

# Machine Learning Review

Paolo Favaro

# Contents

- Revision of basic concepts of Machine Learning
- Based on **Chapter 5** of Deep Learning by Goodfellow, Bengio, Courville

# Resources

- Books and online material for further studies
  - Machine Learning @ Stanford (Andrew Ng)
  - **Pattern Recognition and Machine Learning**  
by Christopher M. Bishop
  - **Machine Learning: a Probabilistic Perspective**  
by Kevin P. Murphy

# Learning Pillars

- Supervised learning
- Semi-supervised learning
- Self-taught learning (unsupervised feature learning)
- Unsupervised learning (+self-supervised learning)
- Reinforcement learning

# Definition

- Mitchell (1997)

*A computer program is said to learn from **experience E** with respect to some class of **tasks T** and **performance measure P**, if its performance at tasks in  $T$ , as measured by  $P$ , improves with experience  $E$ .*

# The Task T

- Example: if we want a robot to be able to walk, then **walking** is the task
- Approaches
  1. We could directly input **directives** for how we think a robot should walk, or
  2. We could provide **examples** of successful and unsuccessful walking (this is machine learning)

# The Task T

- Given an input  $x$  (e.g., a vector) produce a function  $f$ , such that  $f(x) = y$  (e.g., an integer, a probability vector)
- Examples
  - Classification
  - Regression
  - Machine translation
  - Denoising
  - Probability density estimation

# The Performance Measure $P$

- To evaluate a ML algorithm we need a way to measure how well it performs on the task
- It is measured on a separate set (**the test set**) from what we use to build the function  $f$  (**the training set**)
- Examples
  - Classification accuracy (portion of correct answers) or error rate (portion of incorrect answers)
  - Regression accuracy (e.g., least squares errors)



# The Experience $E$

- Specifies what data can be used to solve the task
- We can distinguish it based on the learning pillars
  - **Supervised**: data is composed of both the input  $x$  (e.g., features) and output  $y$  (e.g., labels/targets)
  - **Unsupervised**: data is composed of just  $x$ ; here we typically aim for  $p(x)$  or a method to sample  $p(x)$
  - **Reinforcement**: data is dynamically gathered based on previous experience

# Data

- We assume that all collected data samples in all datasets:

1. come from the same distribution  $\longrightarrow p_{x^{(i)}}(x) = p_{x^{(j)}}(x)$

2. are independent  $\longrightarrow p(x^{(1)}, \dots, x^{(m)}) = \prod_{i=1}^m p(x^{(i)})$

- This assumption is denoted **IID** (independent and identically distributed)

# Example: Linear Regression

- Given IID data inputs  $x \in \mathbb{R}^n$  and outputs  $y \in \mathbb{R}$
- **Task T**: predict  $y$  with the linear regressor  $\hat{y} = w^\top x$   
need to find the weights  $w$
- **Experience E**: training set  $X^{\text{train}} \in \mathbb{R}^{m \times n}$ ,  $Y^{\text{train}} \in \mathbb{R}^m$   
and test set  $X^{\text{test}} \in \mathbb{R}^{m \times n}$ ,  $Y^{\text{test}} \in \mathbb{R}^m$
- **Performance P**: Mean squared error

$$\text{MSE}^{\text{test}}(w) = \frac{1}{m} |X^{\text{test}} w - Y^{\text{test}}|^2$$

