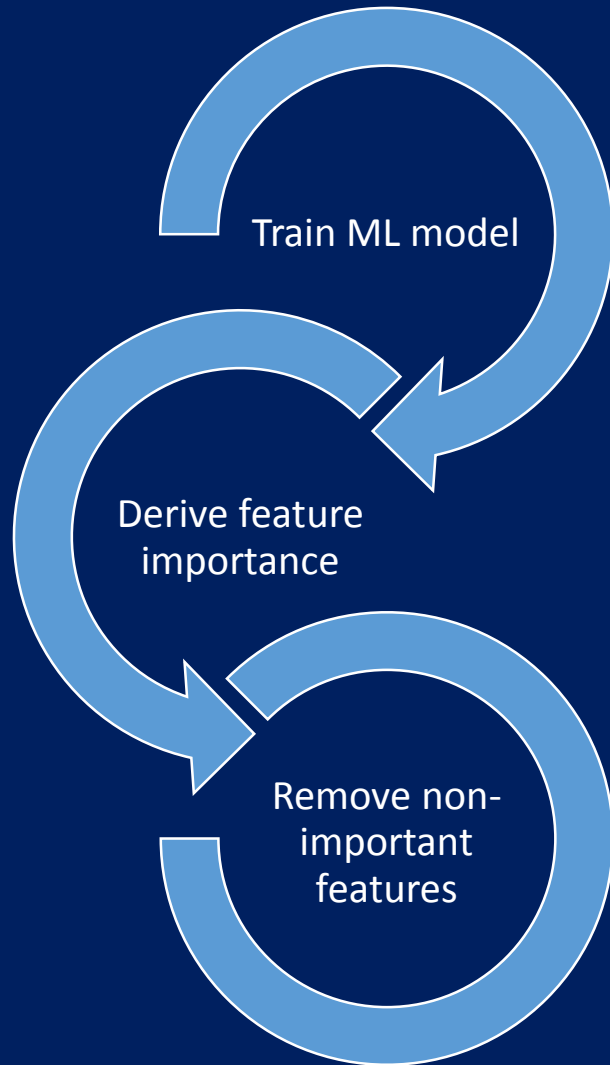
Wrapper methods



Embedded methods



Embedded methods




Data Pre-processing 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**

Knowledge Resources



Feature Selection for Machine Learning
Soledad Galli
★★★★★ 4.6 (156)



Feature Engineering for Machine Learning
Soledad Galli
★★★★★ 4.5 (253)

Udemy.com, includes code

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**

**Feature
Engineering +
Selection**



Feature Engine

Python package for feature engineering
Work in progress

