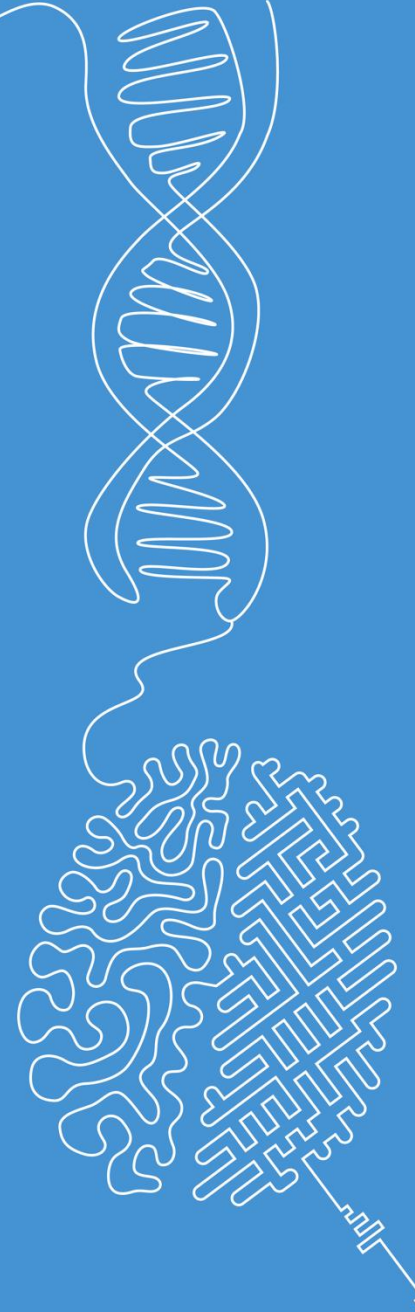


Summary

Machine Learning

Norman Juchler



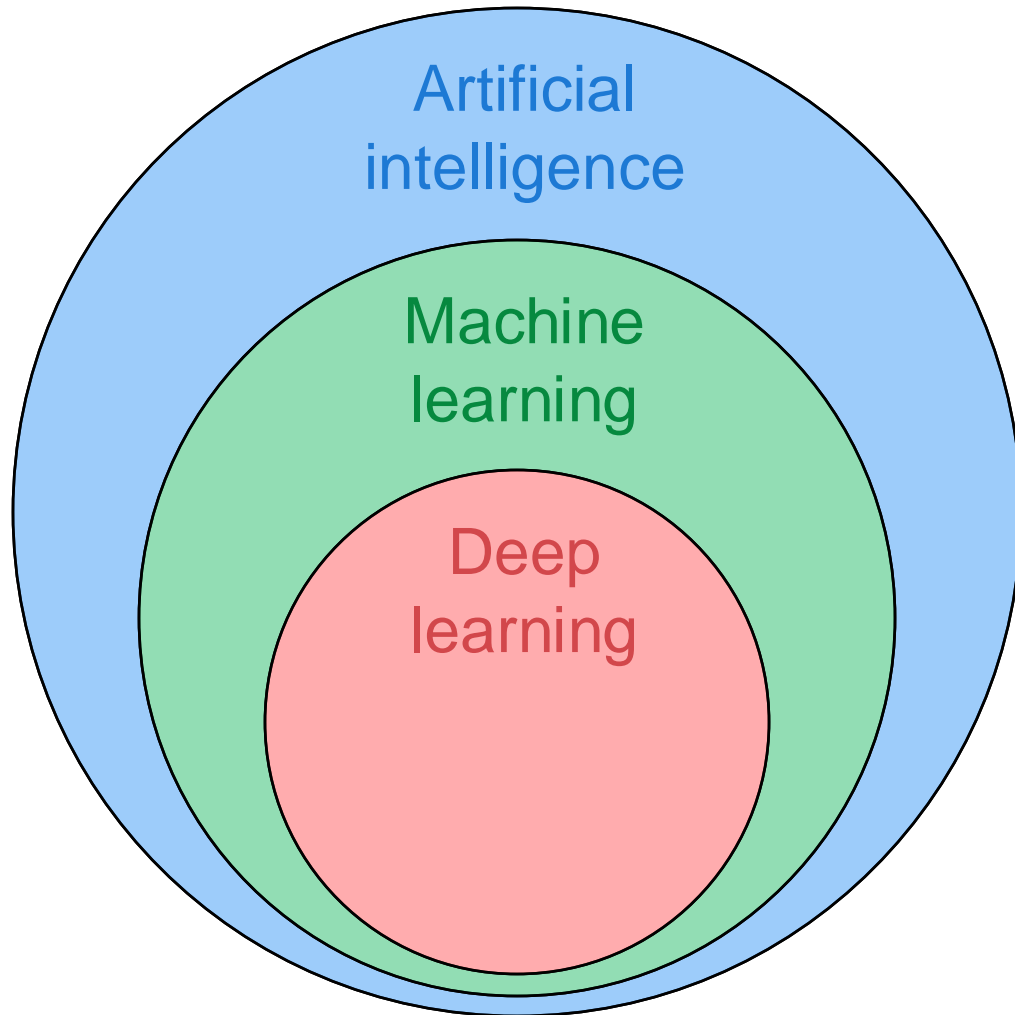
INSTRUCTIONS: Highspeed recap of the semester

Notes:

- The following slides summarize the topics covered in this module
- Although some slides have been modified, they are generally the same as the ones you already know from the corresponding lecture.
- It is probably advisable not to use this batch of slides for studying, as the slides are torn out of their context...

Fundamental concepts

What is artificial intelligence?

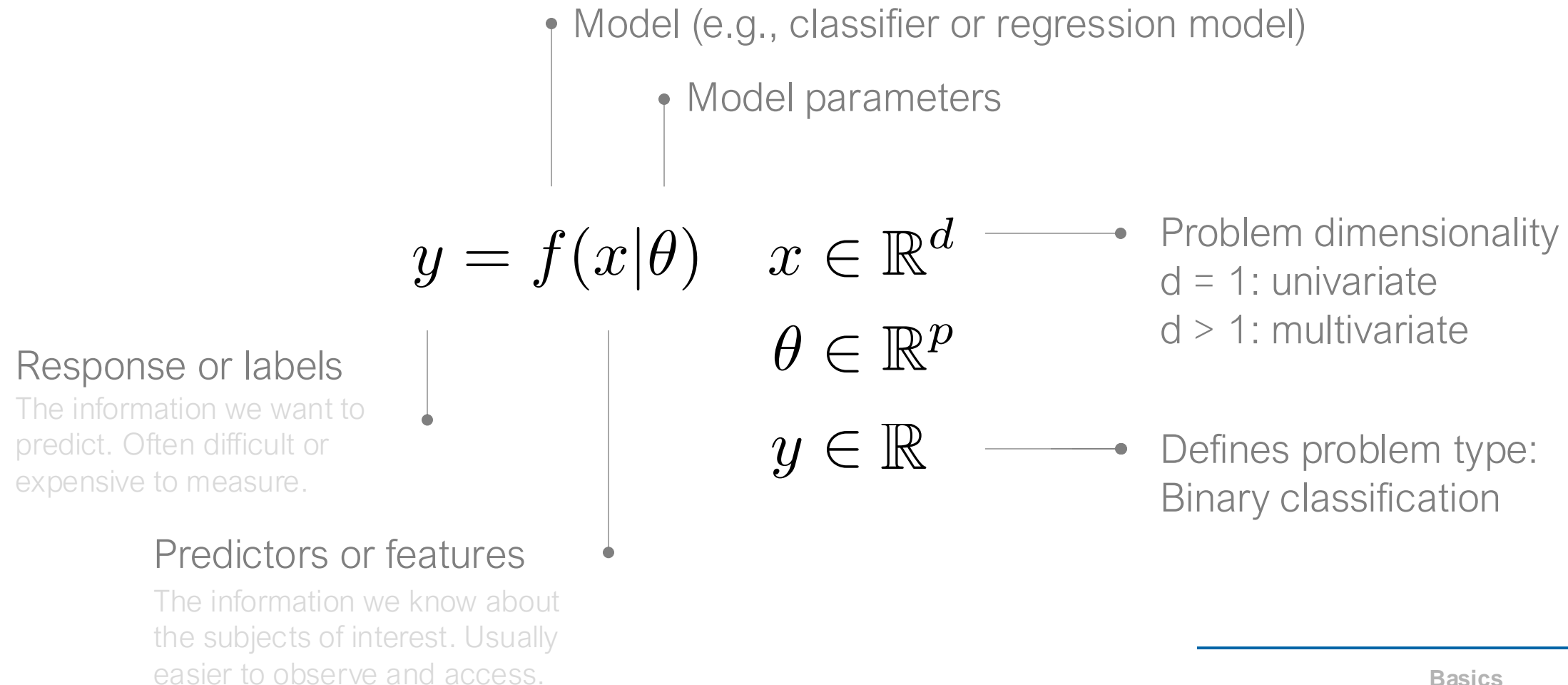


Distinctions:

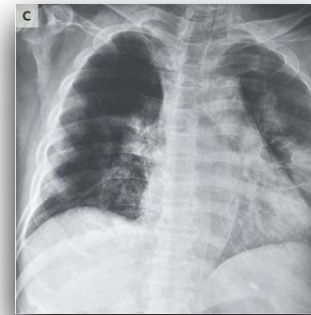
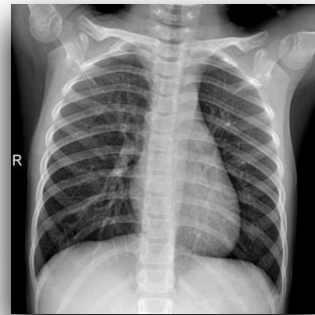
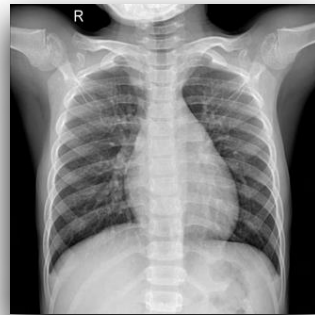
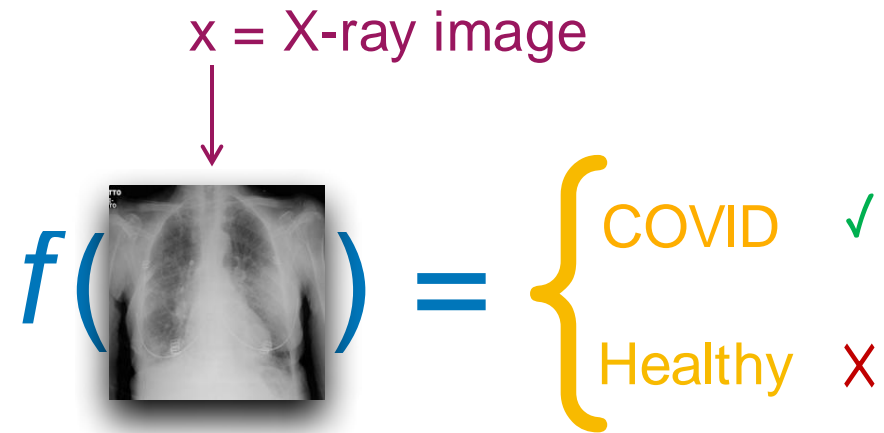
- **Statistics \neq machine learning**
 - ML: the focus is on prediction
 - Statistics: focus is on inference
 - ML theory and algorithms rely on statistics and probability theory
 - ML often makes fewer formal assumptions about data
- **Machine learning \neq statistical learning**
 - The terms are often used interchangeably, but have slightly different meanings
 - SL: includes the formulation of stochastic models for the data generating processes
 - SL: besides the use for prediction, the focus is also on formally understanding the data

Formal introduction of a machine learning model

- In ML, a **model** is a mathematical representation, function or algorithm that defines the relationship between input data and the desired output.