# Linear Regression

- Solve task T by minimizing the  $\text{MSE}^{\text{train}}$

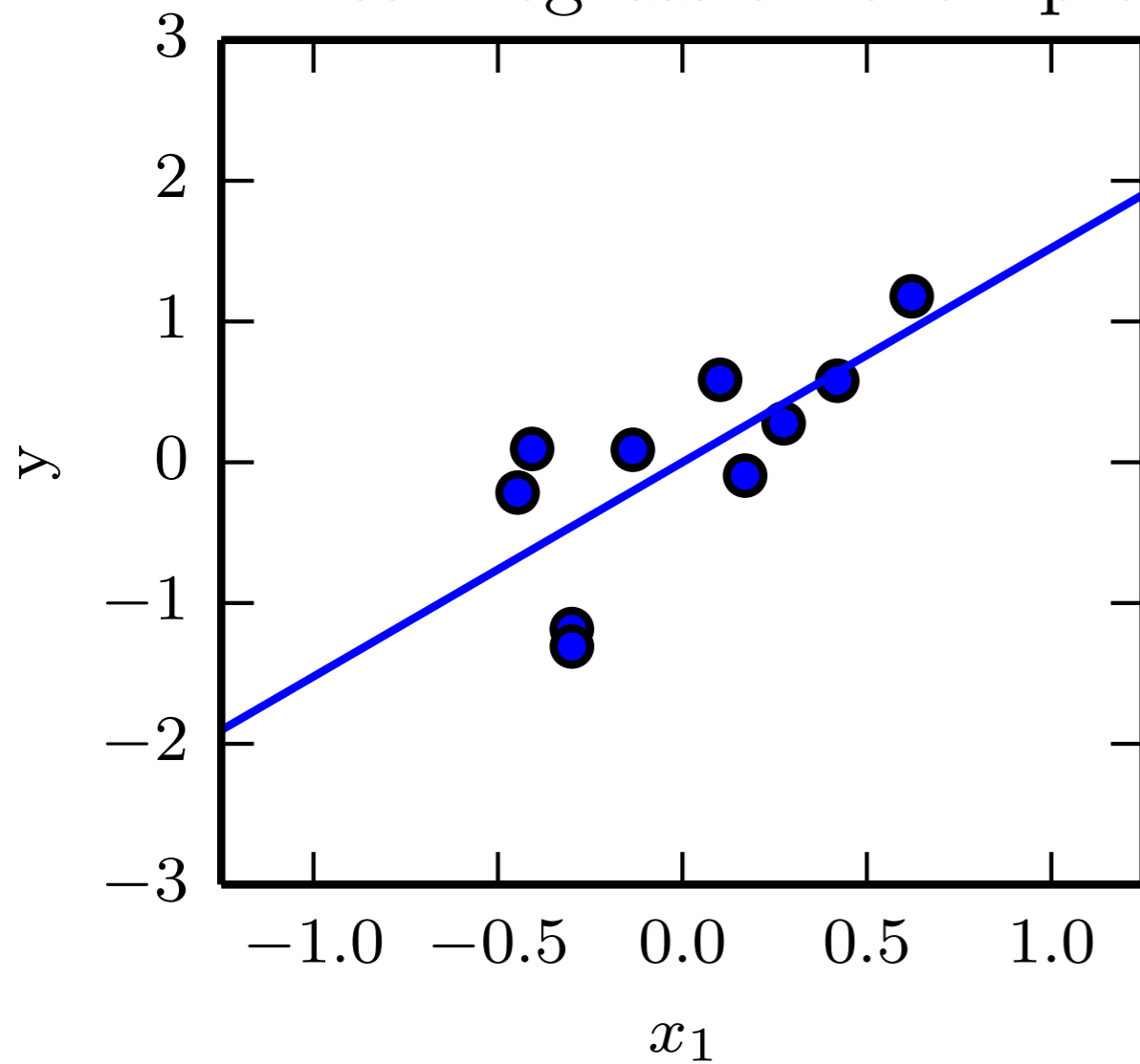
$$\text{MSE}^{\text{train}}(w) = \frac{1}{m} \|X^{\text{train}}w - Y^{\text{train}}\|^2$$

- Compute the gradient of  $\text{MSE}^{\text{train}}(w)$  with respect to  $w$  and set to 0 (normal equations)
- The solution is (pseudo-inverse)

