

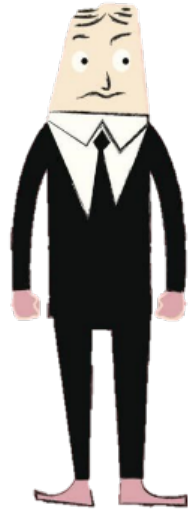
# Machine Learning

An Introduction

08/08/2019

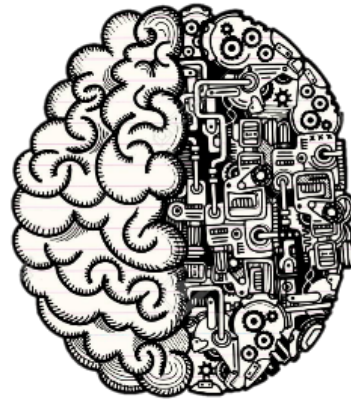
# Machine Learning (ML)

Learn from experience

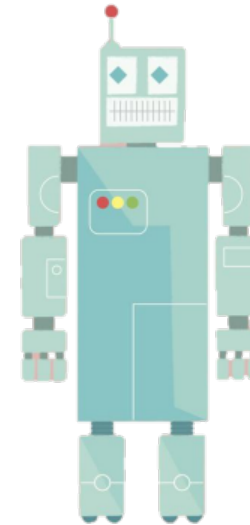


DATA

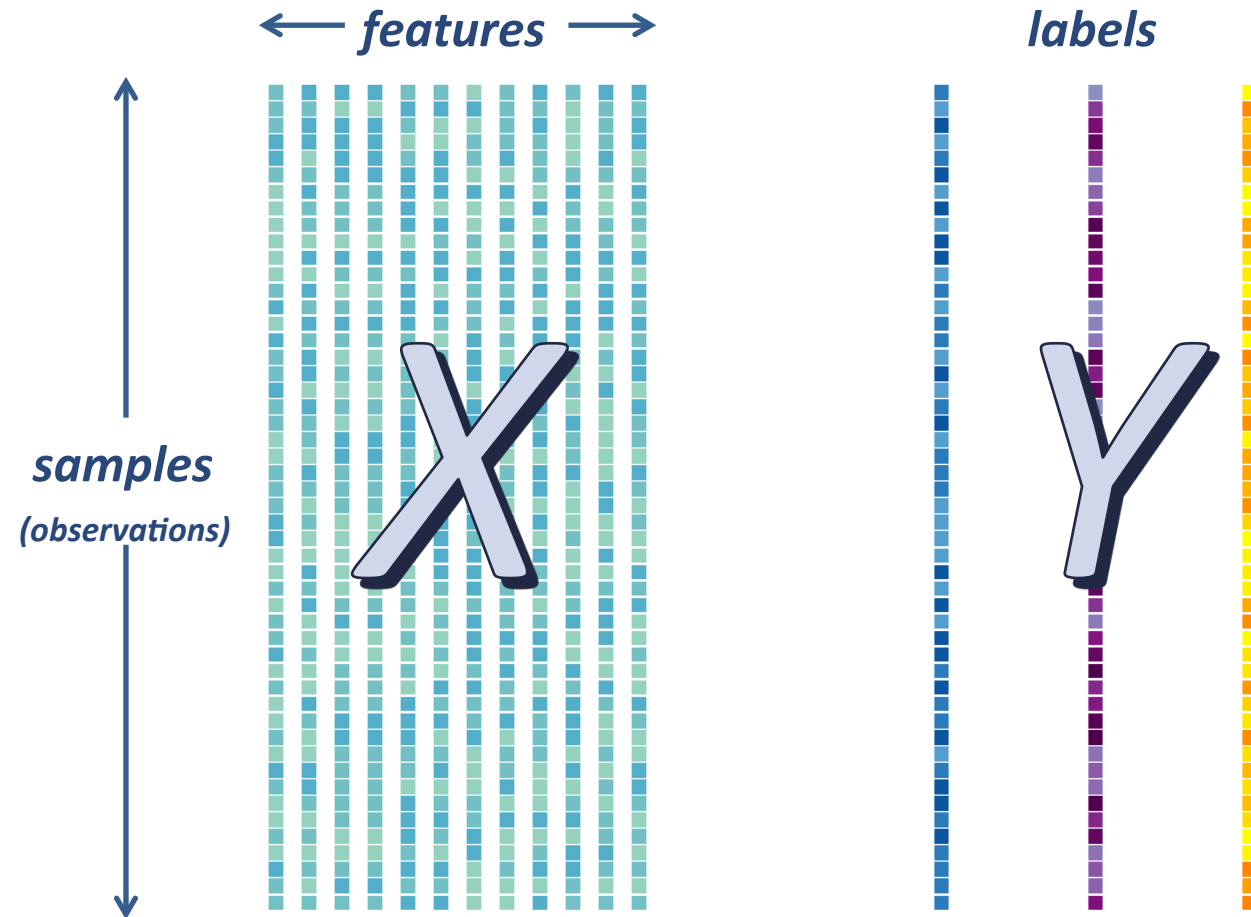
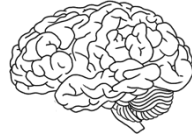
Learn from ~~experience~~