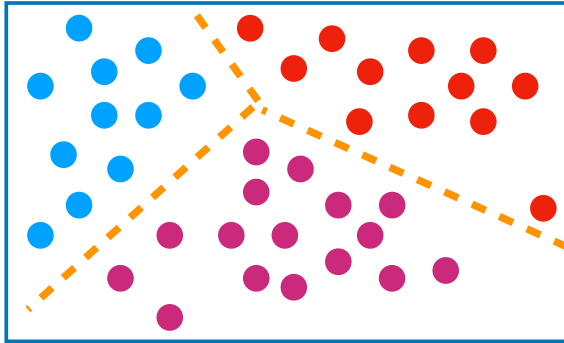


How a machine learning model operates



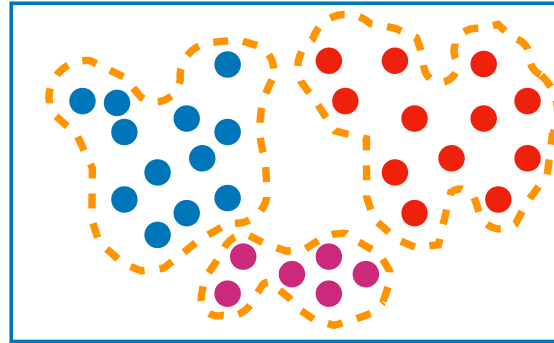
Machine learning paradigms

Supervised Learning



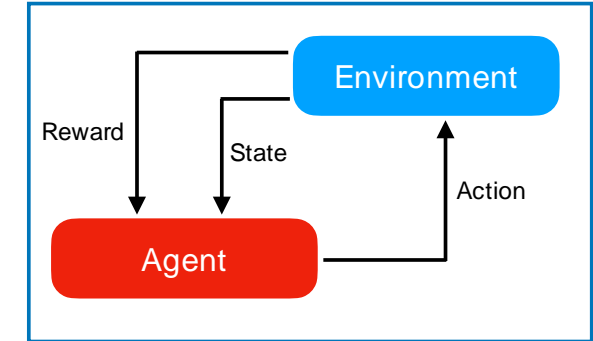
$$\{x_i \in \mathbb{R}^d, y_i \in \mathbb{R}\}_{i=1}^n$$

Unsupervised Learning






$$\{x_i \in \mathbb{R}^d\}_{i=1}^n$$

Reinforcement Learning



Supervised Learning



x_i :   
 y_i : panda lion dog

Unsupervised Learning



x_i : 

Reinforcement Learning



Agent: Kid
 Environment: Neighbourhood street

Bias and variance in machine learning

- The **bias** is a property of the model that causes it to miss relevant relations between the features and the target output, e.g. because of wrong assumptions in the model.
- The **variance** measures how susceptible the model is to small changes in the training data. A model with high variance will change a lot if the training data is changed slightly.

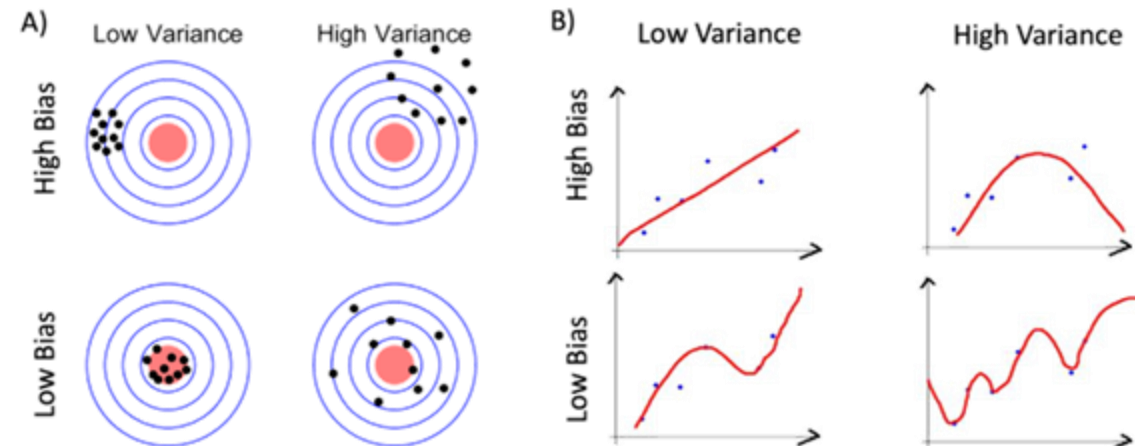
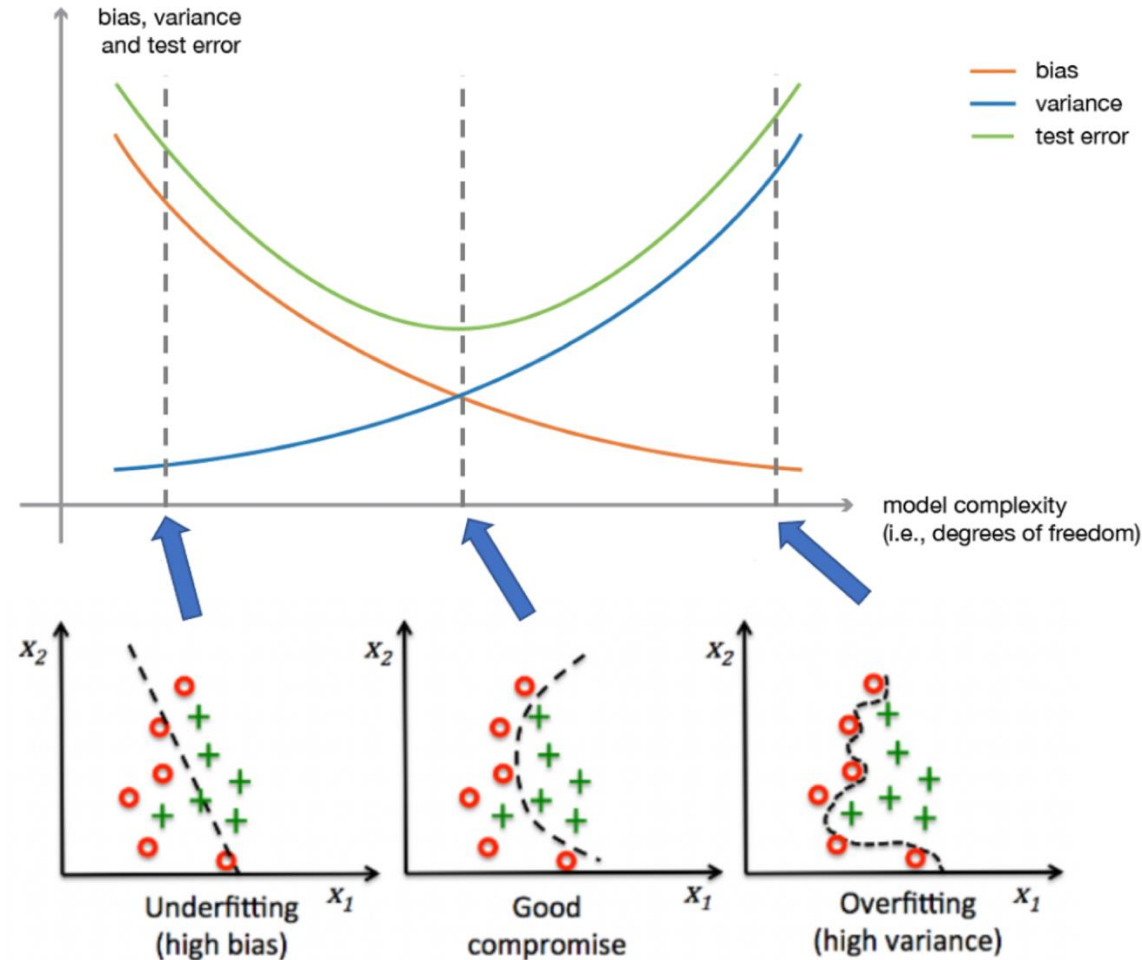


Figure 4: A: Statistical properties of bias and variance. B: Example of regression models for different combinations of bias and variance (assuming that the red curve in the lower left is the true data generating process).