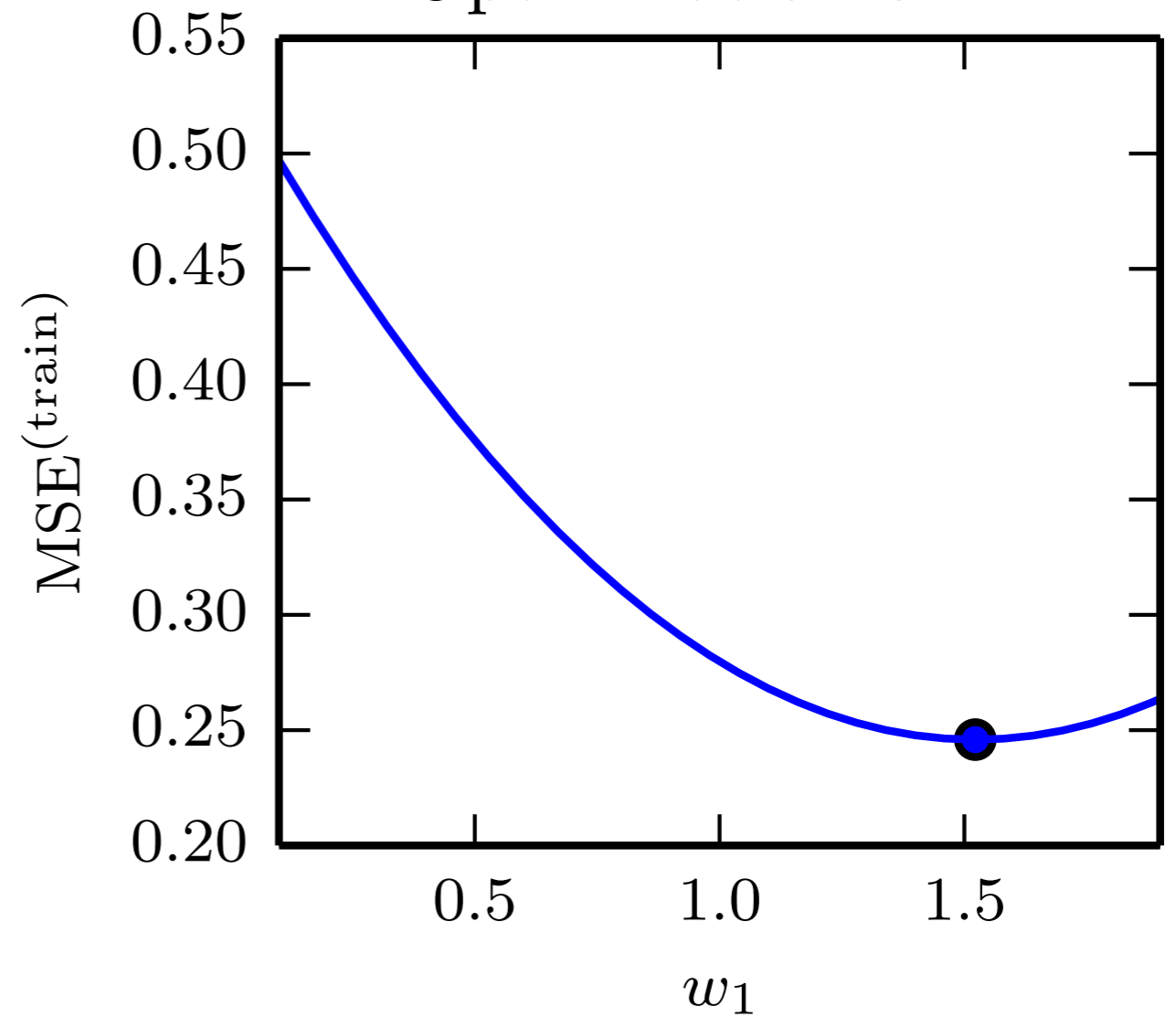
$$w = \left( X^{\text{train}\top} X^{\text{train}} \right)^{-1} X^{\text{train}\top} Y^{\text{train}}$$