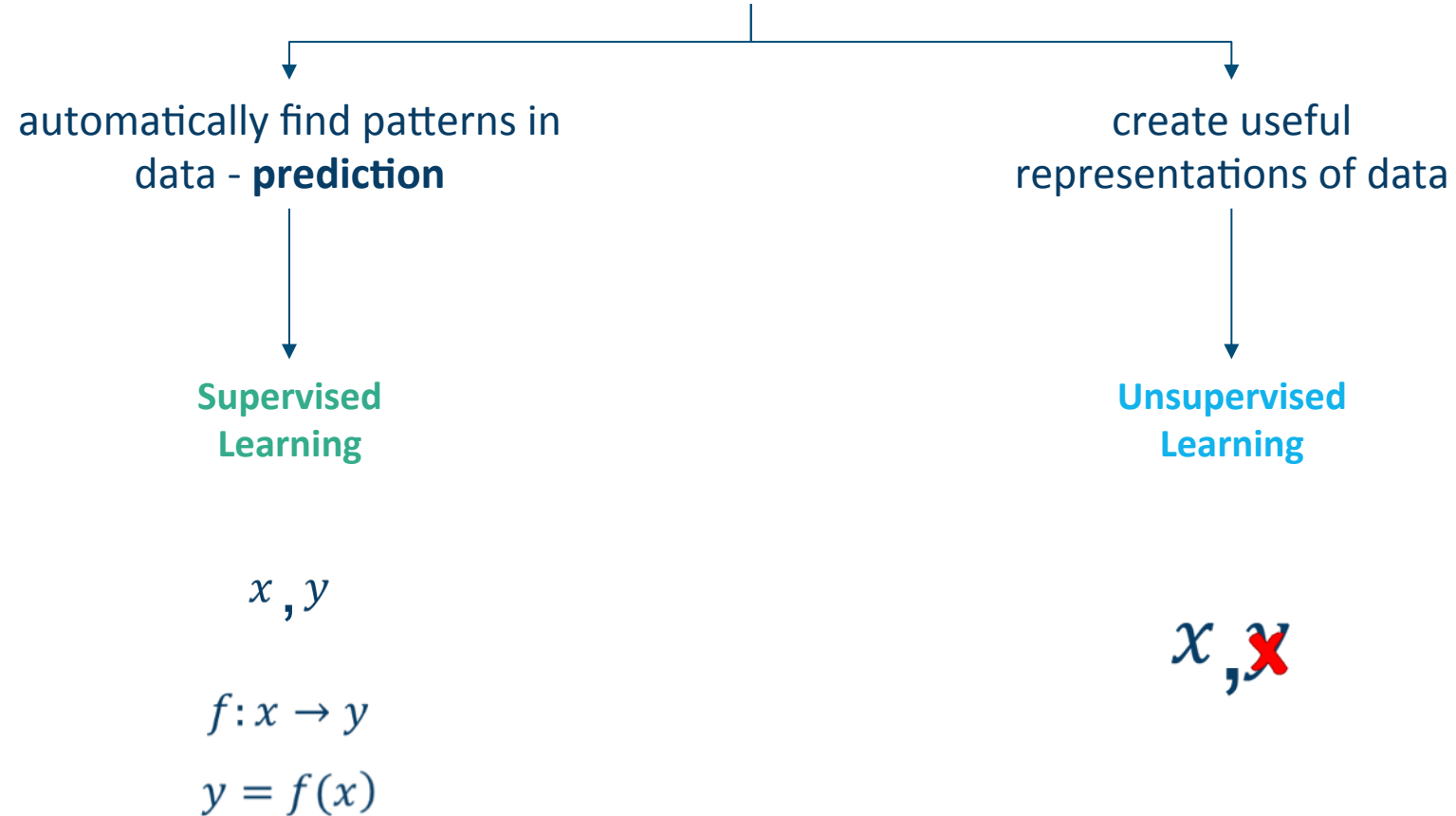
Follow instructions



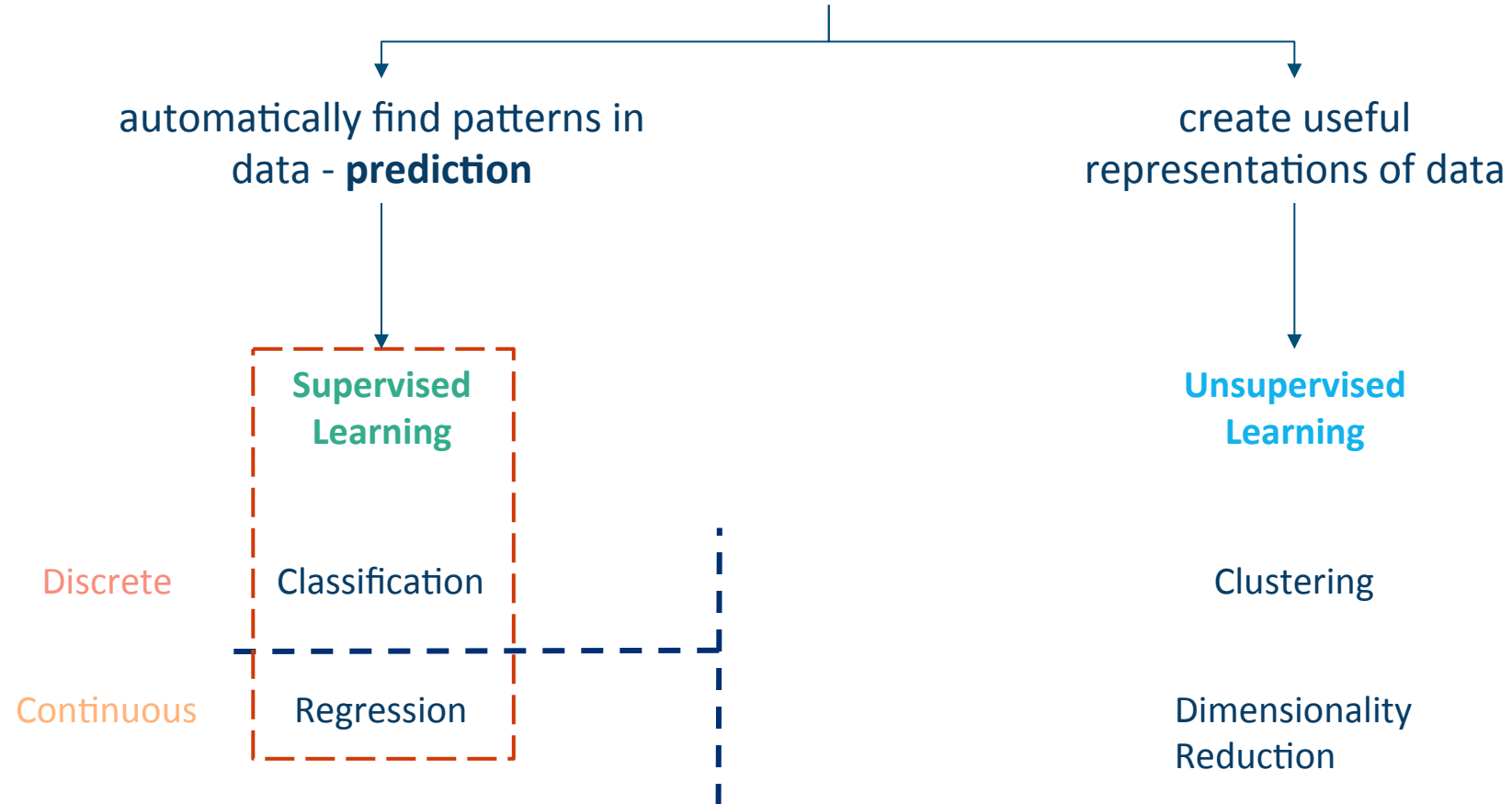
# Terminology



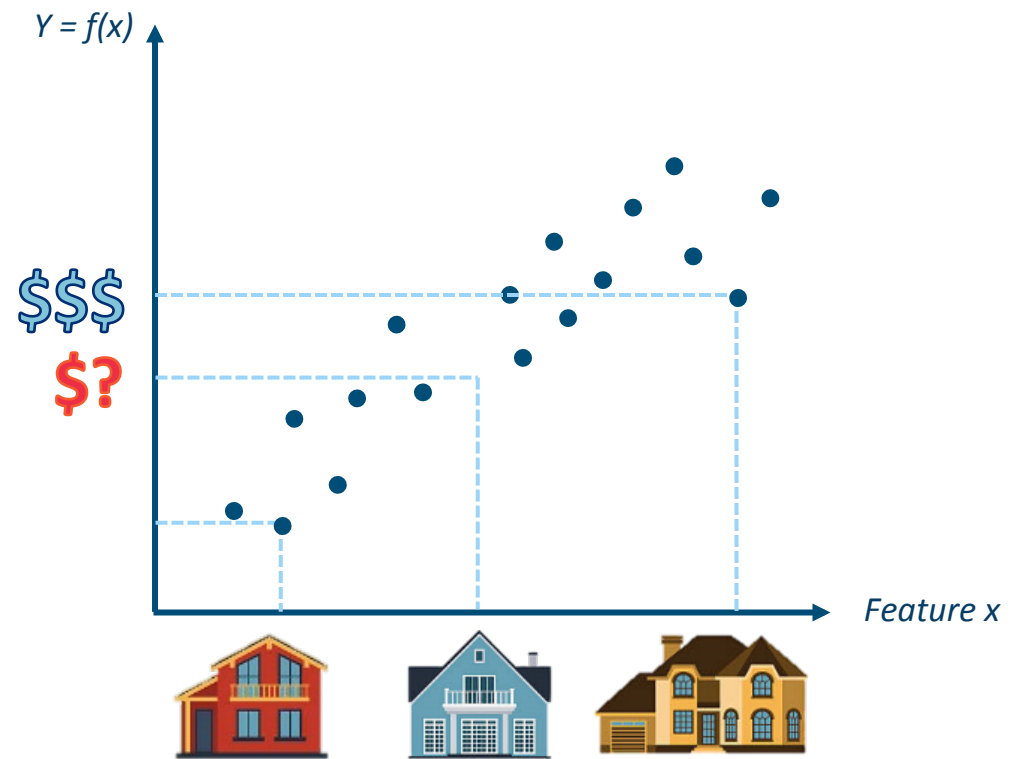
# ML algorithms



# ML algorithms



# Regression

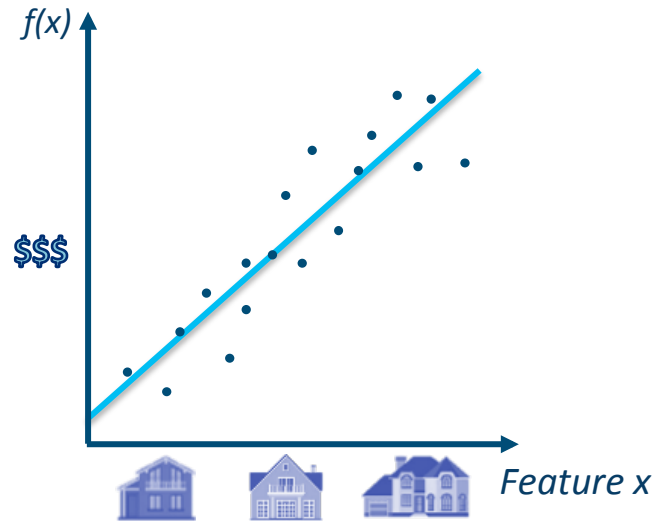


# Regression

$$f: x \rightarrow y$$

$$y = f(x)$$

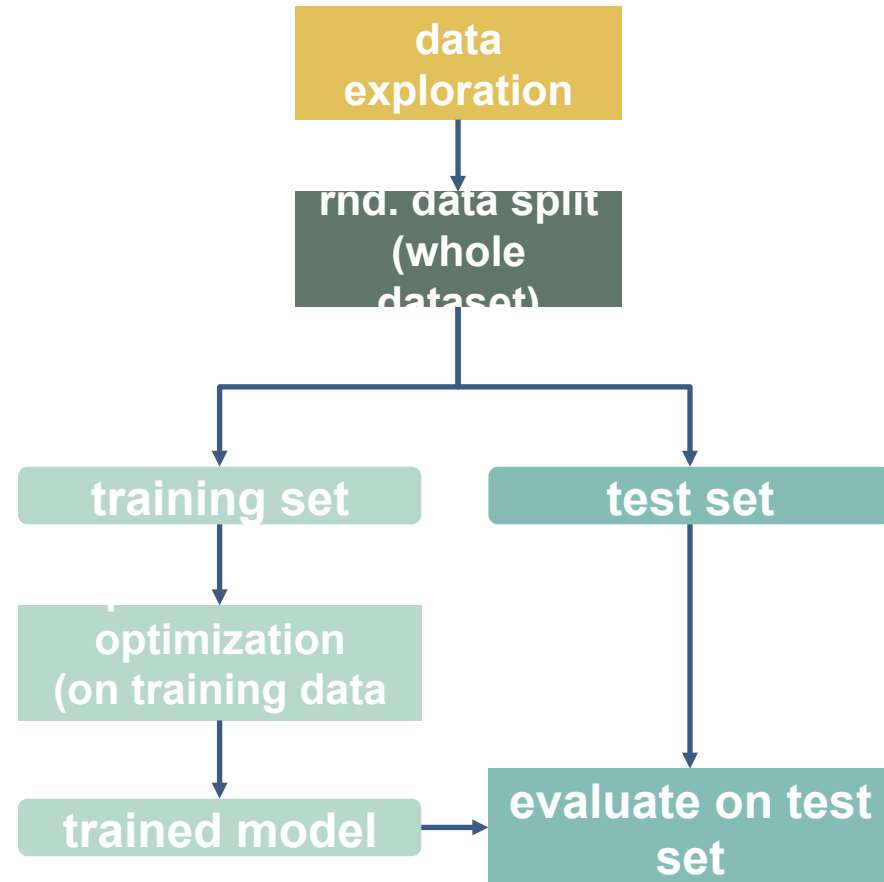
$$y = ax + b$$



$$y = ax^2 + bx + c$$

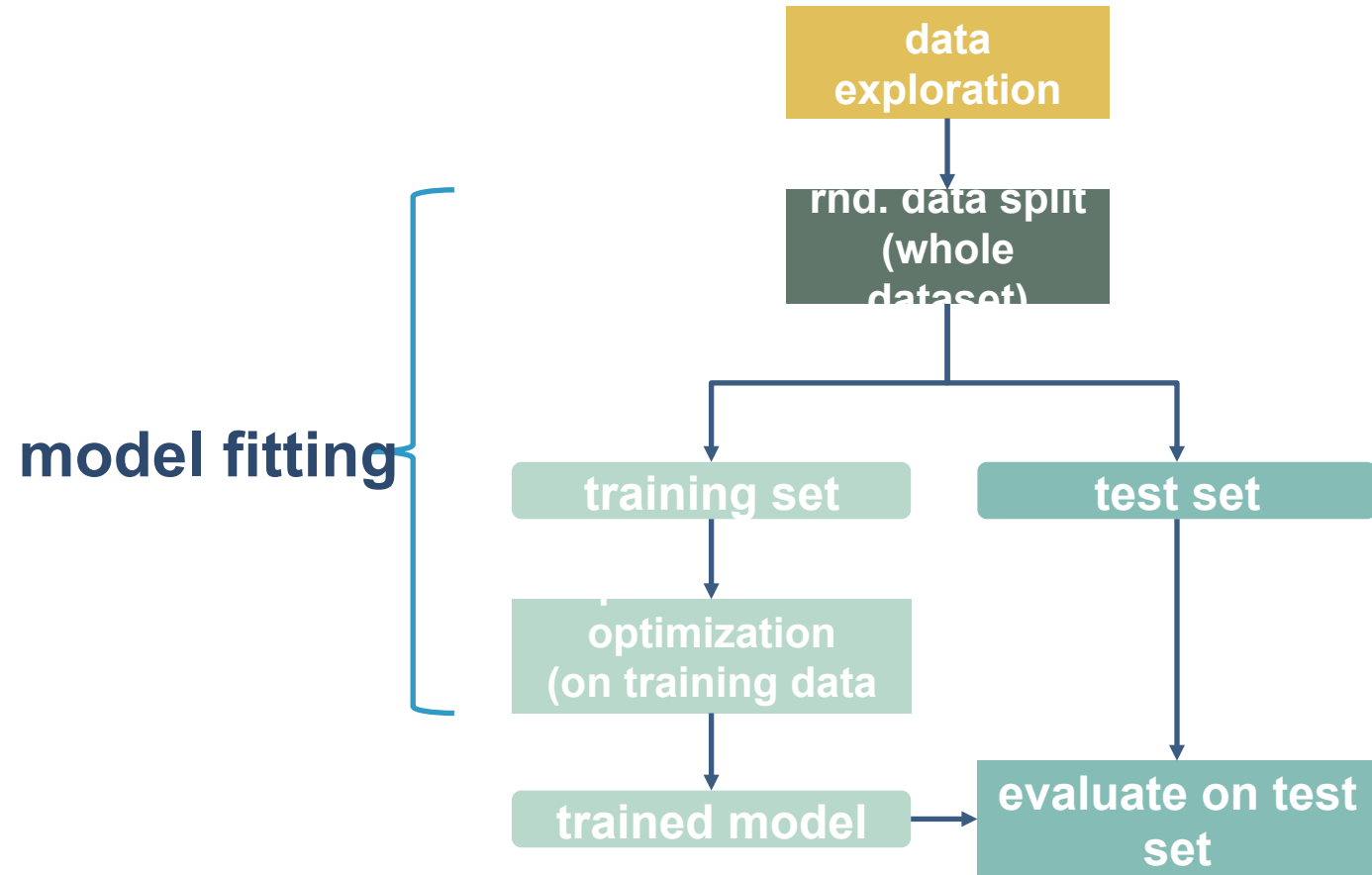


# Typical ML Pipeline





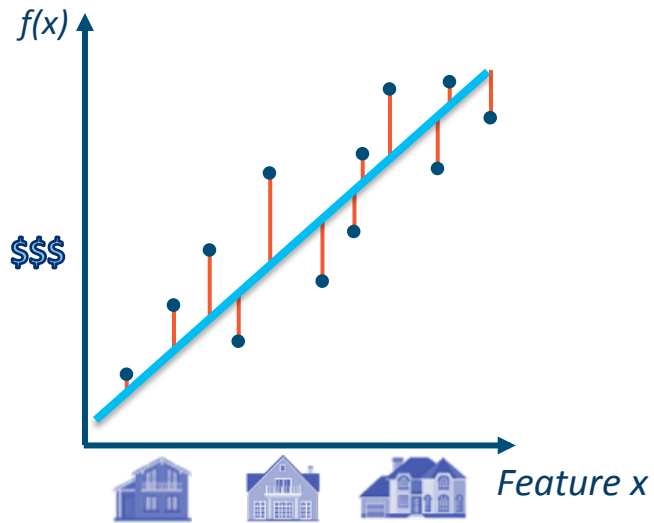
# Typical ML Pipeline



# Model Fitting

$$f: x \rightarrow y$$

$$y = f(x)$$



0



$x \downarrow train, y \downarrow train$

$x \downarrow test, y \downarrow test$

1

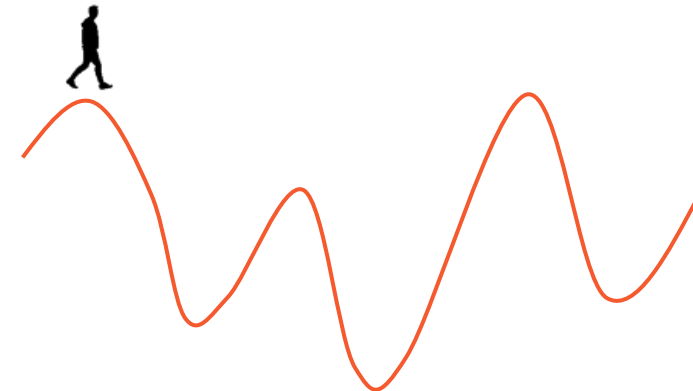
$$f: x \rightarrow y \longrightarrow y = ax + b$$

1

$$Loss = g(a, b) \longrightarrow MSE = \frac{1}{n} \sum_{i=1}^n (y \downarrow i - \hat{y}_i)^2$$

3

optimization algorithm  $\longrightarrow$  gradient descent

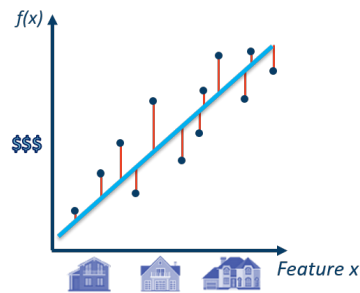


# Model Fitting

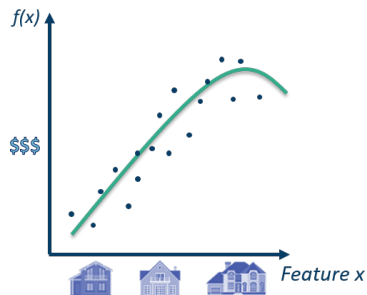
## ② Model (*estimator*) selection

linear regression

$$y = ax + b$$

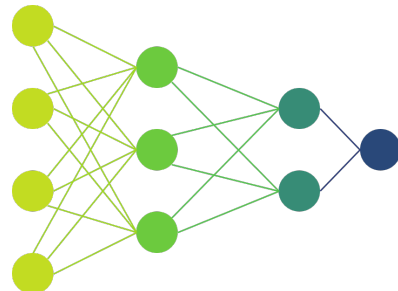


$$y = ax^2 + bx + c$$



polynomial regression

multi-layer  
perceptron  
(ANN)