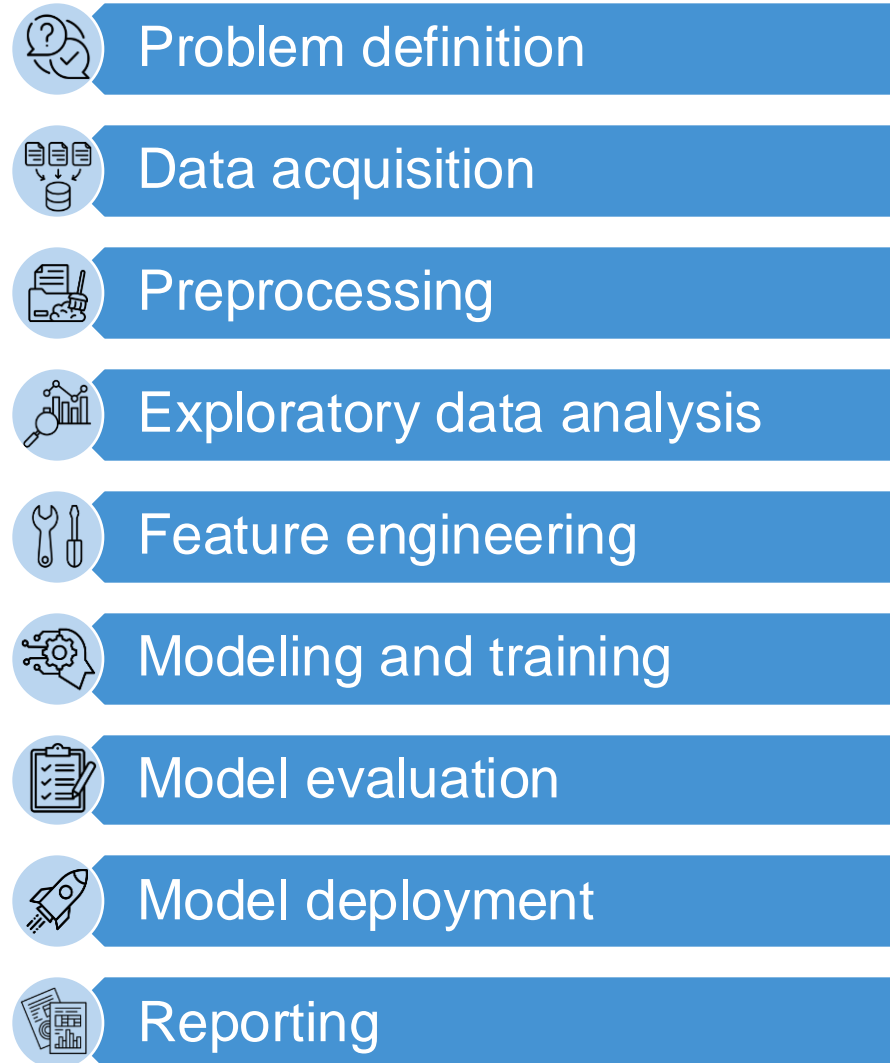
Model complexity vs. bias-variance tradeoff

We need to find the sweet spot in the middle!



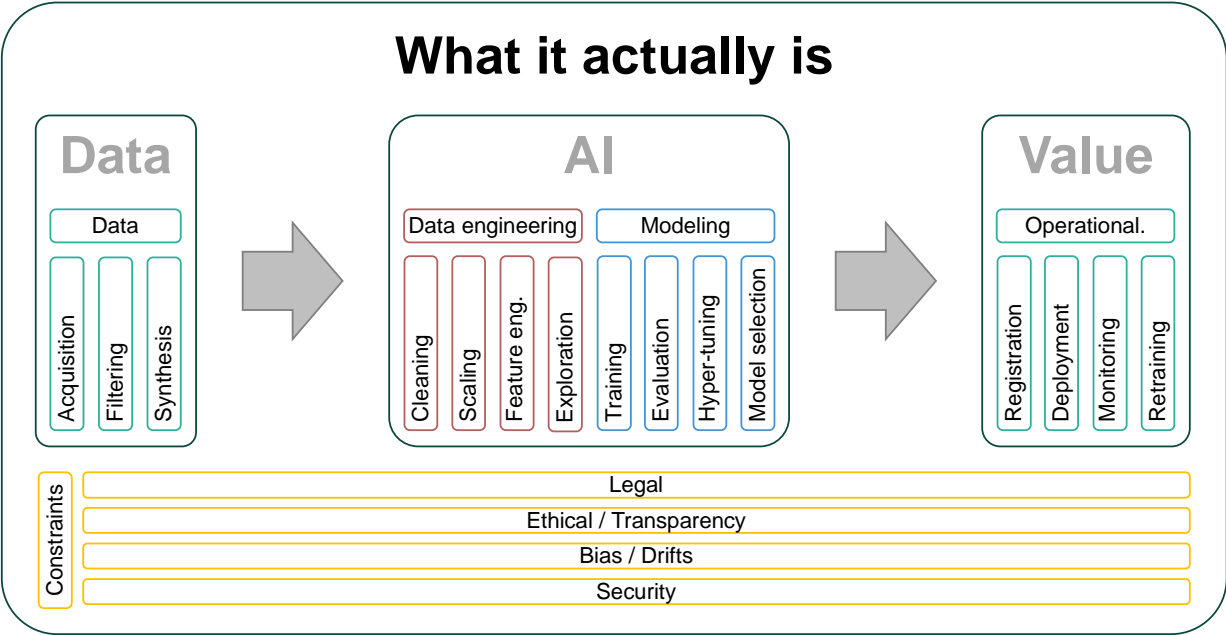
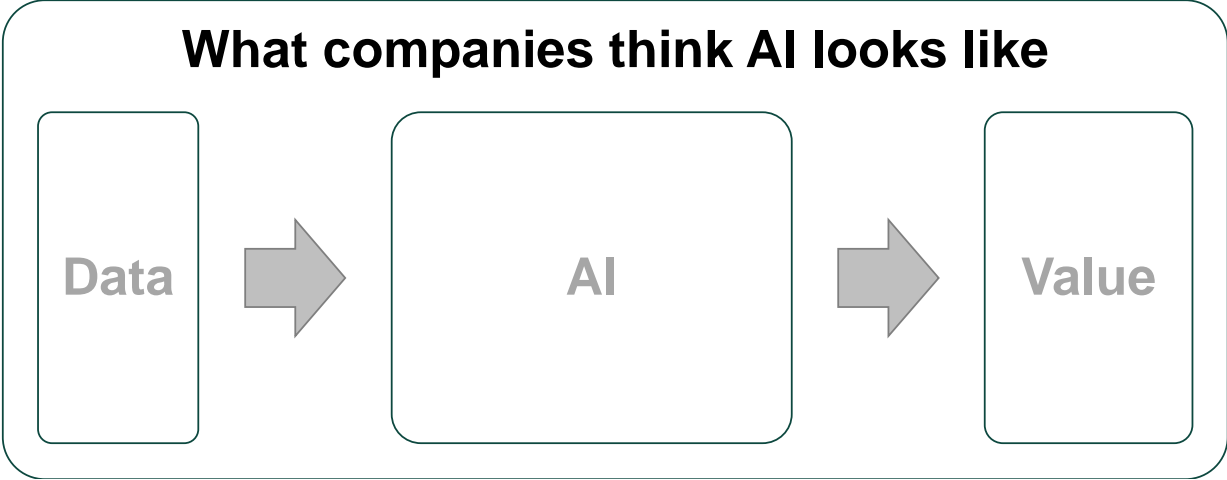
Data science workflow

The data science workflow



- Implementing the workflow is an **iterative process**
- Repeat or revisit stages of the workflow based on new findings or feedback.

Another perspective...

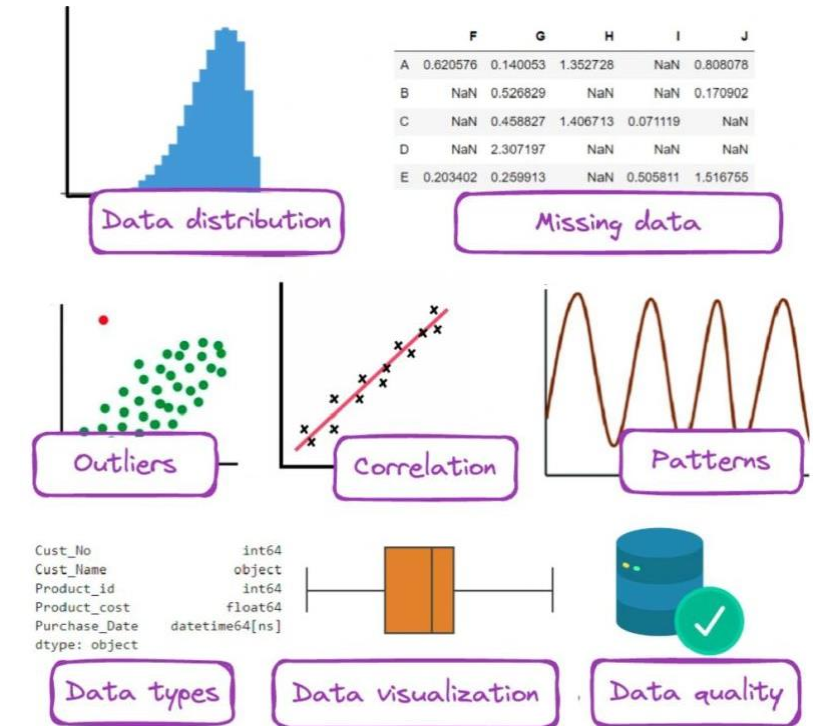


Purpose of exploratory data analysis (EDA)

- **Understand the data's structure:** Identify patterns, relationships, trends, and anomalies in the dataset.
- **Detect missing or incorrect data:** Uncover issues such as missing values, outliers, or inconsistencies that need to be addressed.
- **Guide feature selection and engineering:** Inform decisions about which features are important or need transformation.
- **Generate hypotheses:** Formulate potential hypotheses and insights for further analysis or modeling.
- **Validate assumptions:** Check if the data fits certain assumptions required by statistical or machine learning models.