# Linear Regression

Linear regression example



Optimization of  $w$



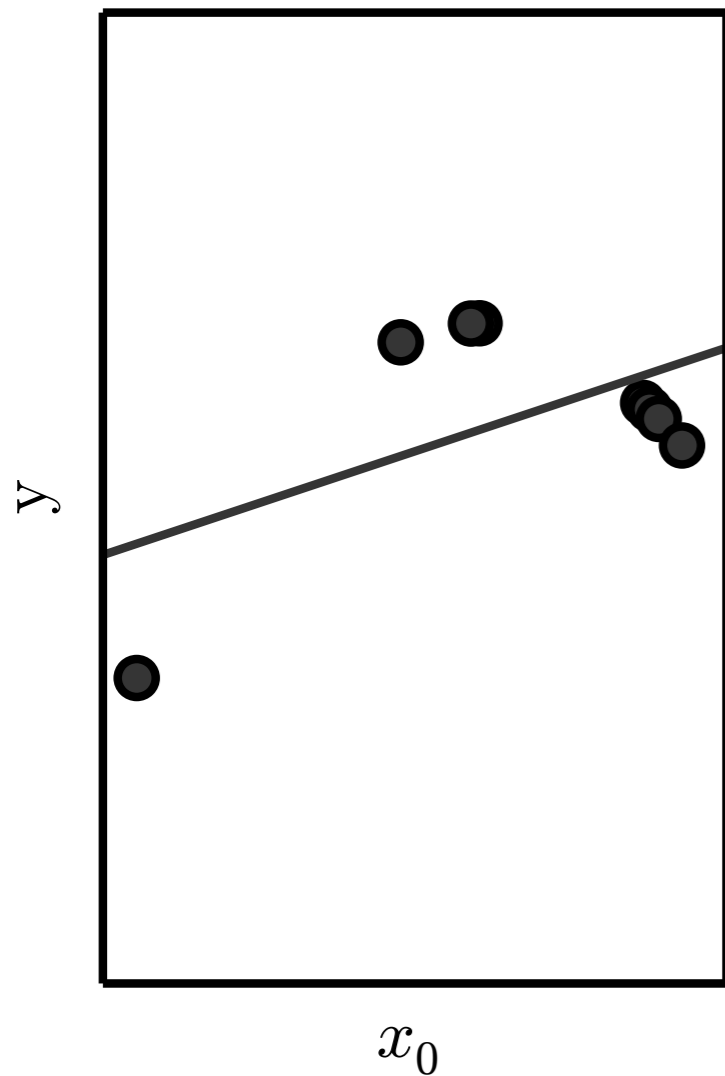
# Overfitting and Underfitting

- Performance  $P$  captures how well the learned model predicts new unseen data
- Ideally we want to select the predictor with the best performance
- What happens when we use predictors of different complexity/capacity?