## ③ loss (*cost*) function

$$\text{mean squared error} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

$$\text{mean absolute error} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

$$\text{max error} = \max(|y_i - \hat{y}_i|)$$

$$\text{explained var.} = 1 - \frac{\text{Var}(y - \hat{y})}{\text{Var}(y)}$$

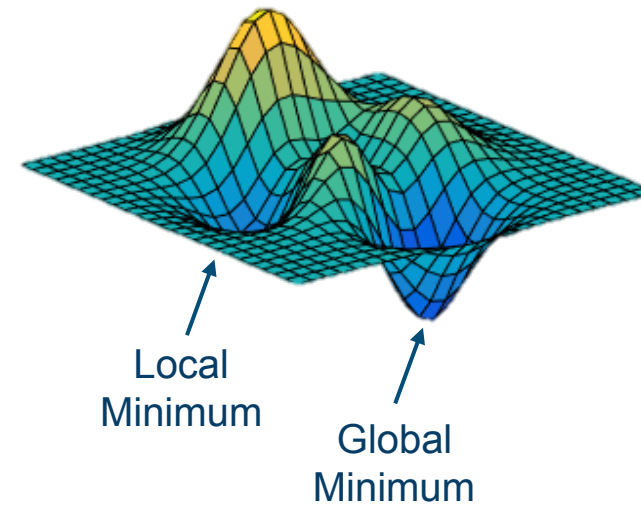
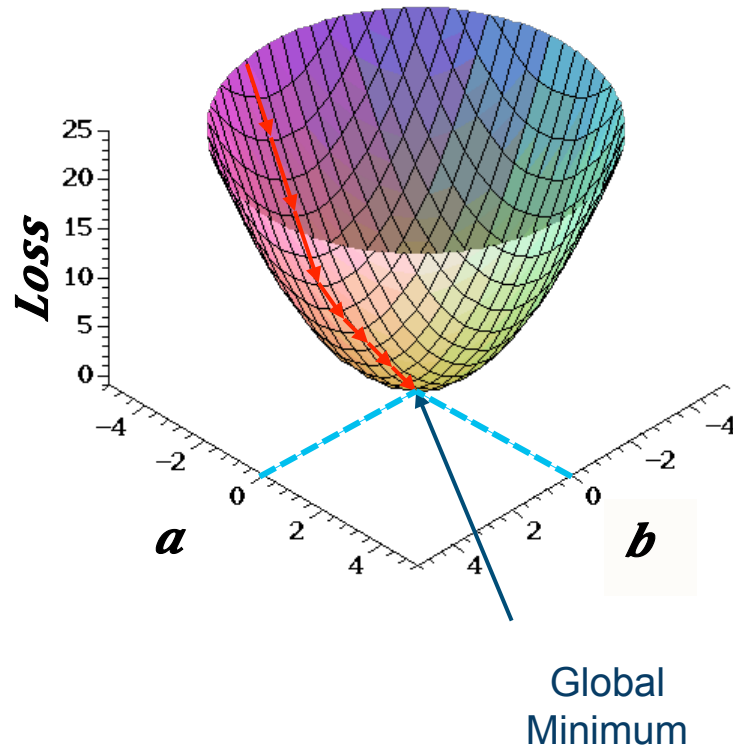
$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

# Gradient Descent

## 4 Optimization of model parameters!

$$y = ax + b$$

$$\text{Loss} = f(a, b)$$



# Gradient Descent

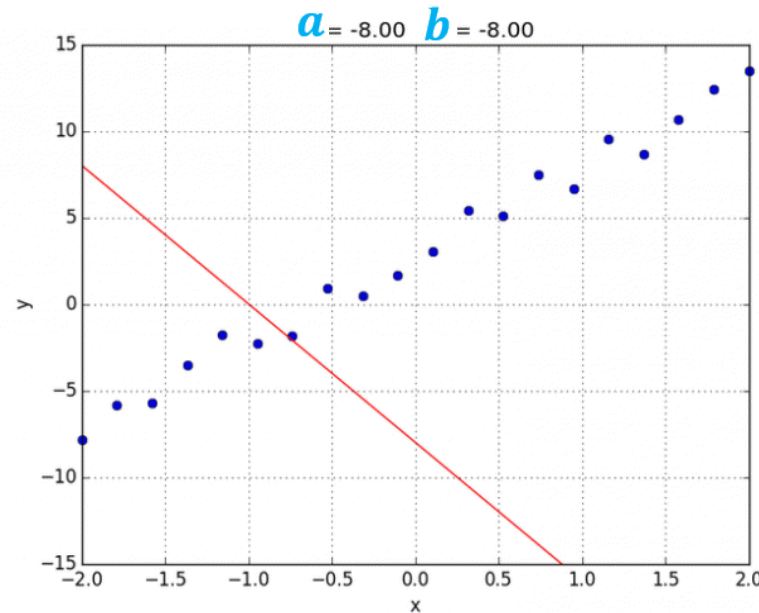
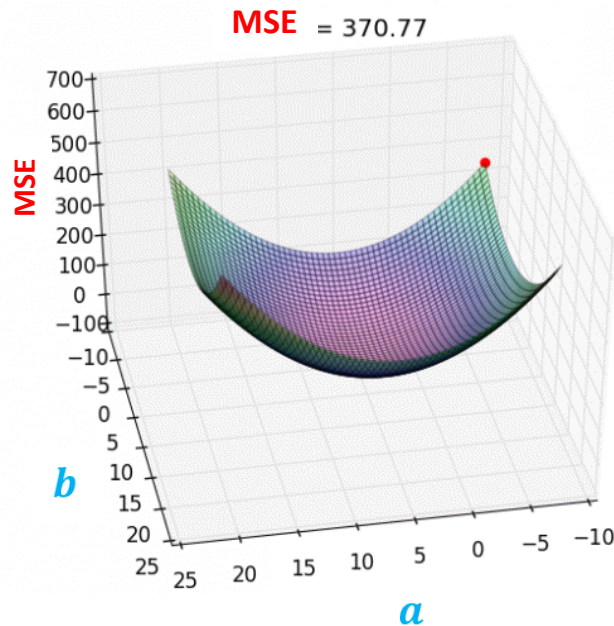
Initialize model parameters ( $a$ ,  $b$ ) randomly

Iterate between:

- 1) Compute *loss* (mean square error – MSE)
- 2) Update model parameters ( $a$ ,  $b$ ) in direction of gradient

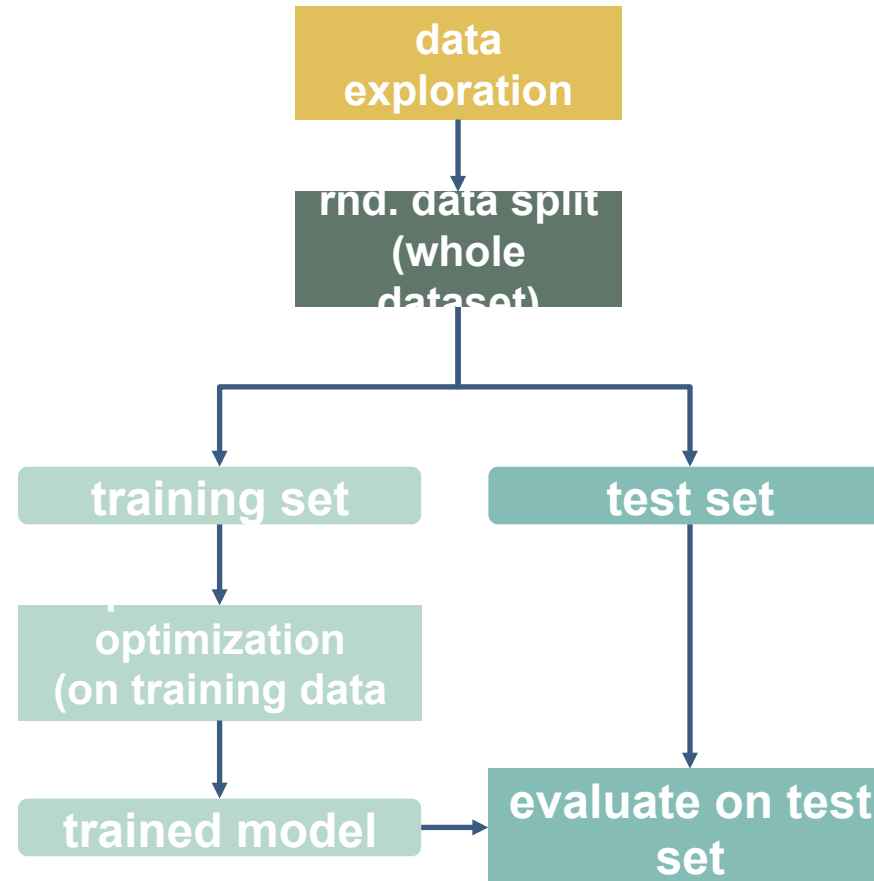
Training data

$$y = ax + b$$

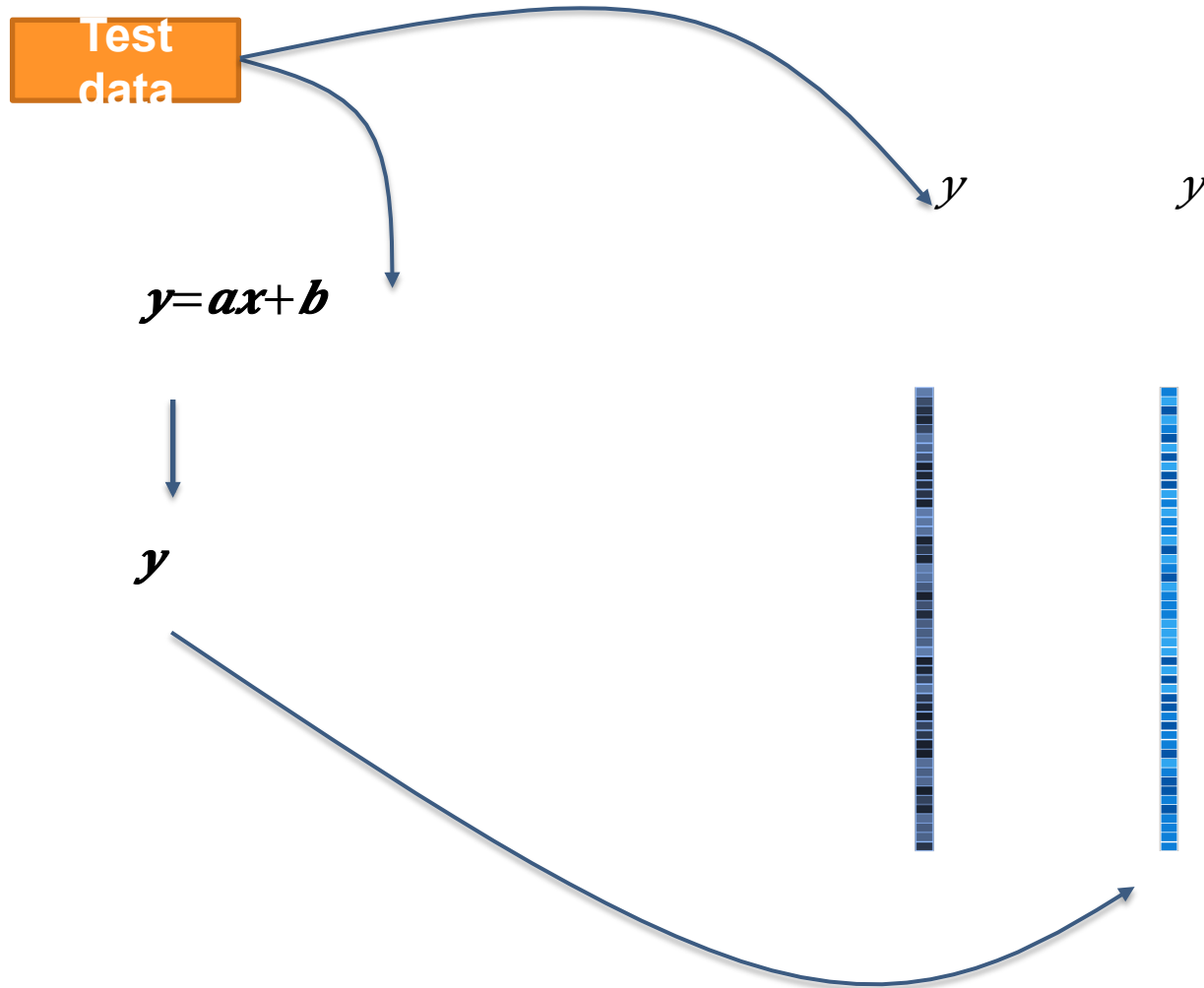


# Typical ML Pipeline

model validation



# Model Validation



## score function

$$\text{mean squared error} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

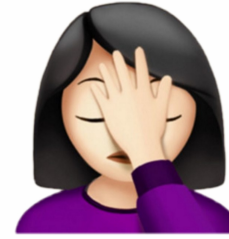
$$\text{mean absolute error} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