Means of EDA

- Summary statistics
- Data visualization
- Correlation analysis (between features)
- Univariate analysis (between features and target)



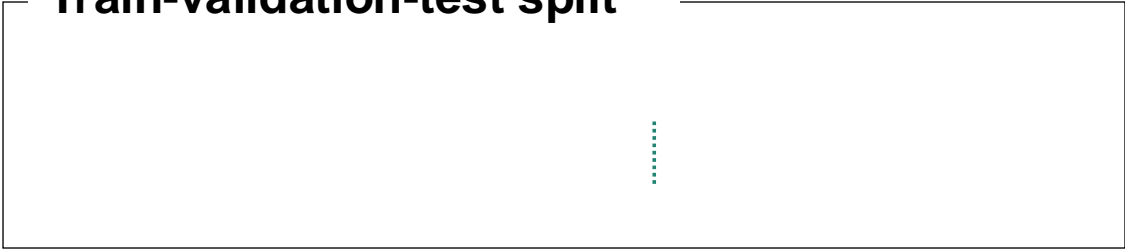
Training and validation

Different data splits

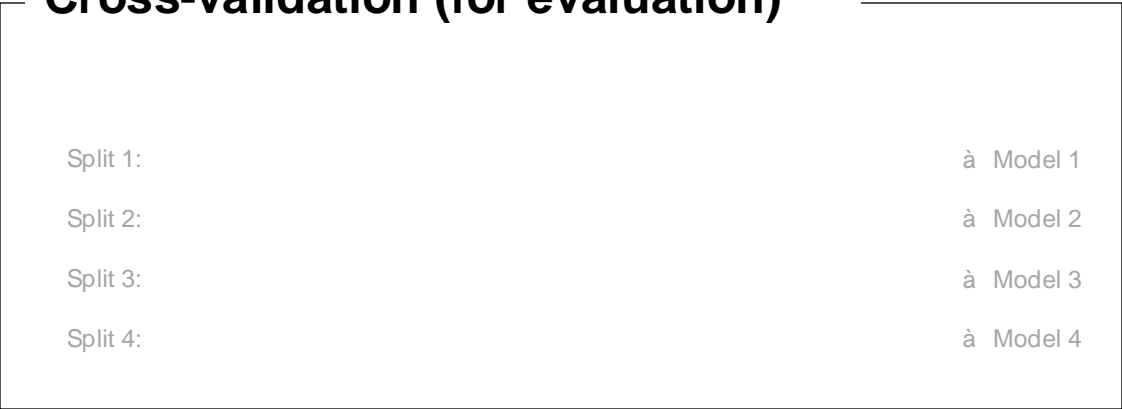
Train-test split



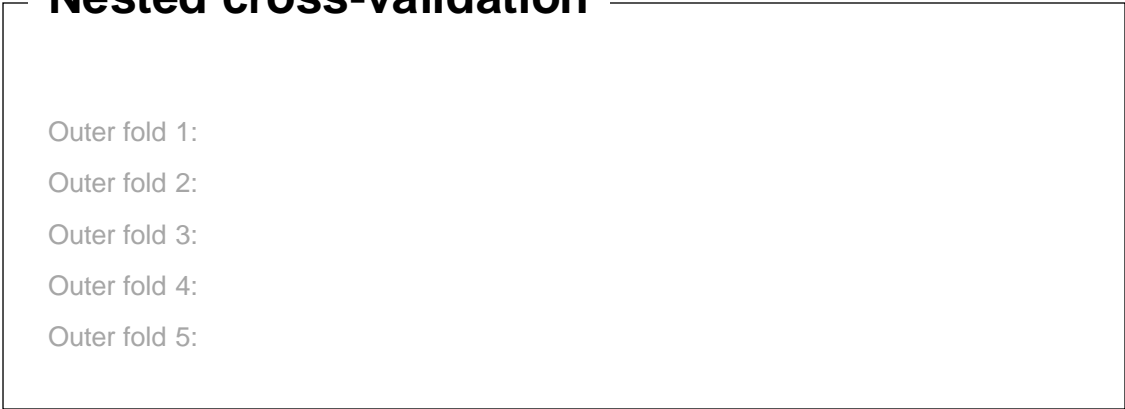
Train-validation-test split



Cross-validation (for evaluation)



Nested cross-validation



Evaluation metrics: for regression

- **Mean squared error (MSE):**

- \hat{y}_i : predicted values
- y_i : true values (from the training data)

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

- **Root mean squared error (RMSE):**

- Has the same unit as y

$$RMSE = \sqrt{MSE}$$

- **Mean absolute error (MAE):**

- Penalizes bigger errors less

$$MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

Evaluation metrics: for classification

■ Accuracy:

$$\frac{TP + TN}{TP + FP + FN + TN} = \frac{TP + TN}{P + N}$$

Fraction of correctly classified events

■ Precision:

$$\frac{TP}{TP + FP}$$

Fraction of relevant instances among the retrieved instances

■ Sensitivity or recall (true positive rate, TPR)

$$\frac{TP}{TP + FN} = \frac{TP}{P}$$

Fraction of relevant instances that were correctly retrieved

■ Specificity (true negative rate, TNR)

$$\frac{TN}{TN + FP} = \frac{TN}{N}$$

Fraction of negative instances that were correctly retrieved

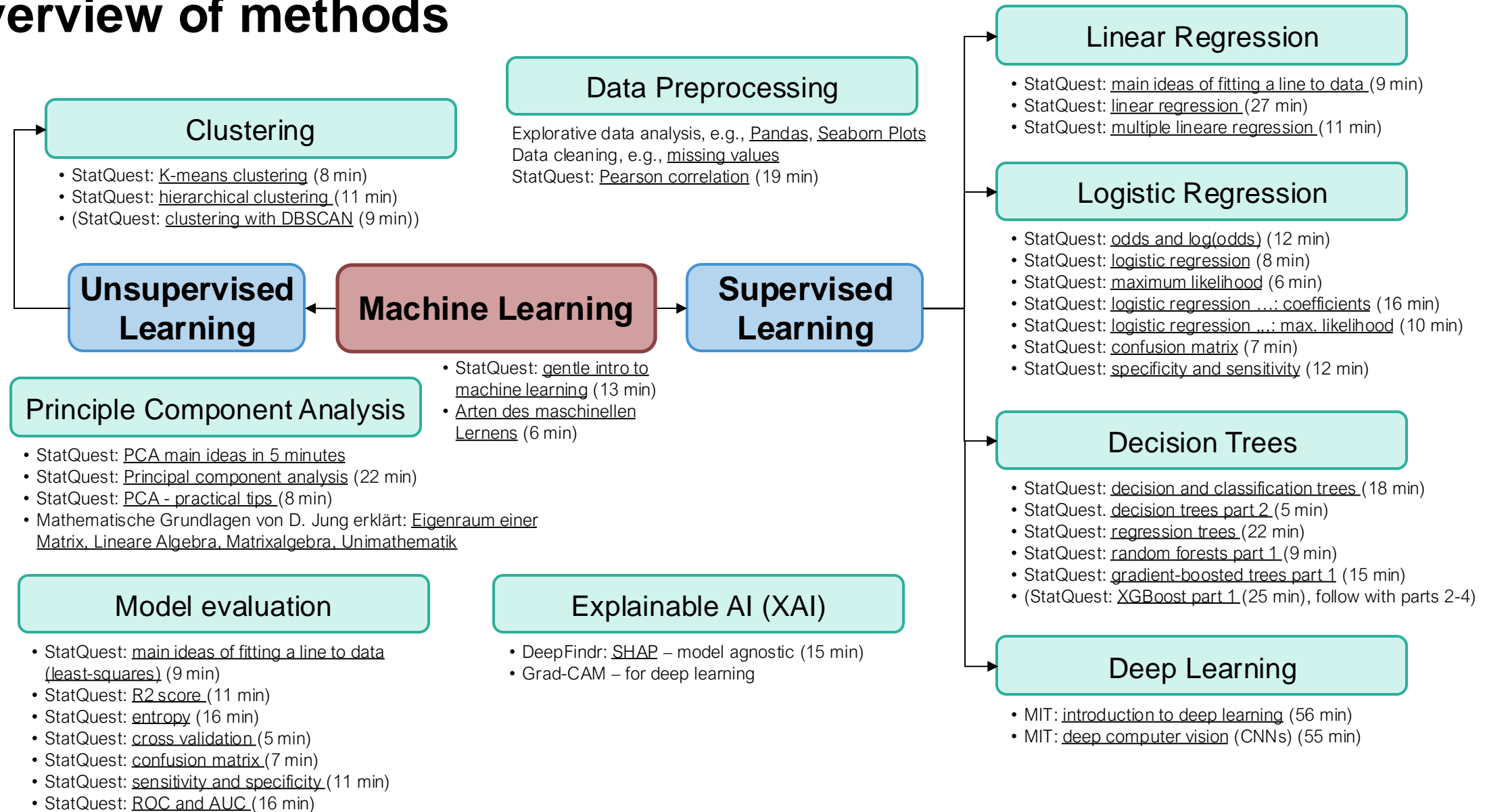
- The sensitivity quantifies the avoidance of false negatives, the specificity the same for false positives.

		True condition ()	
		P Condition positive	N Condition negative
Prediction ()	Predicted condition positive	TP True positive	FP False positive Type I error
	Predicted condition negative	FN False negative Type II error	TN True negative

Contingency table
(aka confusion matrix)

Specific methods

Overview of methods



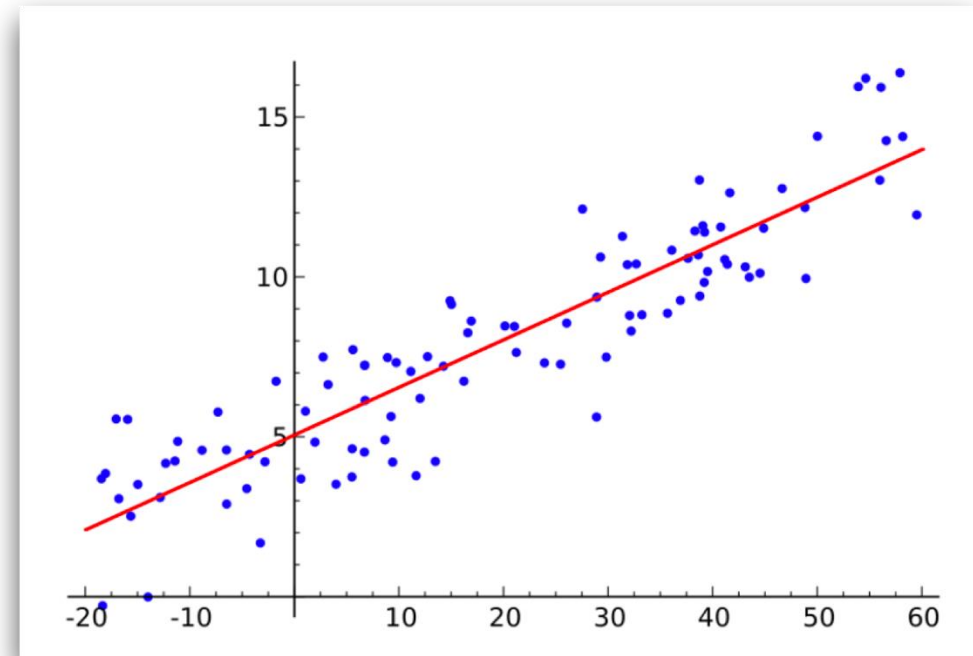
Linear models

- Assume we have a set of p features: $\mathbf{x}_i := (x_{i1}, x_{i2}, \dots, x_{ip})$
- We want to use them to predict a target variable y
- The simple assumption we can make is (model ansatz):

$$y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}$$

$$f_i(\mathbf{x}, \boldsymbol{\beta}) = \boldsymbol{\beta} \mathbf{x}$$

- Where the β_i are unknown parameters that we want to determine from the data



Loss functions

- A loss function (a.k.a. cost function) in the context of machine learning usually measures how well the predicted values match the true target values.

$$\mathcal{L}(y - \hat{y})$$

- For regression, we used the residual sum of squares (RSS) as loss function.