# Overfitting and Underfitting

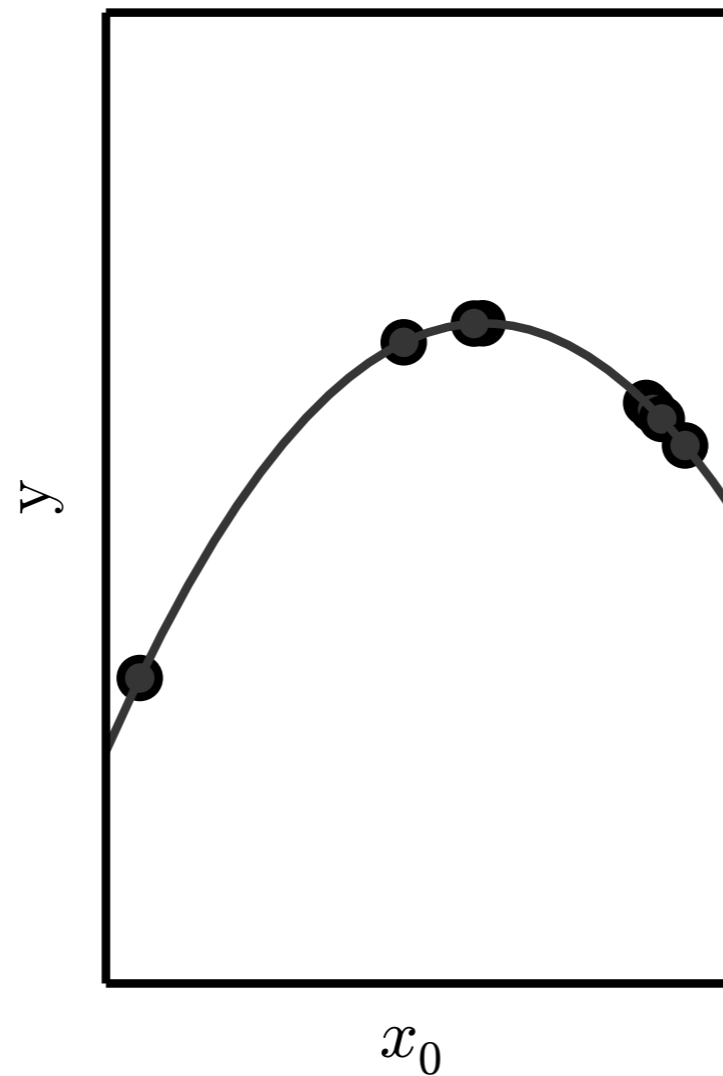
shown data is the training set

Underfitting



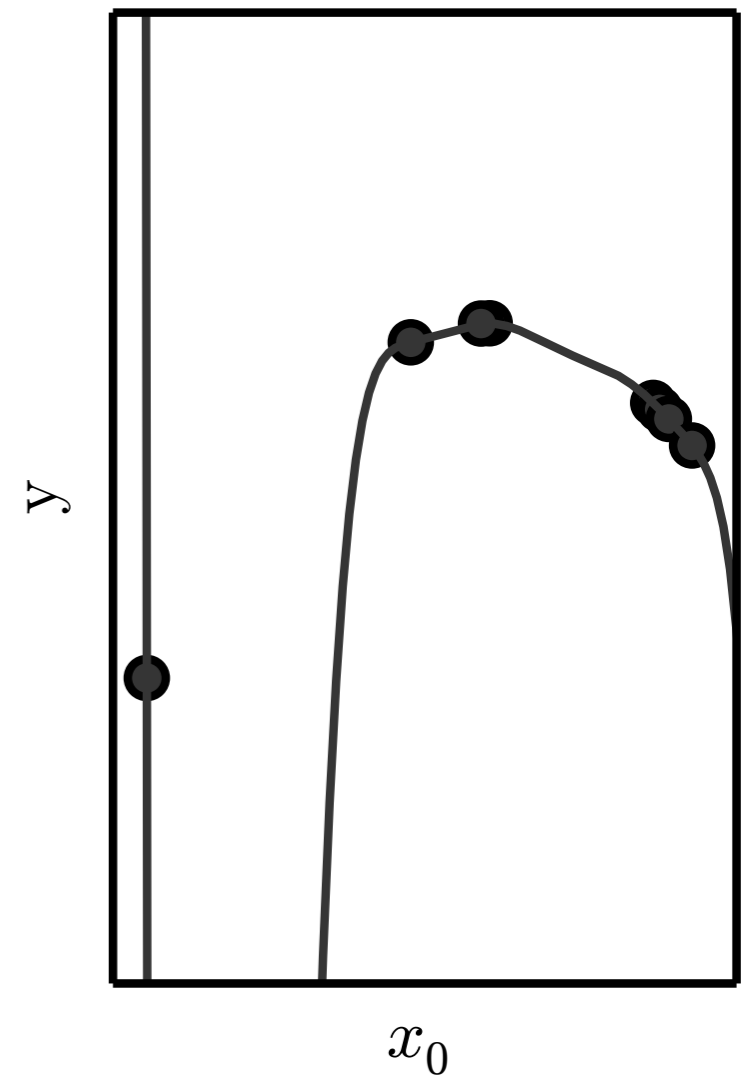
simple predictor

Appropriate capacity



optimal predictor

Overfitting



complex predictor

# Loss function

- Define a **predictor** function  $f : \mathcal{X} \mapsto \mathcal{Y}$
- Define a **loss** function  $l : \mathcal{Y} \times \mathcal{Y} \mapsto \mathbb{R}$  which measures how different the two inputs are

- Examples

- 0-1 loss 
$$l(y, f(x)) = \begin{cases} 0 & \text{if } y = f(x) \\ 1 & \text{if } y \neq f(x) \end{cases}$$

- Quadratic loss 
$$l(y, f(x)) = (y - f(x))^2$$



# Bayes Risk

- **Bayes risk** is defined as (average loss)