$$\text{max error} = \max(|y_i - \hat{y}_i|)$$

$$\text{explained var.} = 1 - \frac{\text{Var}(y - \hat{y})}{\text{Var}(y)}$$

$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

# Houston ... we have a problem!

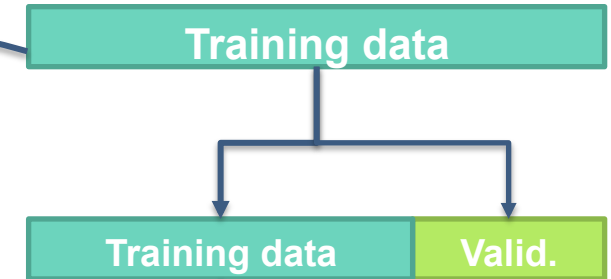
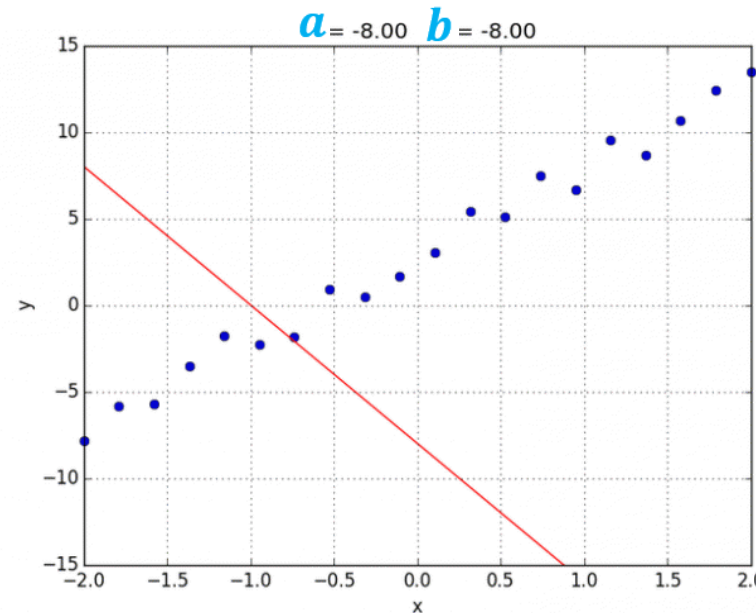
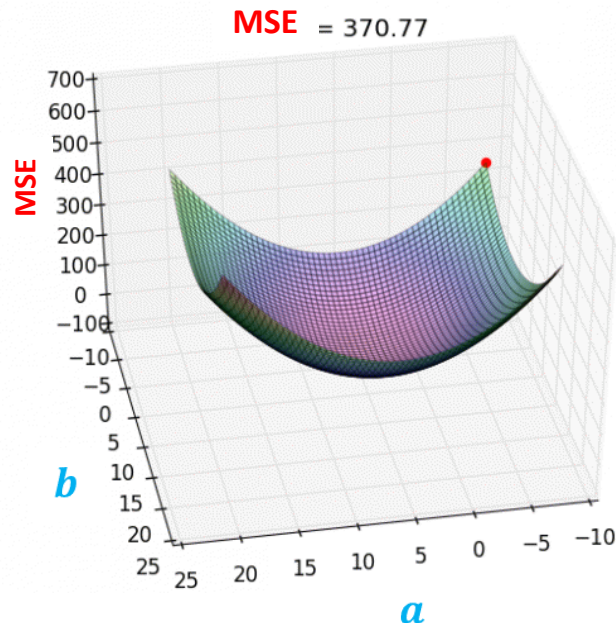


Initialize model parameters ( $a, b$ ) randomly

Iterate between:

- 1) Compute *loss* (mean square error – MSE)
- 2) Update model parameters ( $a, b$ ) in direction of gradient

$$y = ax + b$$



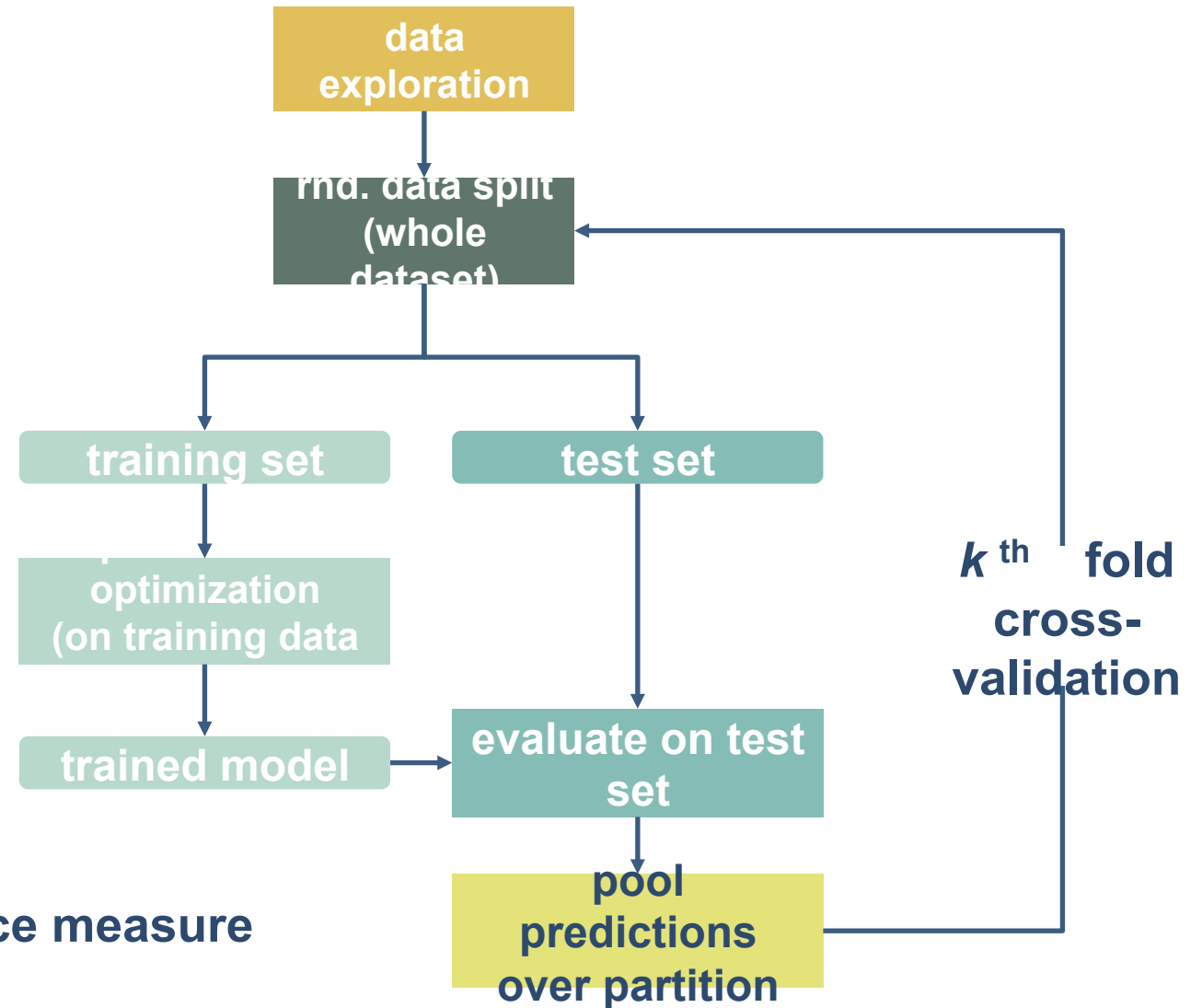
multiple  
training  
epochs!

How many TE?





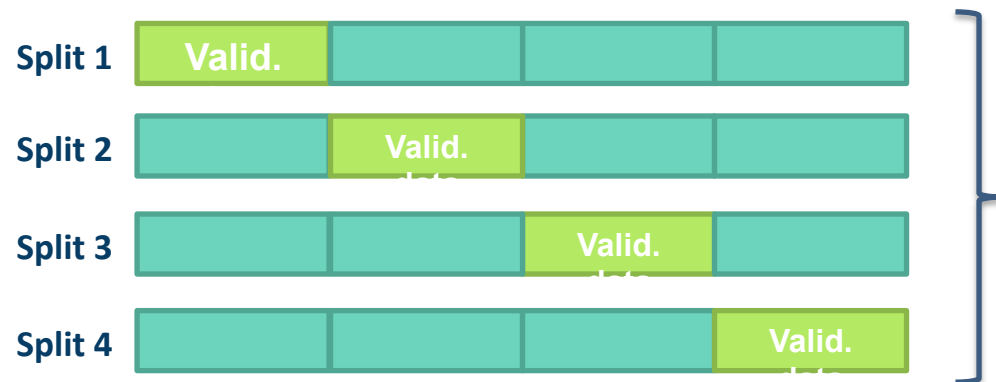
# Typical ML Pipeline



## WHY CV?

- Avoid overfitting
- Improve model generalizability
- Provide a more accurate performance measure
- *Hyper*-parameter tuning

# [*k-fold*] Cross-Validation (CV)



***Hyper-parameter tuning!***

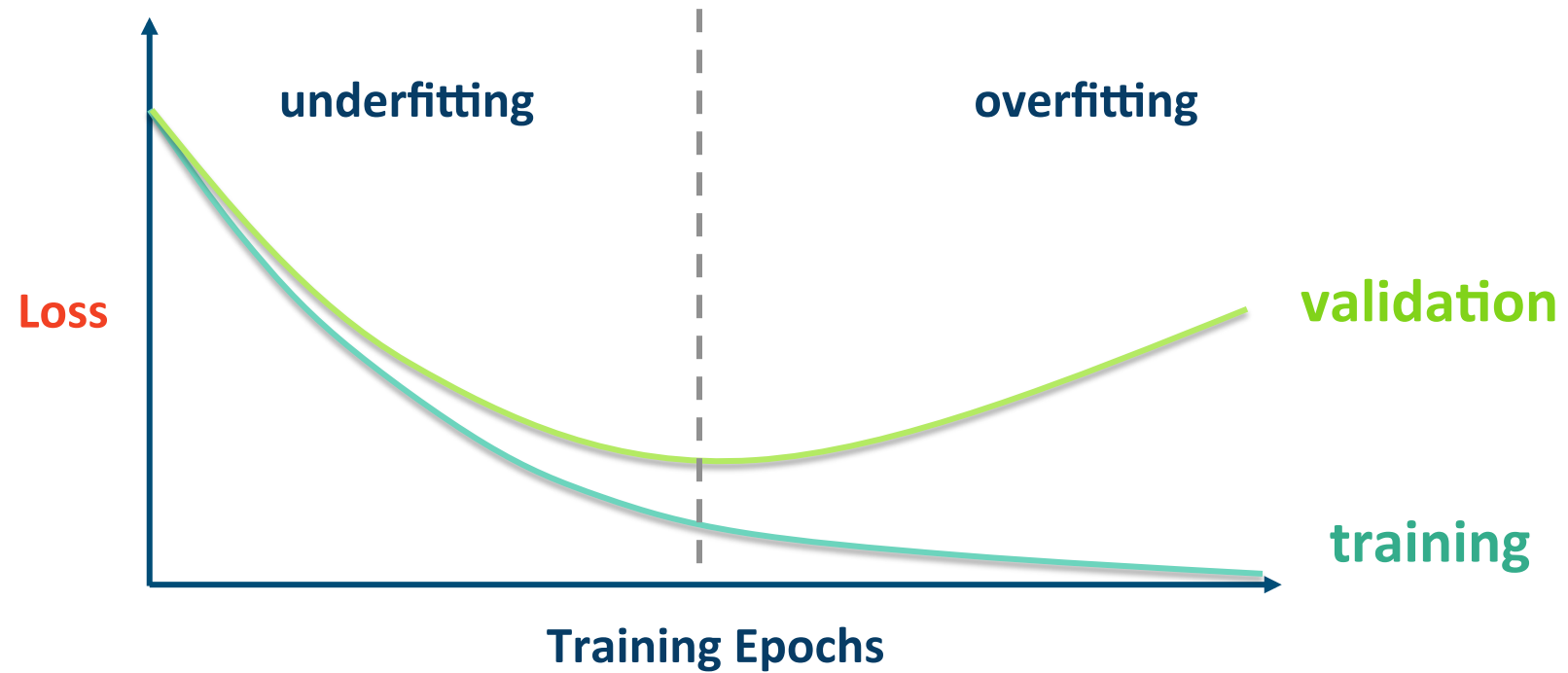
- 1) Train model parameters on **training** set
- 2) Evaluate training with the **validation** set
- 3) Report error on **test** set