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

- Sometimes, it is more meaningful to compute the **mean squared error** (MSE)

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

- (Note: both yield the same optimal solution!)

Numerical methods for optimization: Gradient descent

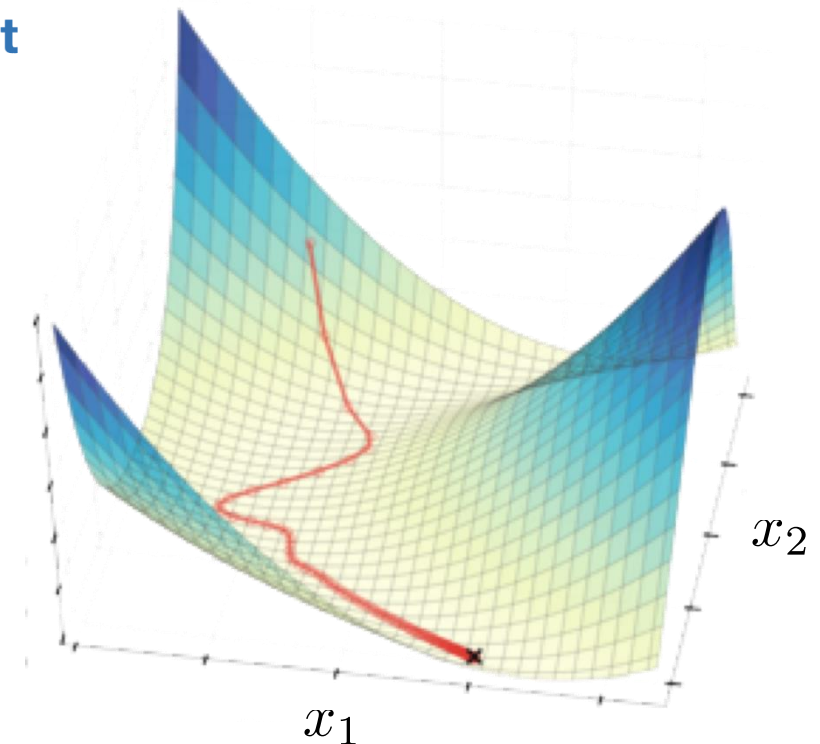
- Gradient descent is another iterative method to find solutions to $f'(\mathbf{x}) = 0$
- The method advances a small step in the direction of the gradient each time
- Gradient descent is applicable under more relaxed conditions

$$\mathbf{x}_{n+1} = \mathbf{x}_n - \gamma_n \nabla f(\mathbf{x}_n)$$

Gradient = direction of steepest ascent

Learning rate (adjustable)

$$\mathcal{L}(\mathbf{x}) = f(\mathbf{x})$$



Regularized regression

- Idea: Prevent overfitting by penalizing large coefficients, encouraging simpler models that generalize better.
- How? Modify the loss function!

Regularized loss = Original loss + Penalty term on coefficients

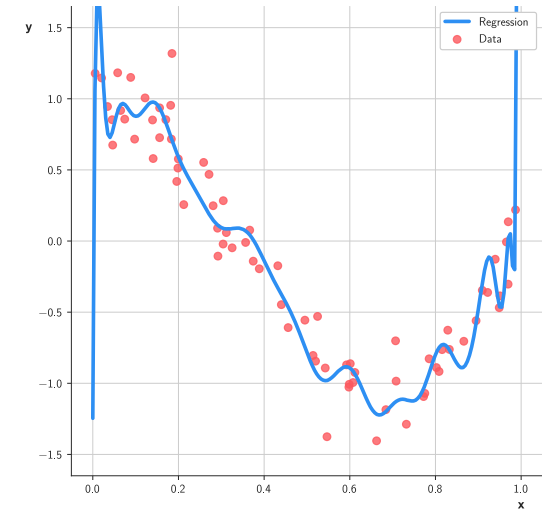
- Example: (Ridge regression)

$$\mathcal{L}(\beta|X, y) = MSE(\beta|X, y) + \lambda \sum_{j=1}^p \beta_j^2$$

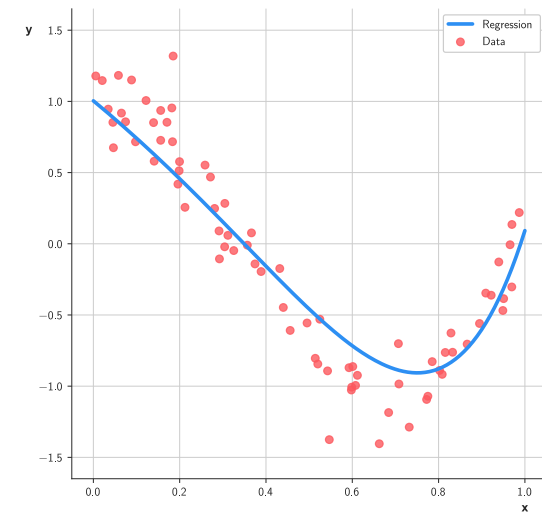
Recap:

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i|\beta))^2$$

- This introduces a **hyperparameter λ** (or α in sklearn):
 - Controls the strength of regularization
 - Larger α increases penalty, leading to smaller coefficients



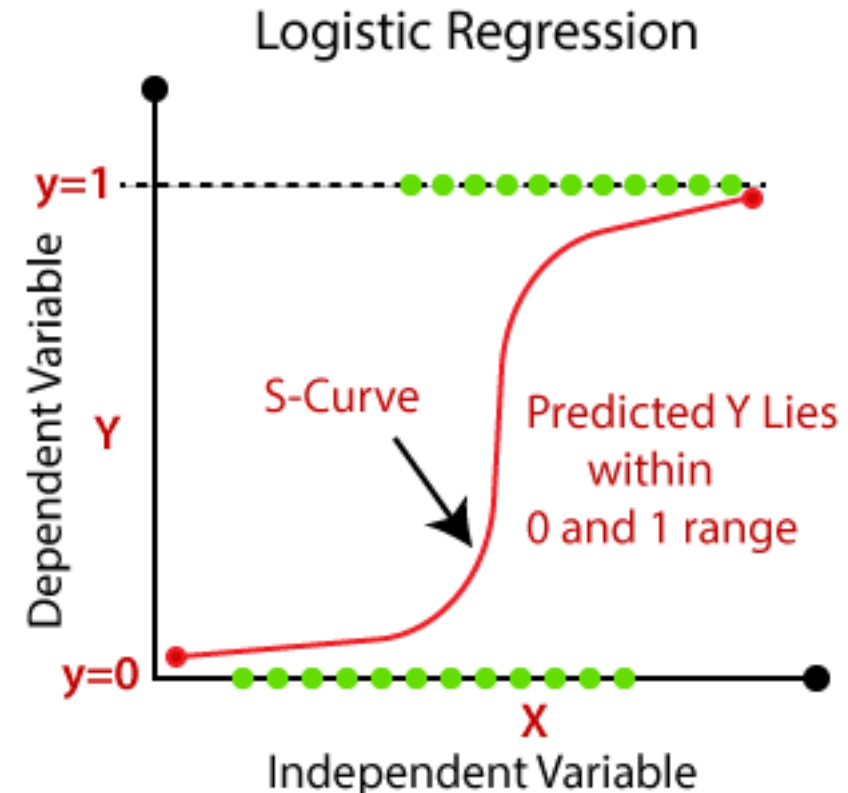
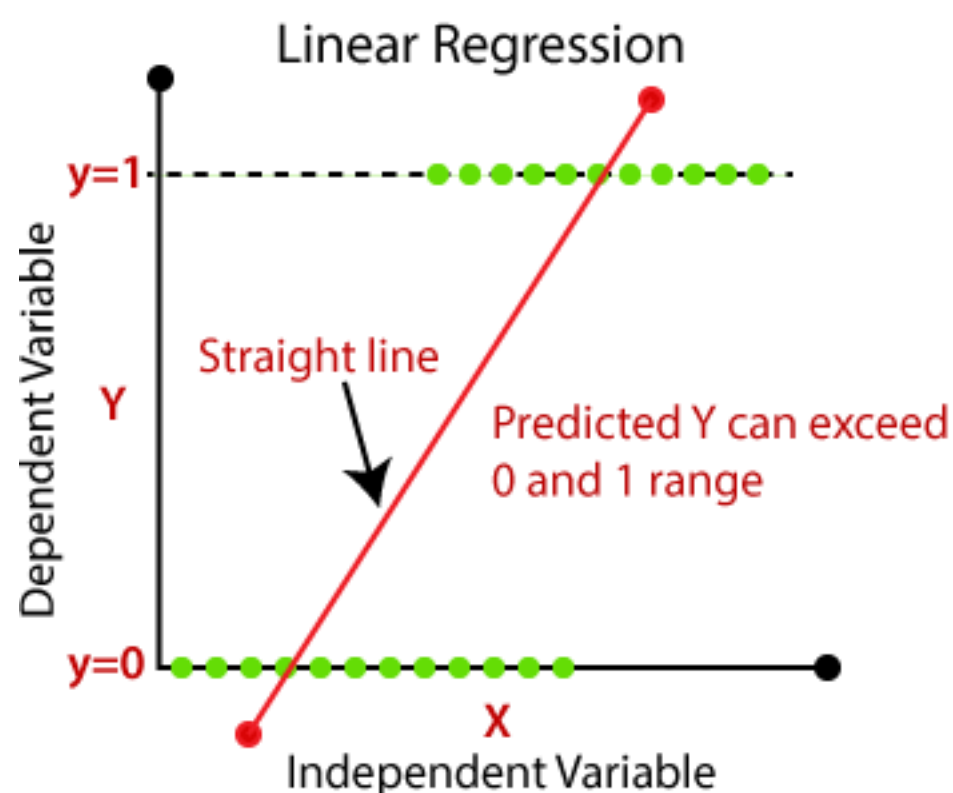
Polynomial regression without regularization showing overfitting



Polynomial regression with regularization (ridge), which in this case prevents overfitting.

Logistic regression for binary classification

- It makes little sense to compute a linear regression for the classification.
- Use a logistic regression instead!