$$R(f) = E_{x,y}[l(f(x), y)] = \int l(f(x), y)p(x, y)dx dy$$

- The optimal predictor function is

$$f^* = \arg \min_f R(f)$$

# Empirical Risk

- Given  $(x_i, y_i)$  with  $i = 1, \dots, m$  the **empirical risk** is

$$\hat{R}(f) = \frac{1}{m} \sum_{i=1}^m l(f(x_i), y_i)$$

- The empirical predictor is

$$\hat{f} = \arg \min_{f \in \mathcal{F}} \hat{R}(f)$$

# Risks

- Bayes risk

$$R(f^*) = E_{x,y}[l(f^*(x), y)]$$

- Empirical risk

$$\hat{R}(\hat{f}) = \frac{1}{m} \sum_{i=1}^m l(\hat{f}(x_i), y_i)$$

- Bayes risk restricted to function family

$$\min_{f \in \mathcal{F}} R(f)$$

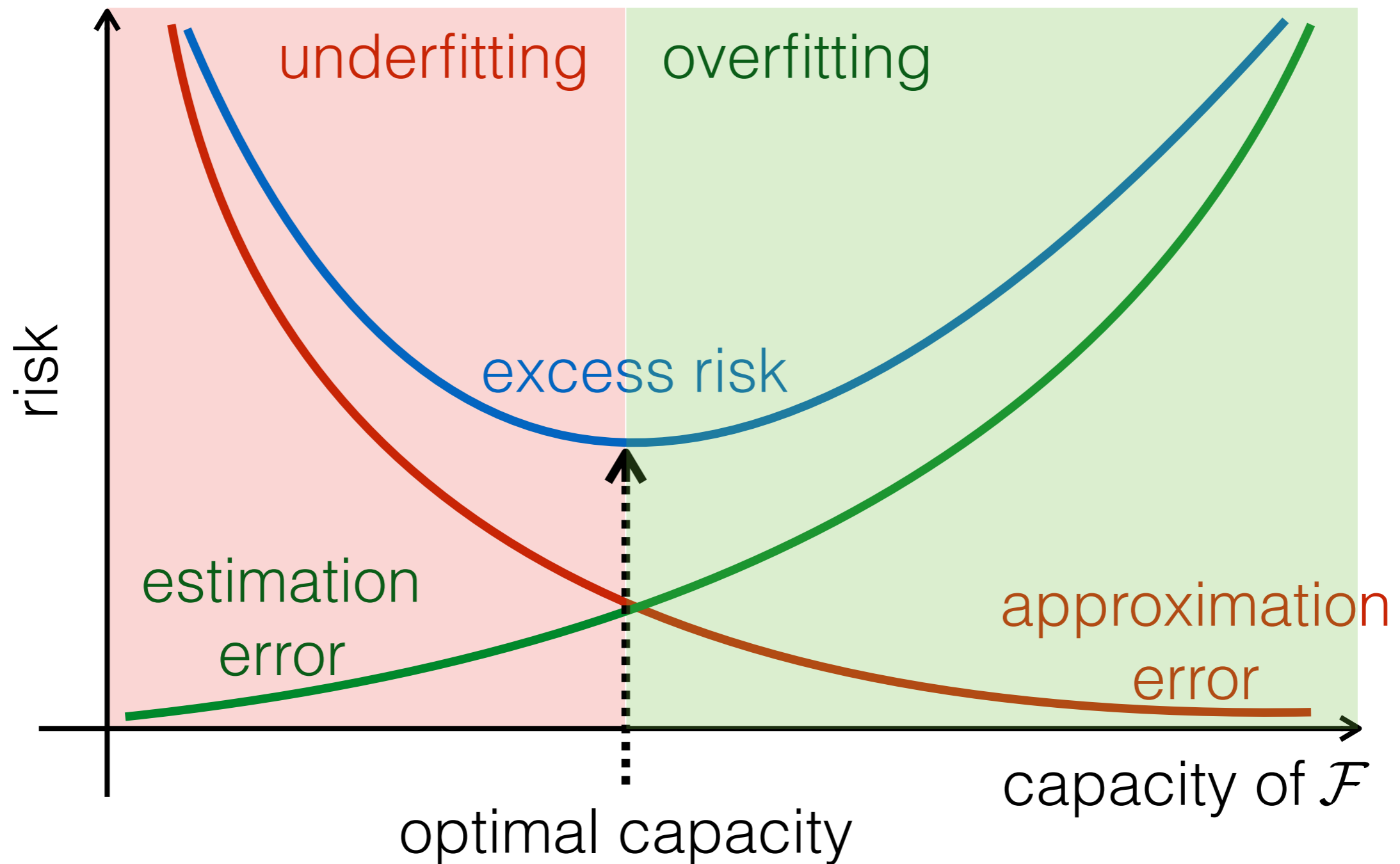
# Estimation vs Approximation

- The **excess risk** is the gap between the empirical risk and the optimal Bayes risk