# Overfitting

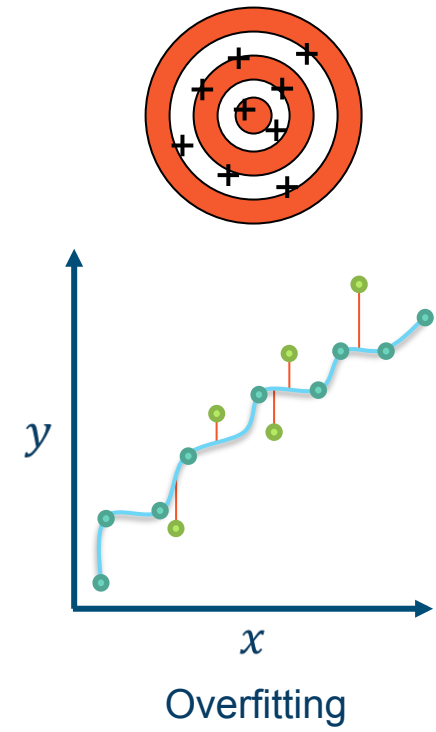
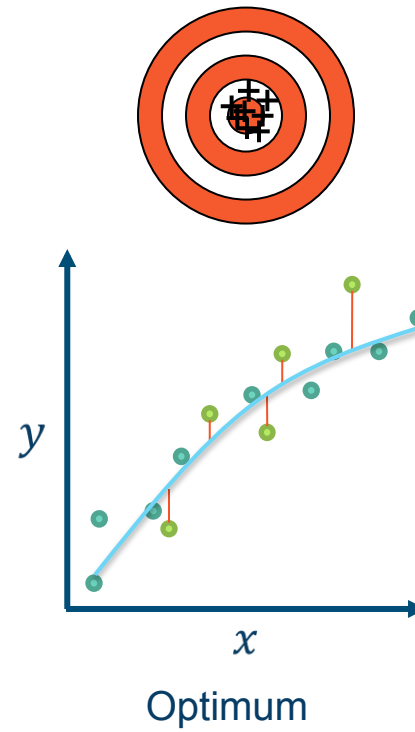
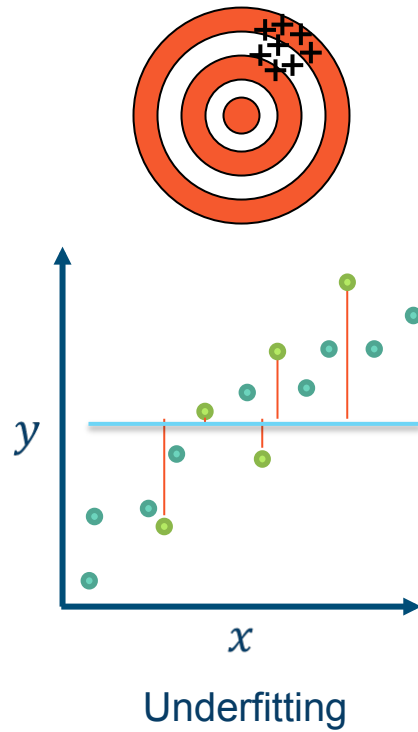
$k^{th}$ -fold



How many TE?



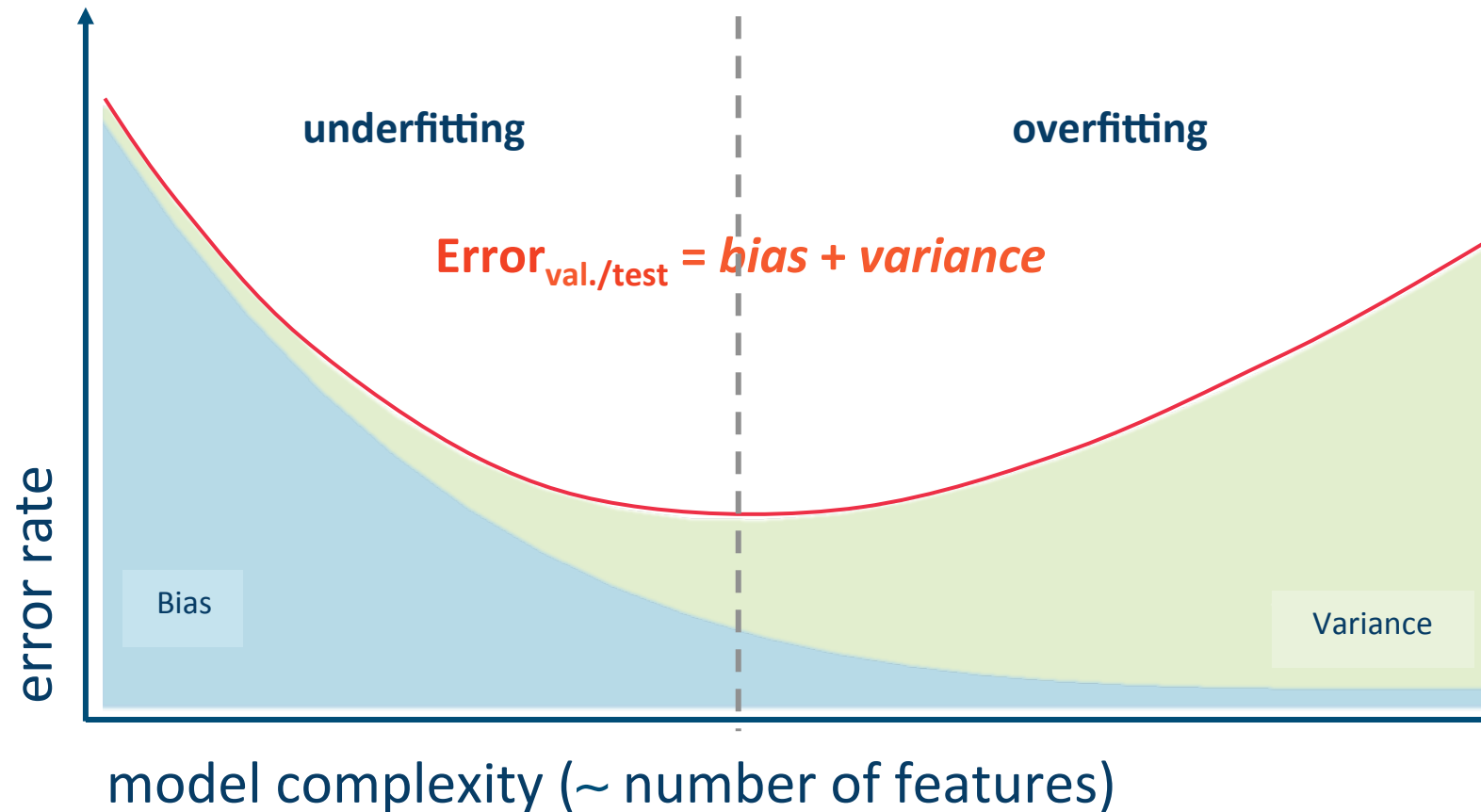
# Bias-Variance Trade-Off



● Training set

● Test set

# Bias-Variance Trade-Off

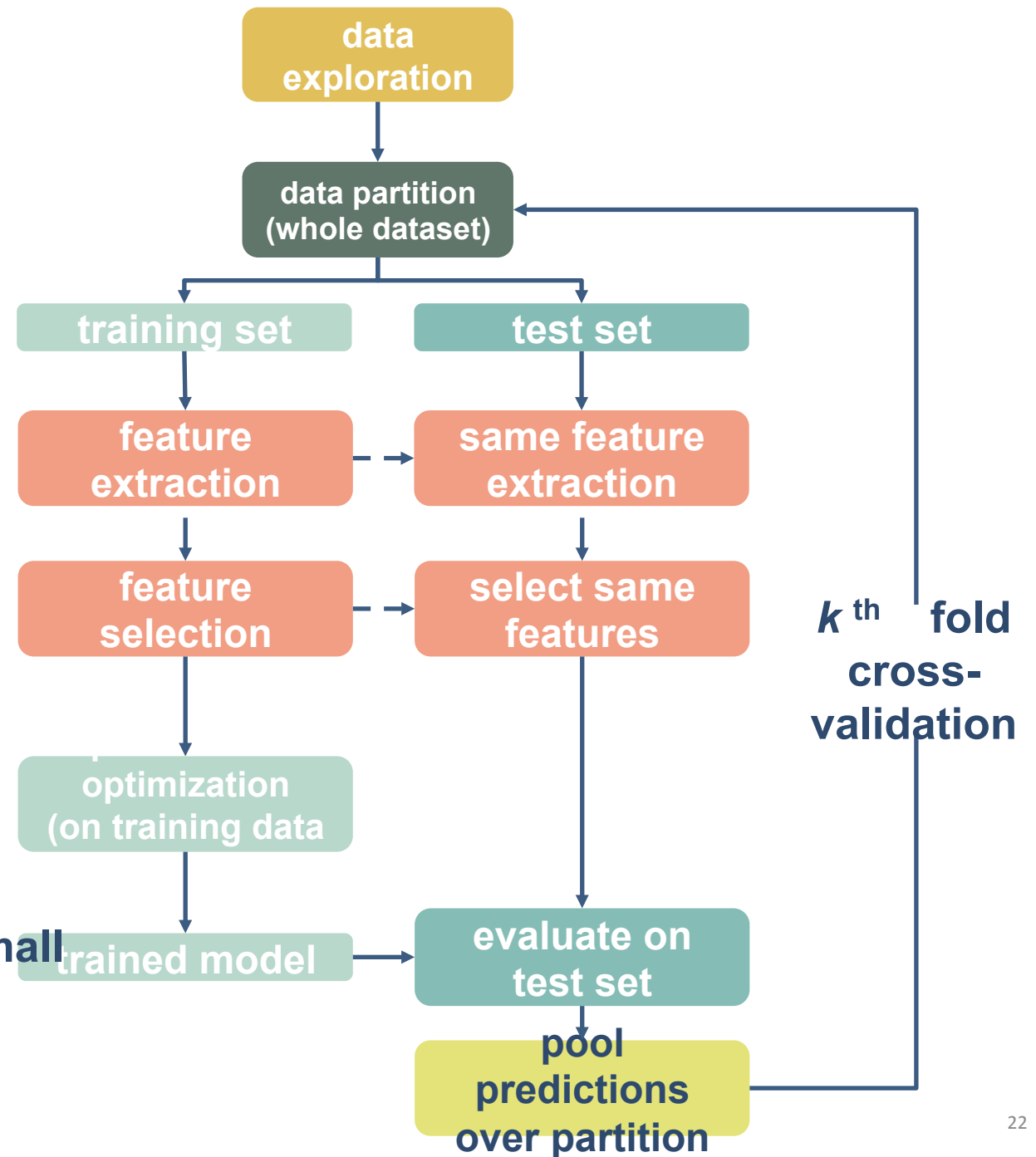


# Typical ML Pipeline

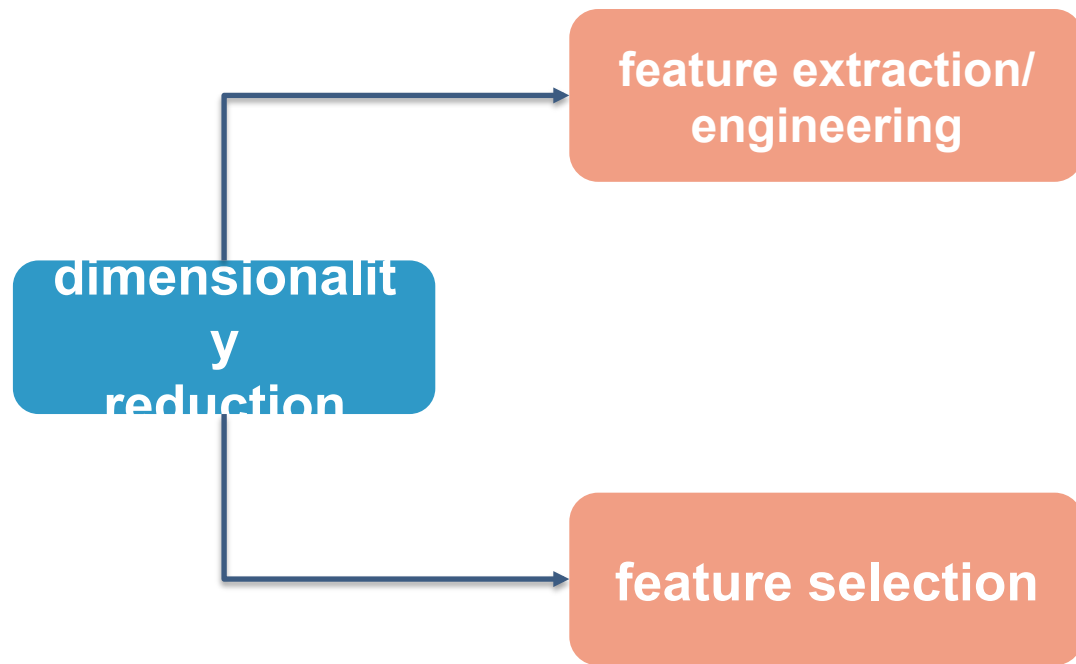
dimensionality  
reduction

## WHY DR ?

- Curse of dimensionality (more features than samples)
- Intrinsic dimension may actually be small
- Extract “salient” features
- Remove noisy and redundant features
- VISUALIZATION!!!!