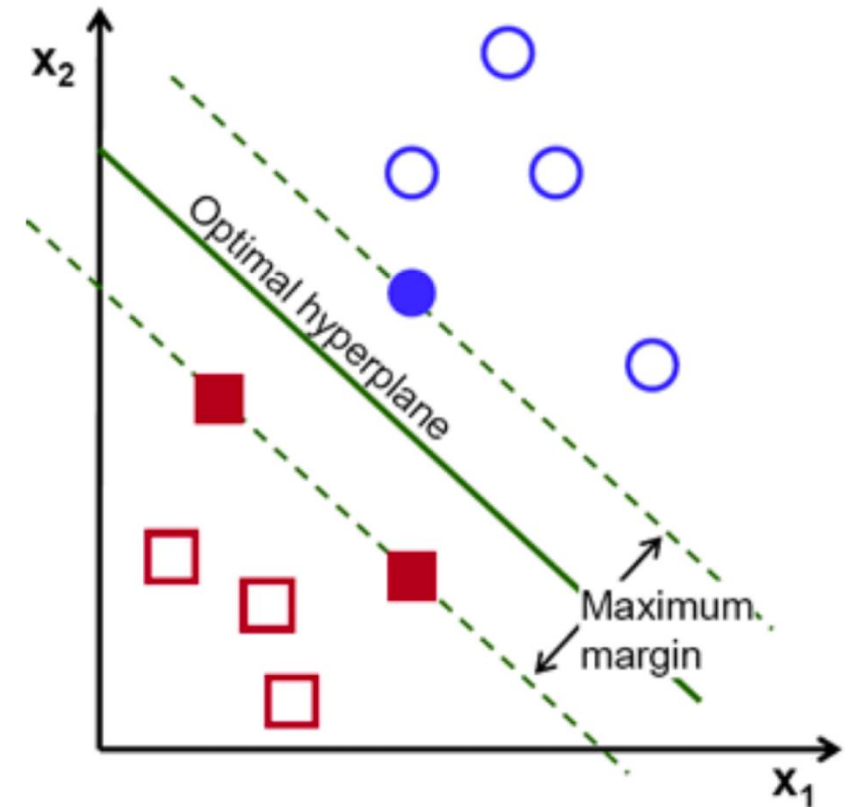


Support vectors and margin (SVMs)

- Linear classifier: Find **hyperplane** that best separates the classes in the feature space
- Any hyperplane can be expressed as

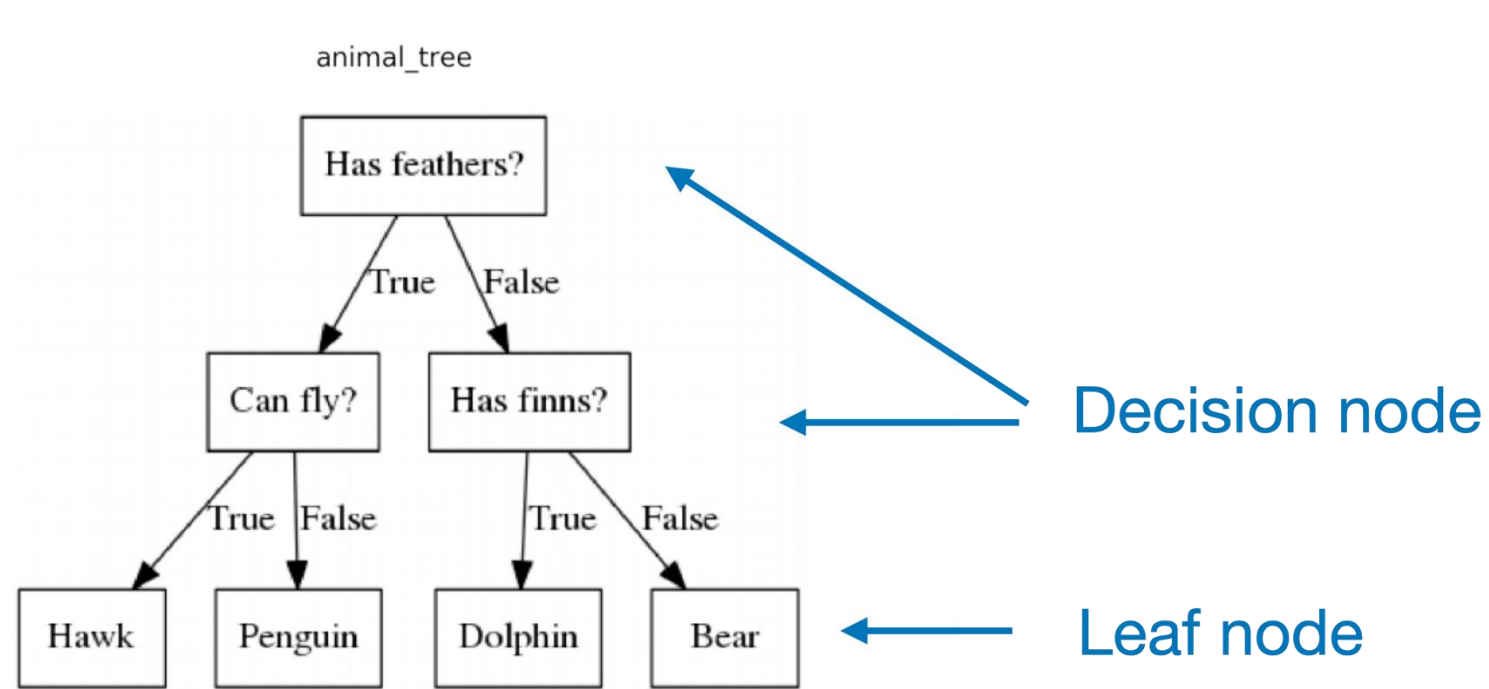
$$\mathbf{w}^T \mathbf{x} - b = 0$$

- Geometric interpretation: \mathbf{w} is the normal of the decision boundary / hyperplane!
- **Definition:** The nearest points of each class to a hyperplane are called **support vectors**
- **Key idea:** SVMs find the hyperplane with maximal distance (**margin**) to the support vectors.



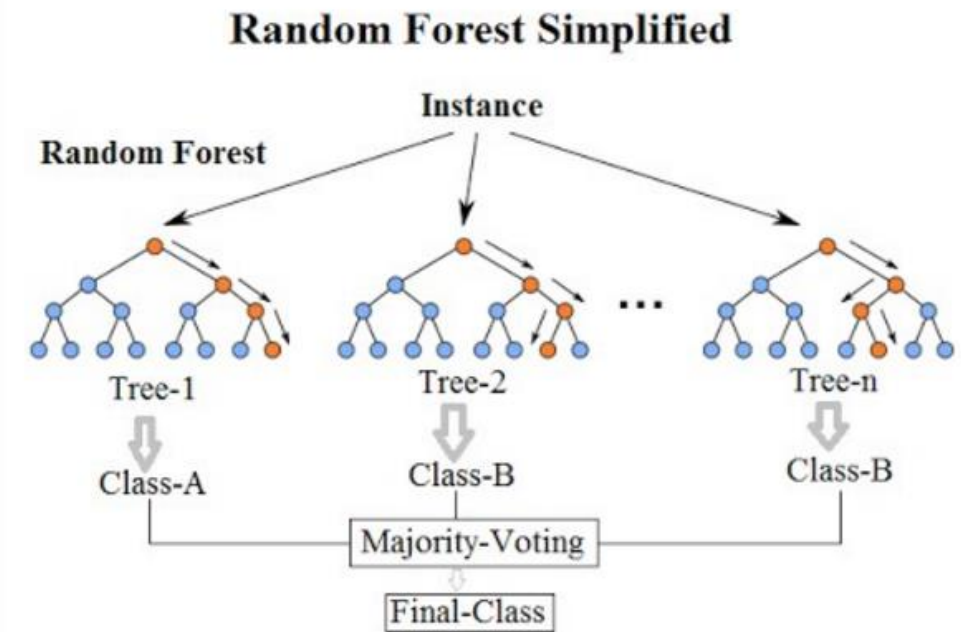
Decision trees

- Main idea: Algorithmically learn to construct a set of decision rules in the form of a tree.



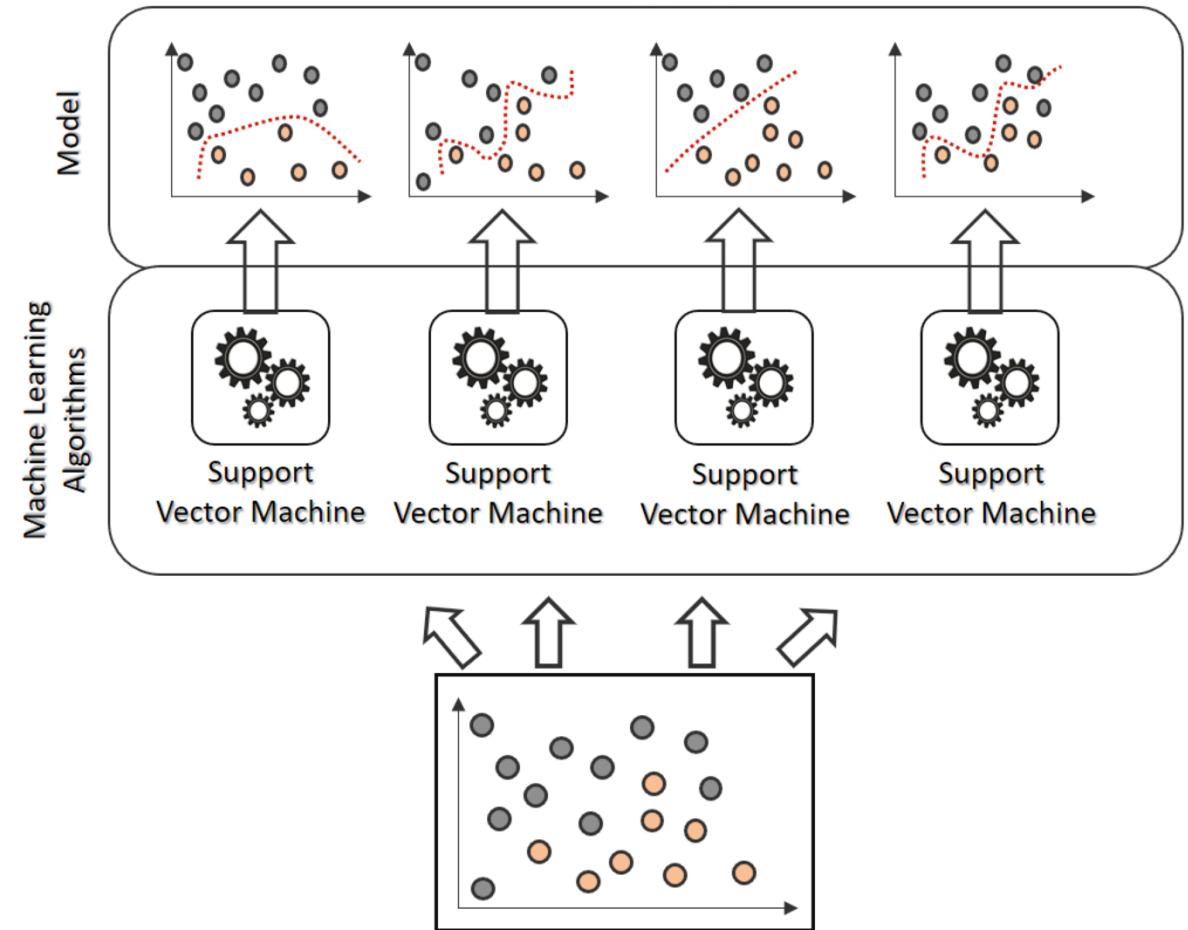
Random forest: Algorithm

- An ensemble of decision trees trained with bagging.
- Many decision trees are generated during training.
- Tree bagging is used to train trees (reduces variance and overfitting).
- Feature bagging: Each tree only uses a random subset of features, which ensures that trees remain decorrelated.
- Used for both regression and classification
- Training is easily parallelizable
- Robust to noise and outliers