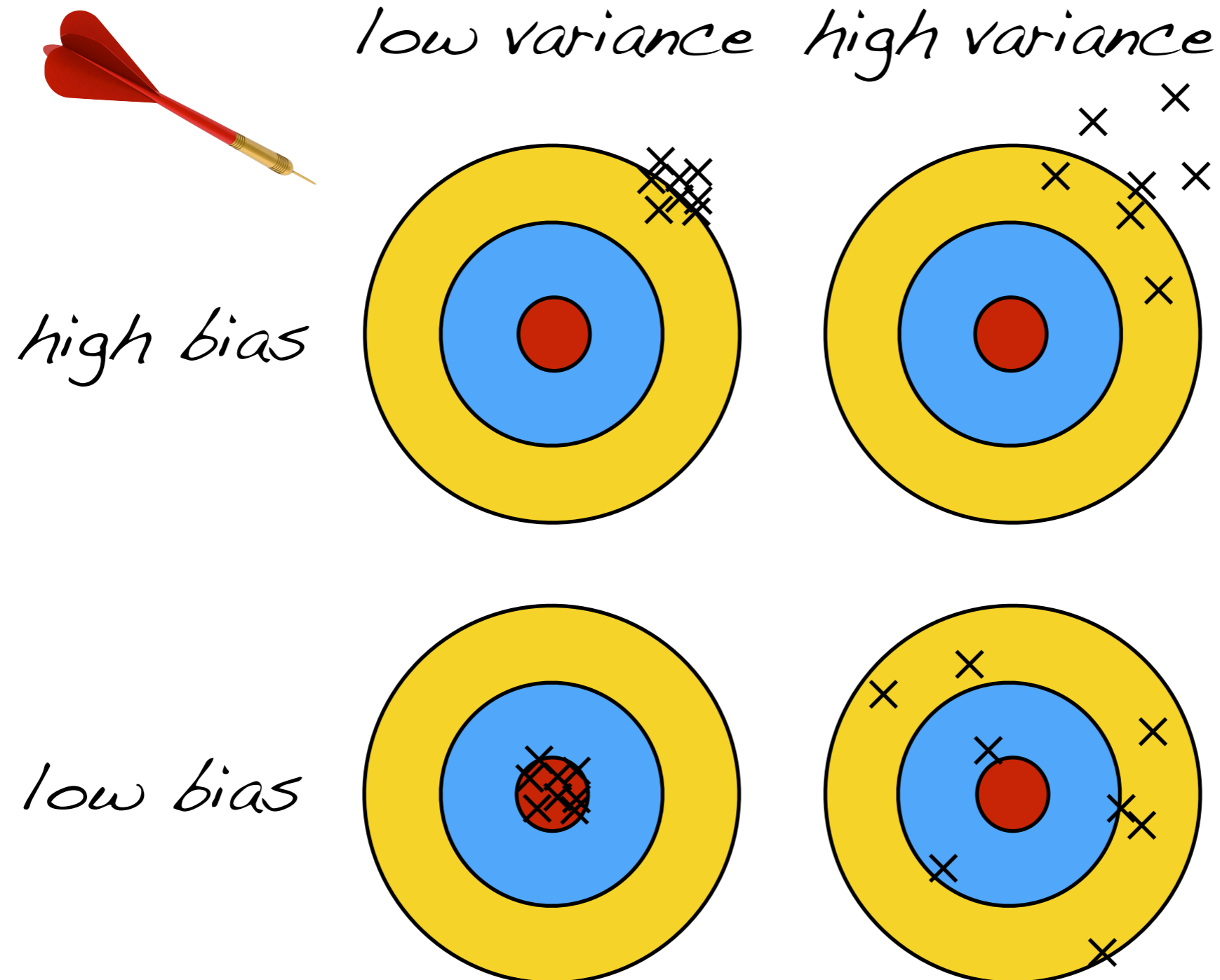
$$\hat{R}(\hat{f}) - R(f^*) = \underbrace{\hat{R}(\hat{f}) - \min_{f \in \mathcal{F}} R(f)}_{\text{estimation error}} + \underbrace{\min_{f \in \mathcal{F}} R(f) - R(f^*)}_{\text{approximation error}}$$

- Estimation (variance)**: due to training set
- Approximation (bias)**: due to function family  $\mathcal{F}$

# Estimation vs Approximation



# Bias and Variance



# Regularization

- Define a parametric family  $\mathcal{F}_\lambda$  of functions, where  $\lambda$  regulates the complexity/capacity of the predictors
- Given the optimal predictor from the empirical risk

$$\hat{f}_\lambda = \arg \min_{f \in \mathcal{F}_\lambda} \hat{R}(f)$$

we would like to choose the capacity based on Bayes risk

$$R(\hat{f}_\lambda)$$

# Regularization

- Bayes risk is not available, thus we write

$$R(\hat{f}_\lambda) = \hat{R}(\hat{f}_\lambda) + \left( R(\hat{f}_\lambda) - \hat{R}(\hat{f}_\lambda) \right)$$

and approximate the second term with a regularization term

$$C(\hat{f}_\lambda) \simeq R(\hat{f}_\lambda) - \hat{R}(\hat{f}_\lambda)$$

then solve  $\hat{\lambda} = \arg \min_{\lambda} \hat{R}(\hat{f}_\lambda) + C(\hat{f}_\lambda)$



# Training, Validation and Test

- In alternative, collect samples into training set  $D_{\text{train}}$  validation set  $D_{\text{val}}$  and test set  $D_{\text{test}}$

- Use the **training set** to define the optimal predictor

$$\hat{f}_{\lambda} = \arg \min_{f \in \mathcal{F}} \hat{R}_{D_{\text{train}}}(f)$$

- Use the **validation set** to choose the capacity

$$\hat{\lambda} = \arg \min_{\lambda} \hat{R}_{D_{\text{val}}}(\hat{f}_{\lambda})$$