# Tip 1: Dimensionality Reduction



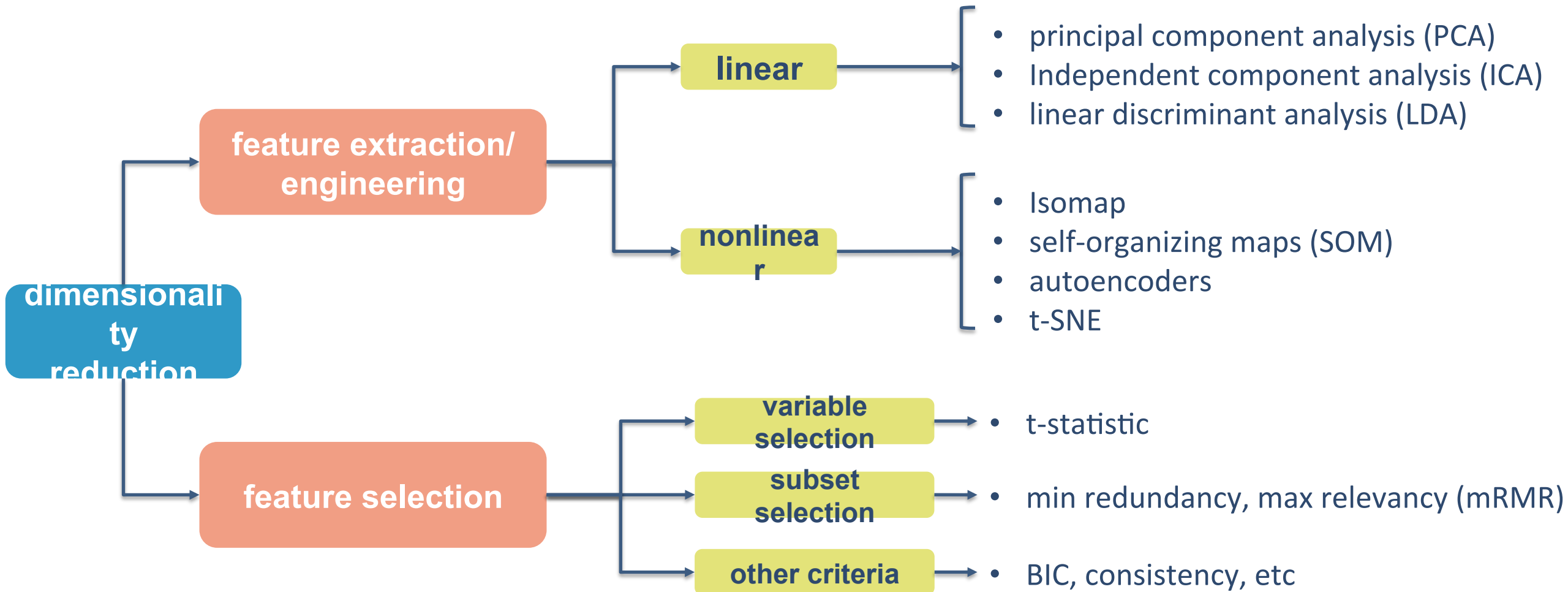
- Compact representation of the data
- Maps input features into a lower dimensional space

e.g. If features are  $X = [x_1, x_2, x_3, x_4]$   
then  $Z = T(X) = [c_1 x_1 + c_2 x_2, x_3 * x_4]$

- Selection of a subset of input features
- Features are still in original space

e.g.  $Z = S(X) = [x_2, x_3]$

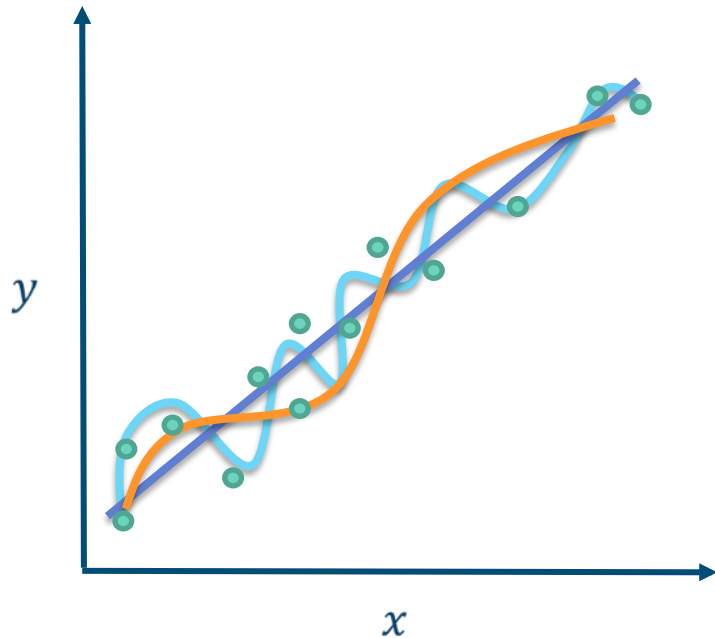
# Tip 1: Dimensionality Reduction





## Tip 2: Regularization

$$y = \cancel{\beta_0} + \beta_1 x + \cancel{\beta_2} x^2 + \dots + \cancel{\beta_p} x^p$$



Penalties on the **LOSS** function to prevent overfitting!

- 1) L1/Lasso: constrains parameters to be **sparse**

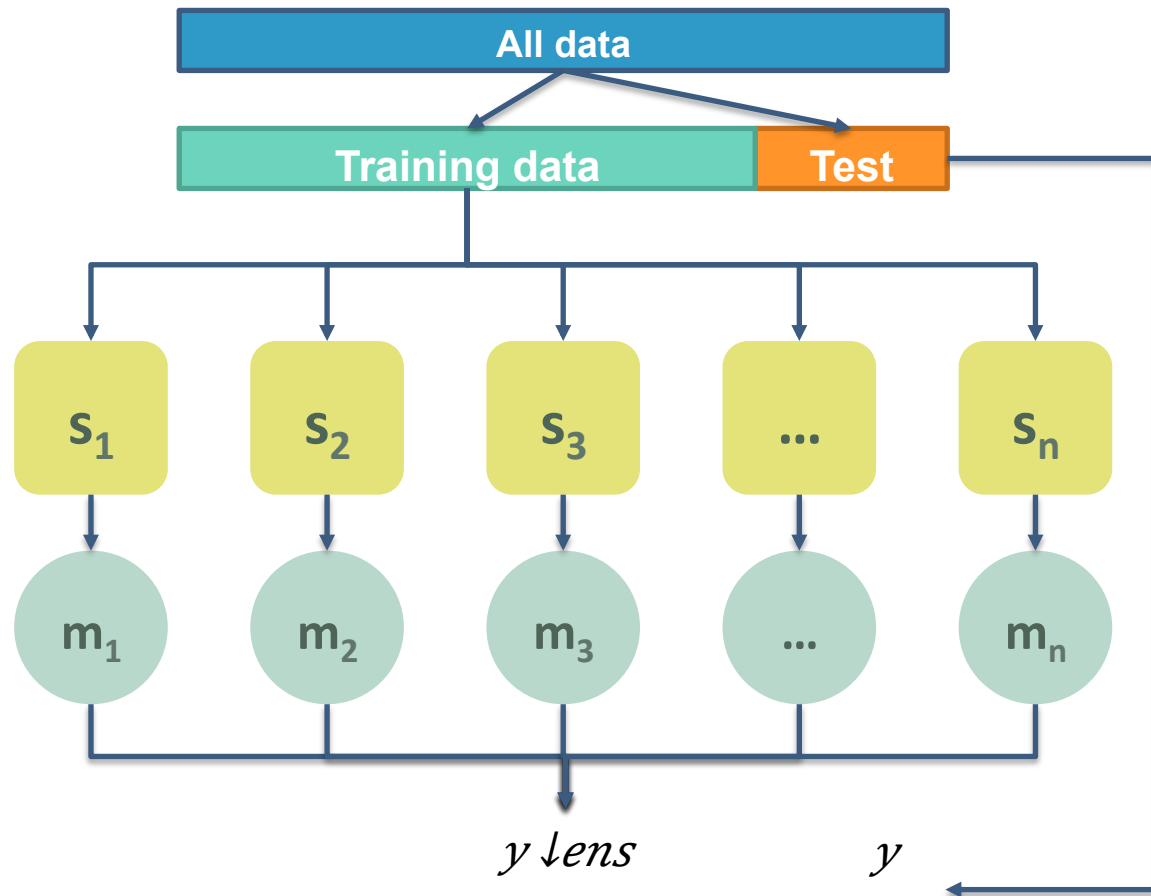
$$MSE = \sum_{i=1}^n \left( y_i - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p |\beta_j|$$

- 2) L2/Ridge: constrains parameters to be **small**

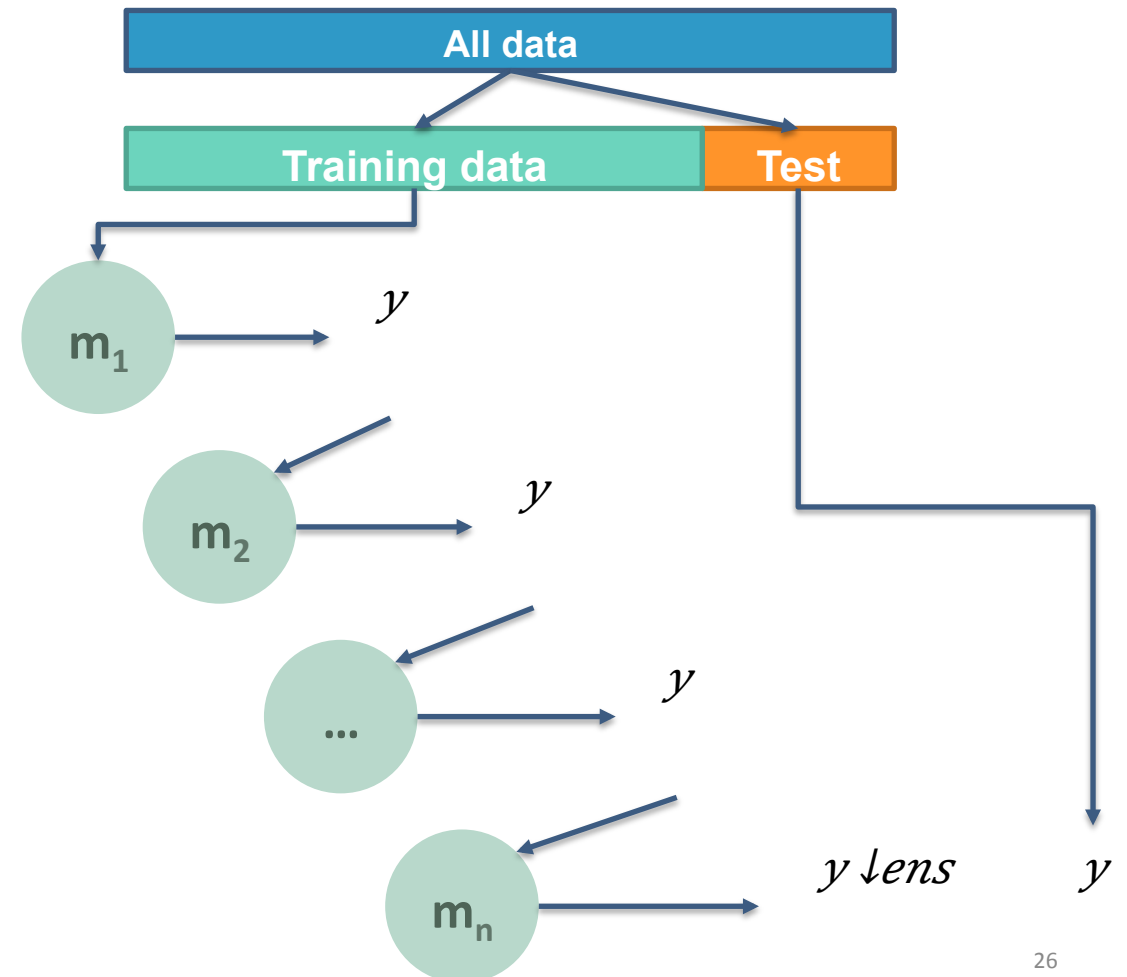
$$MSE = \sum_{i=1}^n \left( y_i - \sum_{j=1}^p x_{ij} \beta_j \right)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