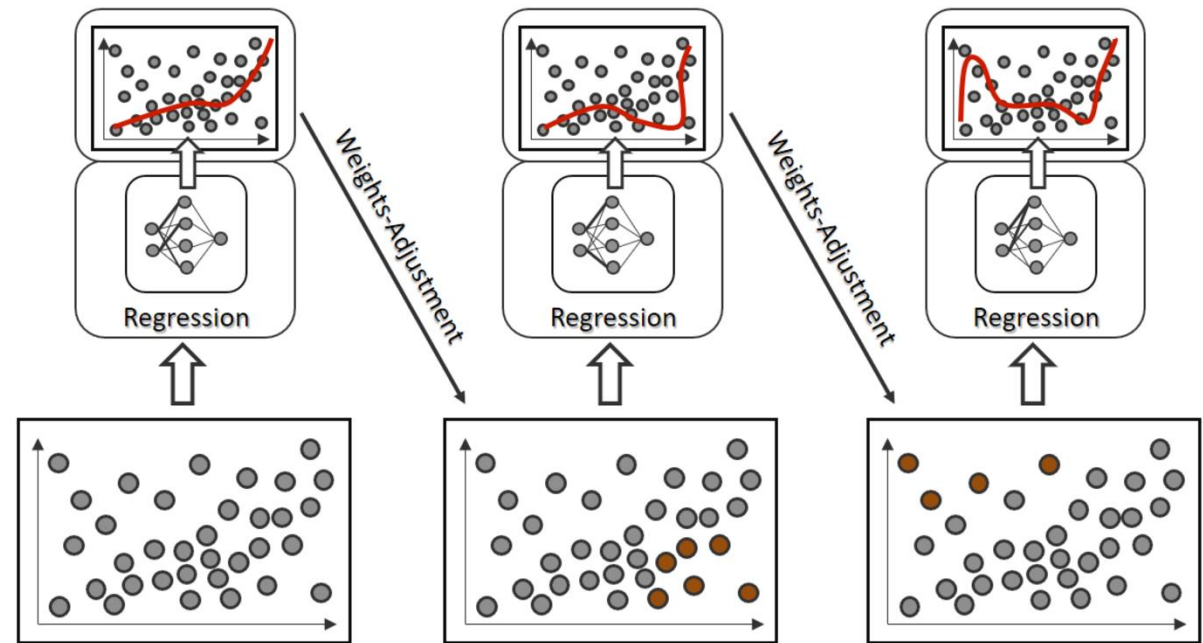
Ensemble methods: Bagging

- Short for bootstrap aggregating
- An ensemble technique where the constituent models are
 - of the same type, but
 - trained on different randomly sampled data subsets.
 - Bootstrapping \Leftrightarrow sampling with replacement
 - Final prediction by aggregation (clf: majority vote, reg: or averaging)
- Overall model becomes more robust/reduces overfitting.
- All constituent models can be trained in parallel.



Ensemble methods: Boosting

- An ensemble method in which the models are trained sequentially, with more weight given to misclassified samples.
- Advantage: Performance often higher than bagging
- Disadvantages (compared to bagging):
 - More susceptible to overfitting
 - Not parallelizable
- Popular methods:
 - AdaBoost (focus on misclassified cases)
 - Gradient boosting (focus on residuals)
 - XGBoost (popular instance of GB)



k-nearest neighbors for classification (and also regression)

- kNN uses local information from *nearby* training examples to predict new labels.

Notes:

- It is common to weight neighbors with the inverse of their distance, such that closer points have greater influence:

$$\text{weight} = \frac{1}{\text{distance}}$$

- Weights should also be applied in case of imbalanced data.
- Since feature-space distances combine different units, normalizing the training data is key!

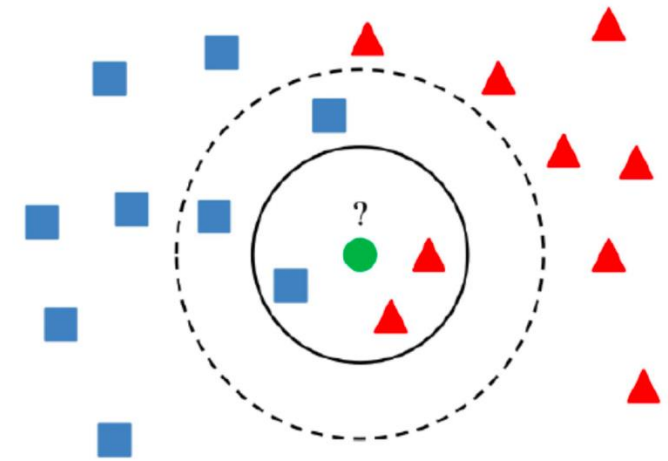
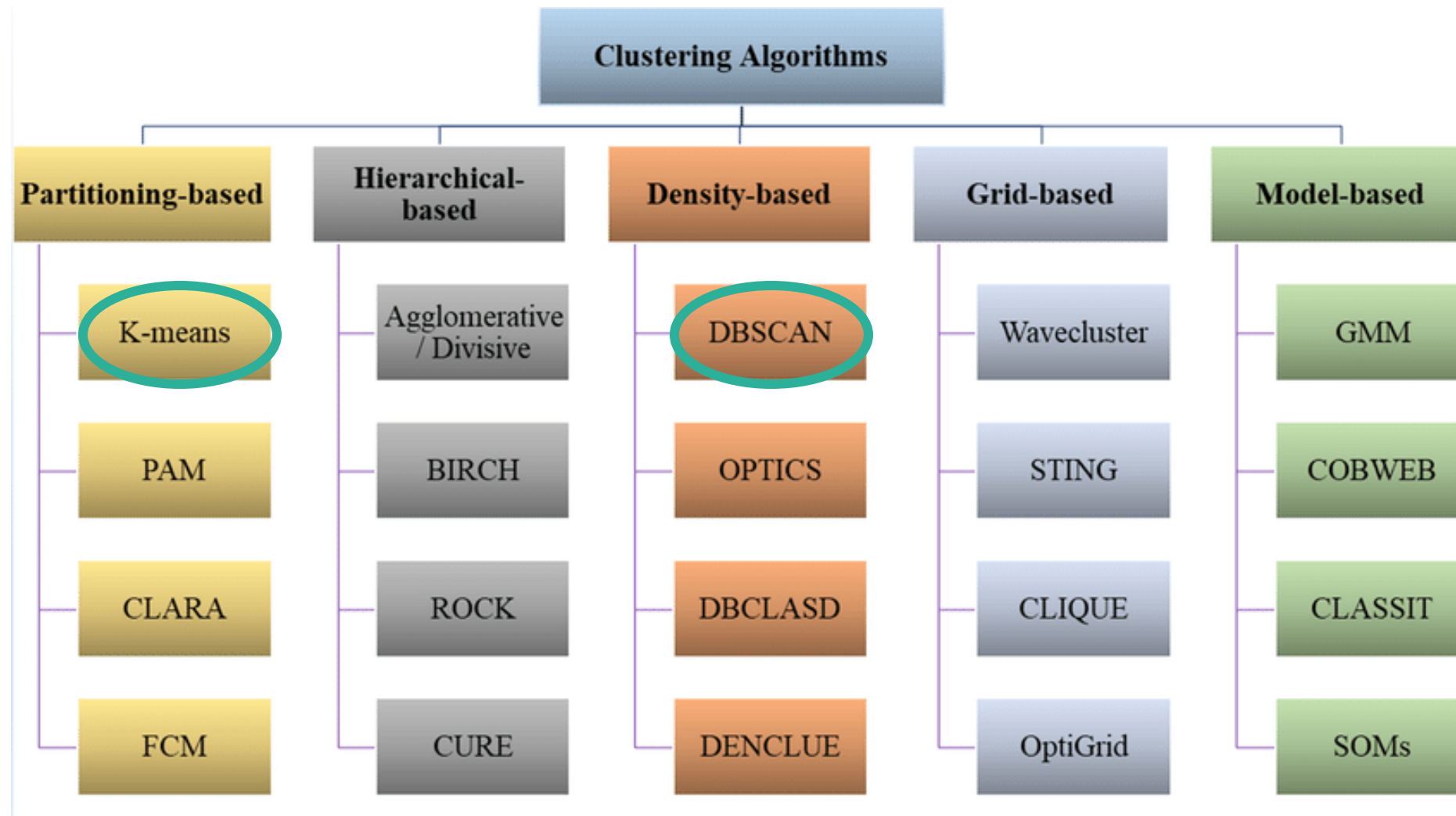


Illustration of kNN for supervised learning:
Search for the k nearest neighbors of the **inference point** for:

— k = 3

- - k = 5

Clustering: Method overview



Distance metrics: How to quantify similarity / proximity?

- Euclidean

$$\sqrt{\sum_{i=1}^n (x_i - y_i)^2}$$

- Manhattan

$$\sum_{i=1}^n |p_i - q_i|$$

- Minkowski

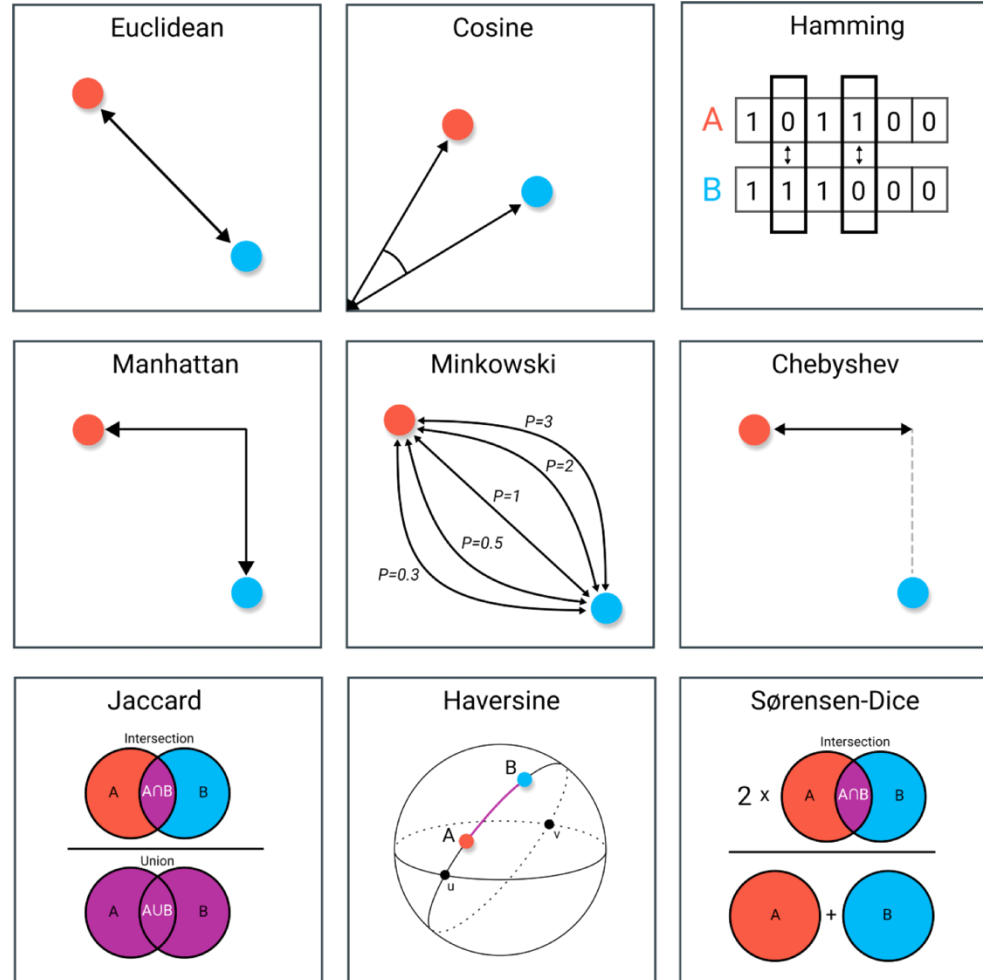
$$\left(\sum_{i=1}^n |x_i - y_i|^p \right)^{\frac{1}{p}}$$

- Cosine

- Hamming

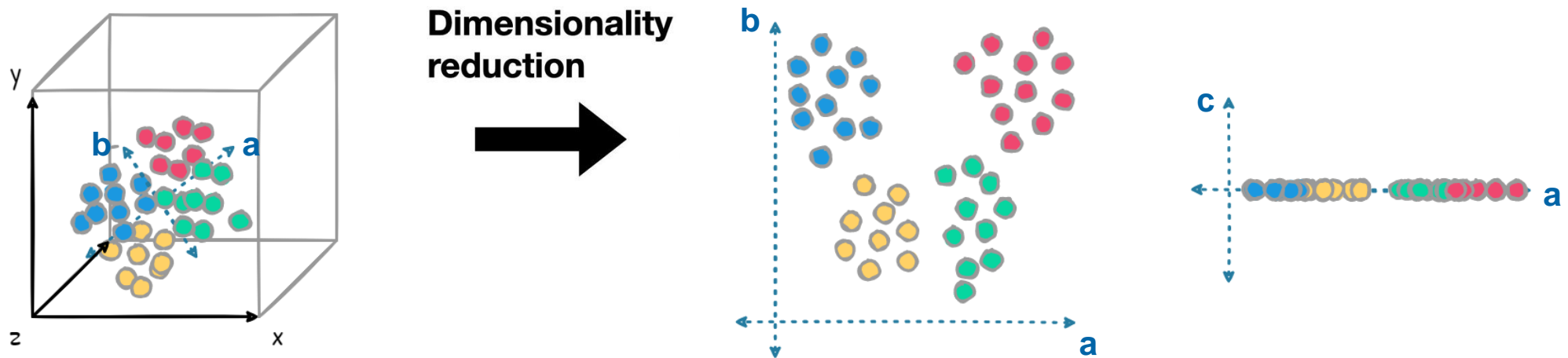
- ...

- and any application specific distances...



Dimensionality reduction through projection

- Projection-based dimensionality reduction involves mapping high-dimensional data onto a lower-dimensional subspace while preserving as much of the data's structure as possible.

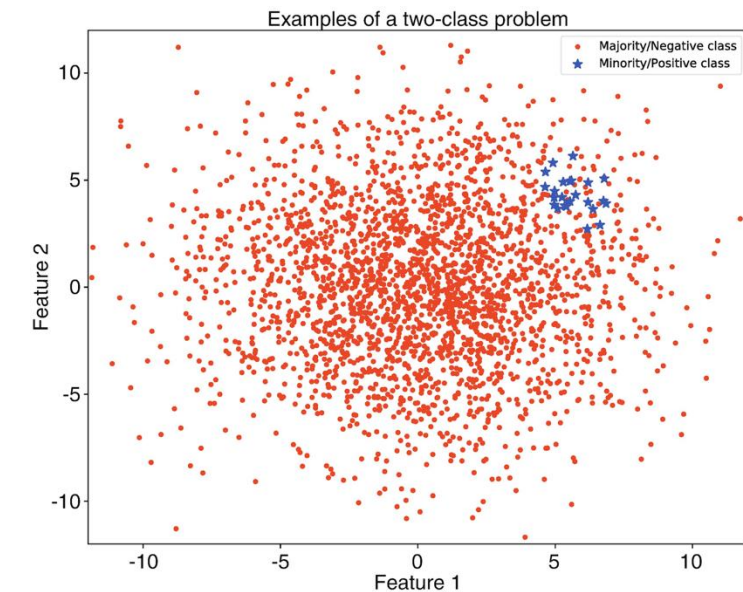
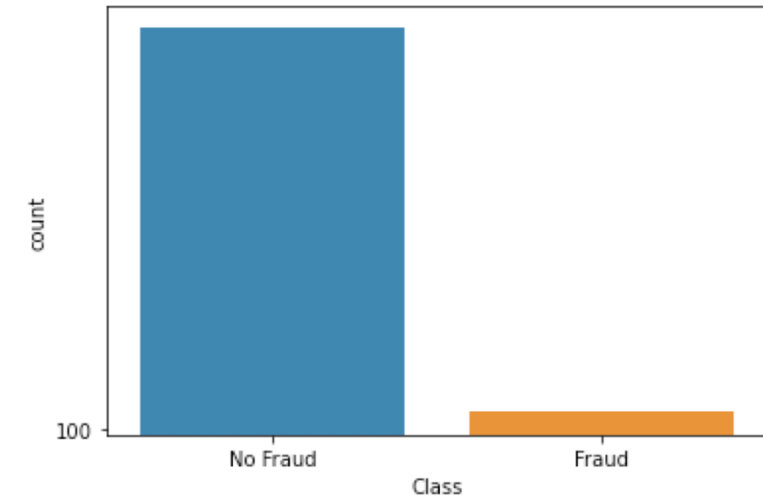


Two possible projections from a 3D space onto a 2D subspace. Here, the first subspace is the plane in which the data points lie (plane with the axes a-b). The second one is the plane a-c, which is perpendicular to the a-b plane.

Miscellaneous

Imbalanced data

- **Definition:** Imbalanced data refers to a situation where certain categories, values, or groups in a dataset are underrepresented compared to others.
- **Examples:**
 - Credit card fraud
 - Natural disasters
 - Tumor cells
 - ...
- Many (if not most) real-world datasets exhibit imbalance
- Often, we are interested in the outlier (or rare) events



The three levels of model parameters



(Model Design + Hyperparameters) → Model Parameters

Examples

Function class
(e.g. cosign
or a polynomial)

Degree of the
polynomial

The polynomial's
parameters (a_0 , a_1 , etc.)

Core ideas of active learning

- **Iterative process:**

The active learner is trained on an initial labeled dataset, then iteratively selects and requests labels for the most uncertain or impactful samples.

- **Uncertainty sampling:**

The model identifies data points where it has the least confidence in its predictions (e.g., probabilities close to 0.5 in binary classification).

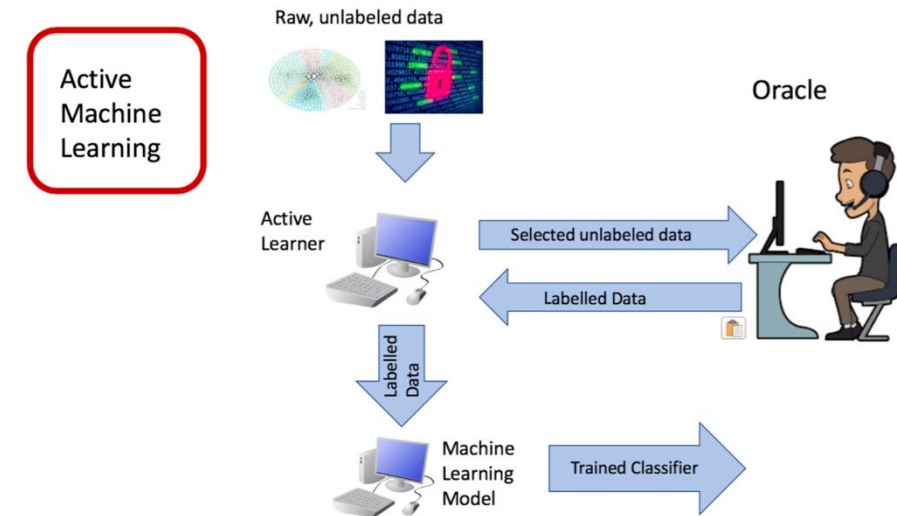


Illustration of active learning. The use of machine learning during a data labeling process. Image source: [Link](#)

The vision of explainable ML