# Tip 3: Ensemble Methods

**bagging** = bootstrapp resampling + aggregation



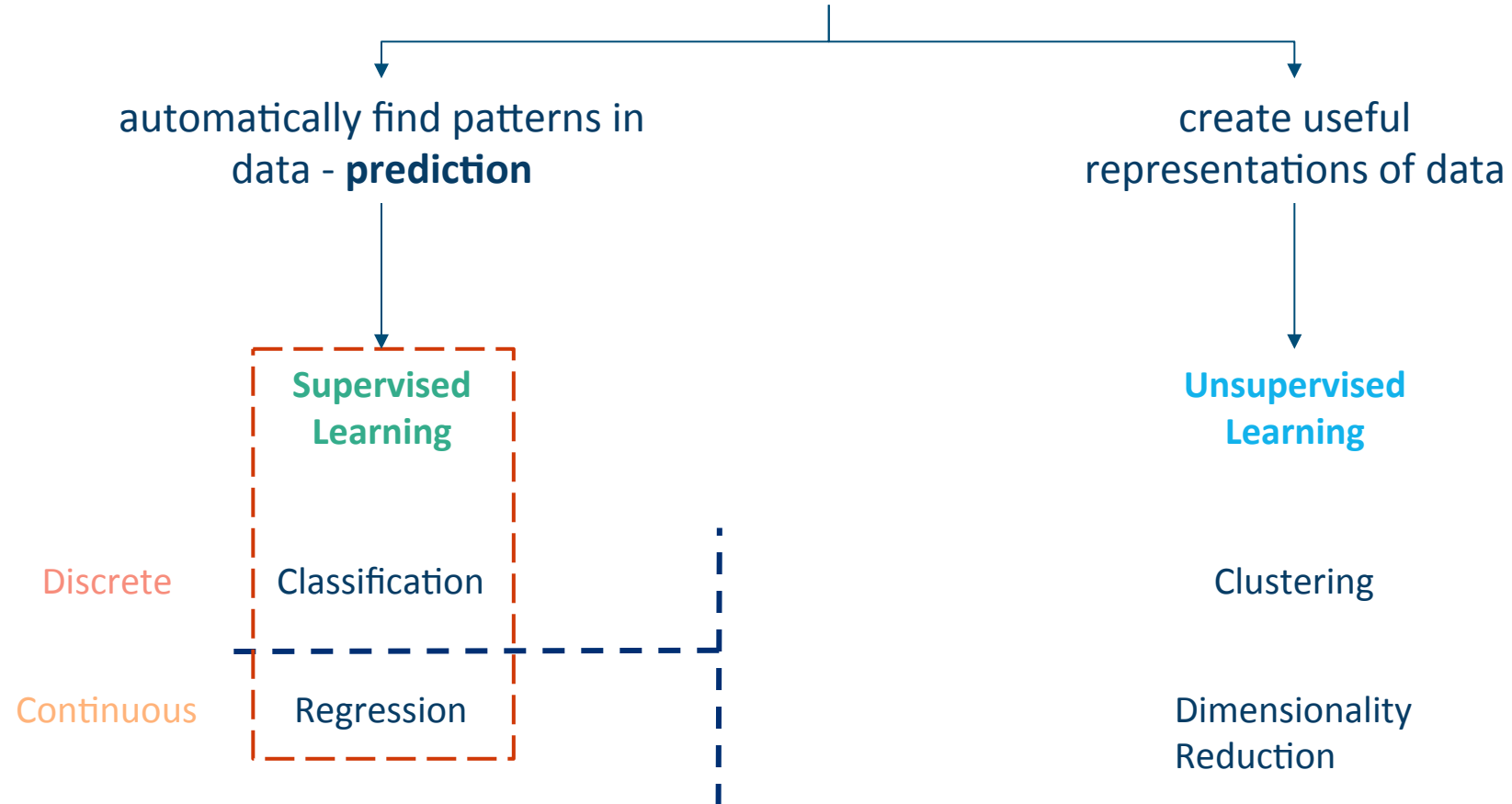
**boosting**



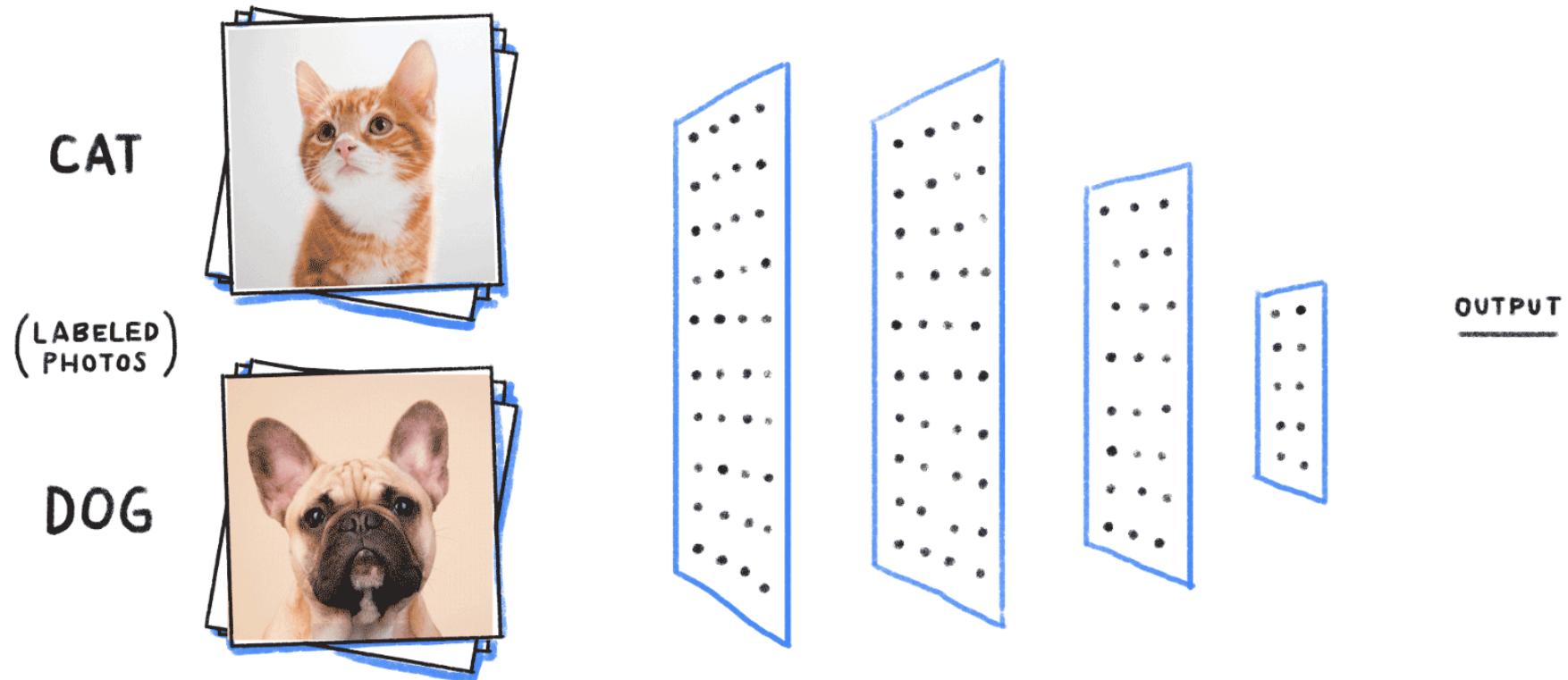
# Sources of Bias

Type	DO NOT ☹️	Sexy name	DO 😊
<i>k</i> -hacking	Try many <i>k</i> 's in <i>k</i> -fold CV (or different training %) and report only the best	<i>k</i> -hacking	Pick <i>k</i> =10, repeat It many times (>200 or as many as possible!), and report the full distribution (NOT boxplots!)
<i>metric</i> -hacking	Try different performance metrics (e.g., accuracy, F1, AUC, error rate, etc.) and report the best	<i>m</i> -hacking	Choose the most appropriate and recognized metric for the problem (e.g., AUC for binary classification)
<i>feature/dataset</i> -hacking	Try subsets of feature(s) or subsamples of dataset(s), but report only the best	<i>d</i> -hacking	Use and report on everything: all analyses on all datasets

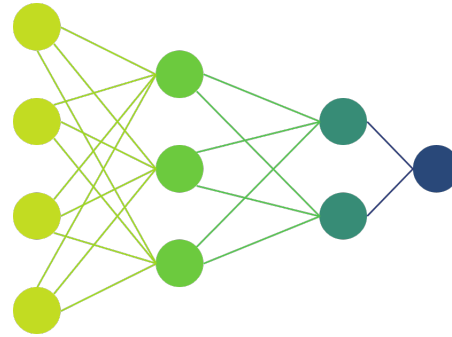
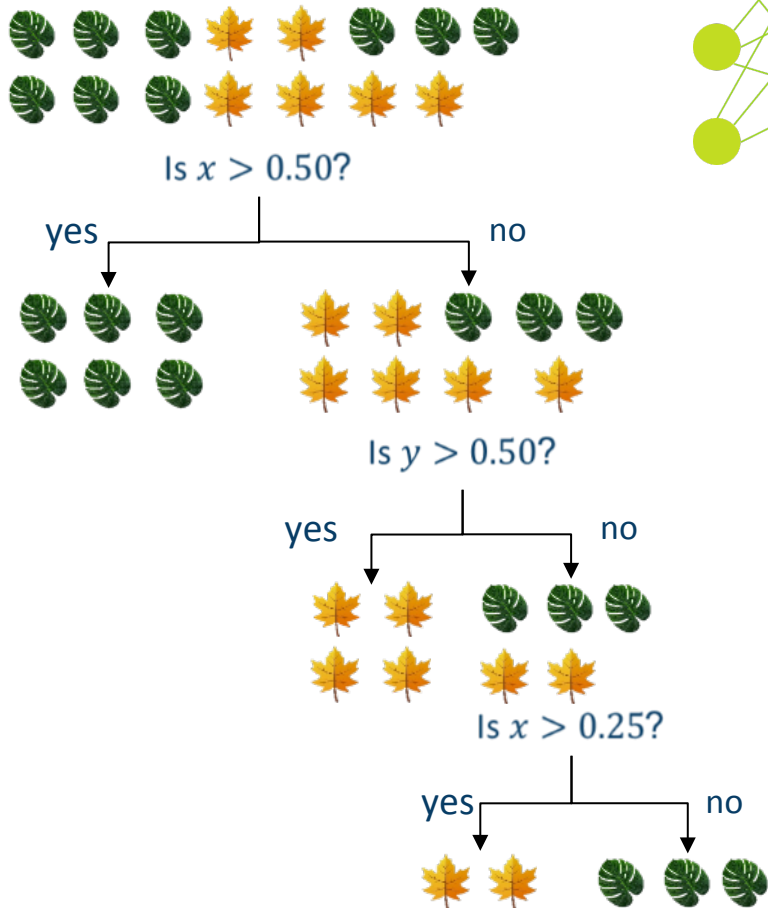
# ML algorithms



# Classification



# Classification

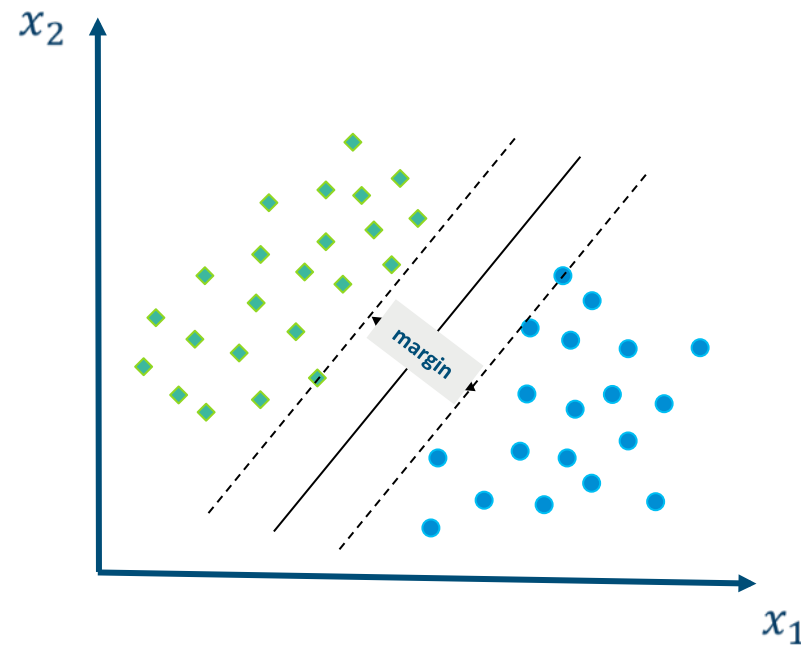


- Support Vector Machine

(SVM)

- Artificial Neural Networks
- Logistic regression
- Decision Trees
- Random Forests

probabilistic  
classifiers



# Performance Metrics – Binary Classification

confusion matrix

		True condition	
		POSITIVE	NEGATIVE
Predicted condition	POSITIVE	TRUE positive - $T\downarrow p$	FALSE positive - $F\downarrow p$ (Type I error)
	NEGATIVE	FALSE negative - $F\downarrow n$ (Type II error)	TRUE negative - $T\downarrow n$

score function

$$accuracy = (T\downarrow p + T\downarrow n) / (T\downarrow p + T\downarrow n + F\downarrow p + F\downarrow n)$$

$$precision = T\downarrow p / (T\downarrow p + F\downarrow p)$$

$$recall = T\downarrow p / (T\downarrow p + F\downarrow n)$$

Appropriate  
when classes  
are imbalanced!

$$F1 \text{ score} = 2T\downarrow p / (2T\downarrow p + F\downarrow p + F\downarrow n)$$

# Performance Metrics – Binary Classification

		True condition	
		POSITIVE	NEGATIVE
Predicted condition	POSITIVE	TRUE positive - $T\downarrow p$	FALSE positive - $F\downarrow p$ (Type I error)
	NEGATIVE	FALSE negative - $F\downarrow n$ (Type II error)	TRUE negative - $T\downarrow n$
Sensitivity (or recall)		$TPR = T\downarrow p / P = T\downarrow p / (T\downarrow p + F\downarrow p)$	
		$SPR = T\downarrow n / N = T\downarrow n / (T\downarrow n + F\downarrow n)$	
		$Specificity (or selectivity)$	

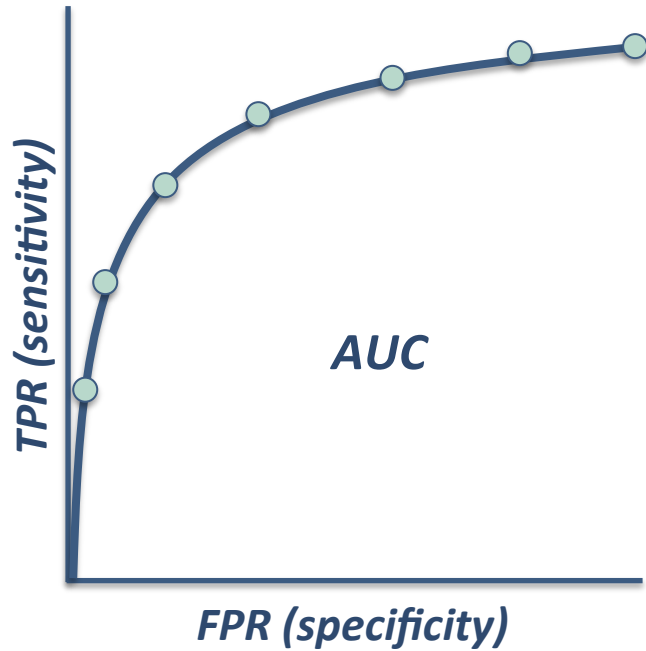
$$FNR = F\downarrow n / P = F\downarrow n / (T\downarrow p + F\downarrow p) \quad FNR = T\downarrow n / N = T\downarrow n / (T\downarrow n + F\downarrow n)$$



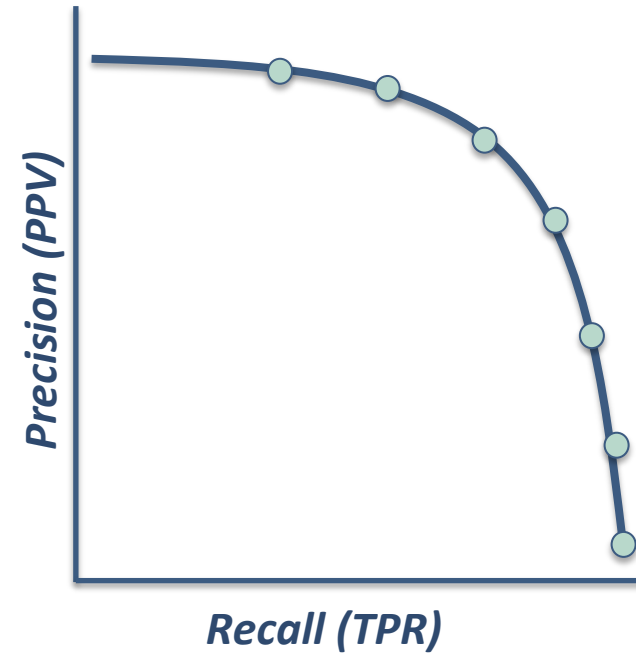
# Performance Metrics – Binary Classification

*(probabilistic classifiers)*

*ROC curve*



*precision-recall curve*

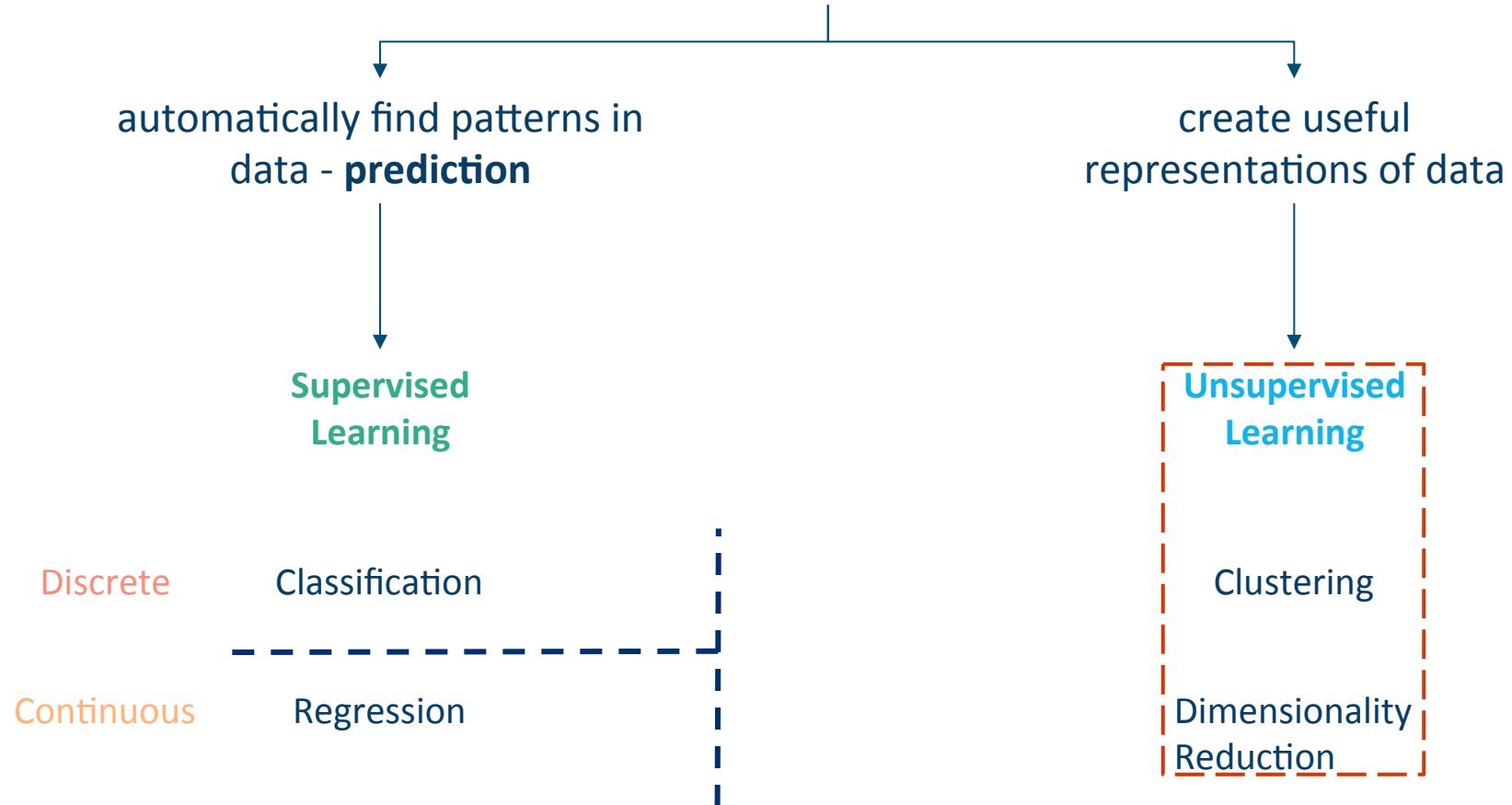


# Multiclass Prediction ( $\neq$ Multilabel Prediction)

		True class		
		CAT	DOG	BIRD
Predicted class	CAT	13	0	0
	DOG	0	10	6
	BIRD	0	0	9

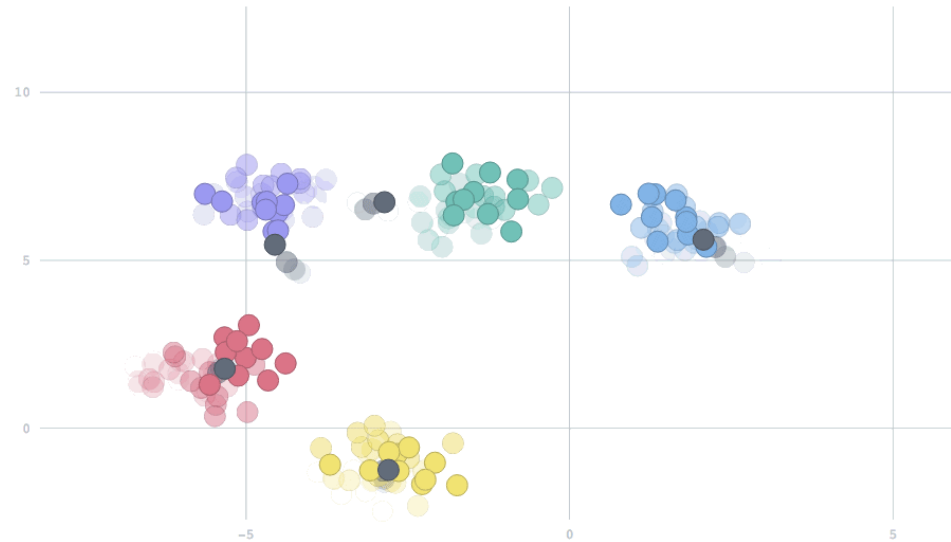
- To extend a binary metric to multiclass problems, the data is treated as a collection of binary problems, one for each class.
- The binary metric is then averaged across the set of classes, each of which may be useful in some scenario.

# ML algorithms



# Clustering

- *K*-means
- Hierarchical clustering
- Mixture of Gaussians



# *Inference vs Prediction* – Linear Model $Y = \beta X + \epsilon$

## statistical inference

### Goal:

- Identify significant contributing variables (statistical null-hypothesis testing,  $p$ -values)

### Uses:

- Scientific discovery. Ideal to uncover characteristics or true properties of the biological processes of the studied phenomenon.
- Useful to judge the individual relevance of each quantitative measure in impacting the response of interest.

## pattern recognition

### Goal:

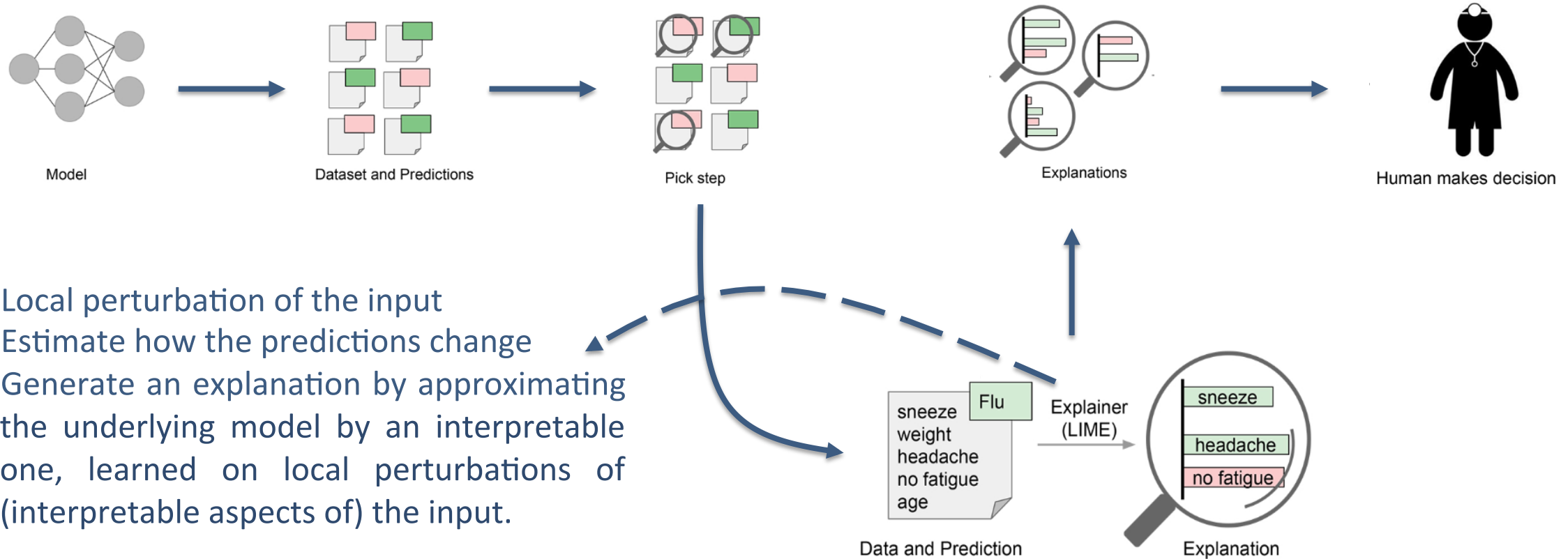
- Identify most predictive variable sets (out-of-sample prediction performance)

### Uses:

- Pragmatic forecasting of biological processes.
- Tends to concern less regarding the data-generating process.

# Diagnosing Features ( $\approx$ Interpretability)

## Local Interpretable Model-Agnostic Explanations [*LIME*]



- Local perturbation of the input
- Estimate how the predictions change
- Generate an explanation by approximating the underlying model by an interpretable one, learned on local perturbations of (interpretable aspects of) the input.