- Use the **test set** to evaluate the performance

$$\text{performance } P = R_{D_{\text{test}}}(\hat{f}_{\hat{\lambda}})$$

# Supervised Learning

- Make a prediction of an output  $y$  given an input  $x$
- Boils down to determining the conditional probability

$$p(y|x)$$

- Formulate problem as that of finding  $\theta$  for a parametric family (Maximum Likelihood)

$$p(y|x; \theta)$$

# Maximum Likelihood

- Given IID input/output samples  $(x^i, y^i) \sim p_{\text{data}}(x, y)$

the **conditional maximum likelihood** estimate is

$$\begin{aligned}\theta_{\text{ML}} &= \arg \max_{\theta} \prod_{i=1}^m p_{\text{data}}(y^i | x^i; \theta) \\ &= \arg \max_{\theta} \sum_{i=1}^m \log p_{\text{data}}(y^i | x^i; \theta)\end{aligned}$$

# Supervised Learning

- **Example:** Binary classification  $y \in \{0, 1\}$
- We aim at determining  $p(y = 1|x; \theta) = \sigma(\theta^\top x)$

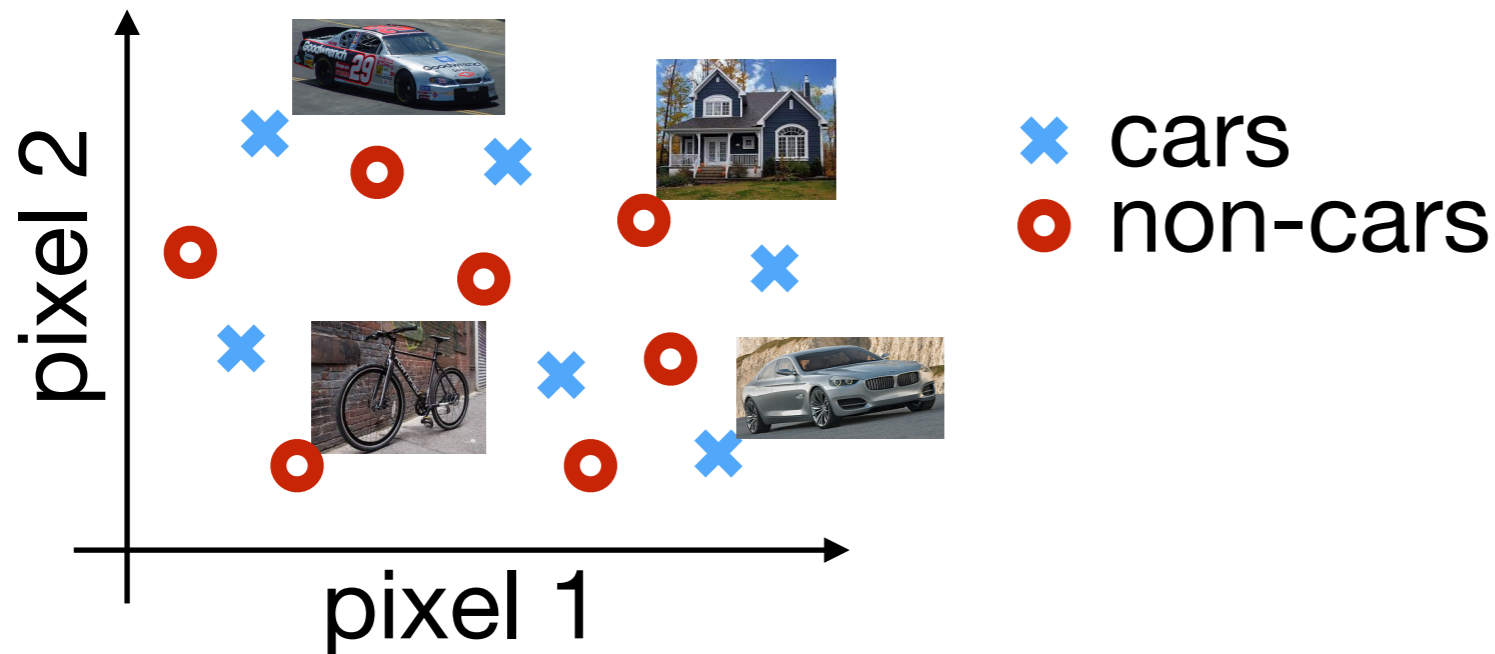
where  $\sigma(z) = \frac{1}{1 + e^{-z}}$  is the sigmoid function

- Class  $y=1$  can be picked when

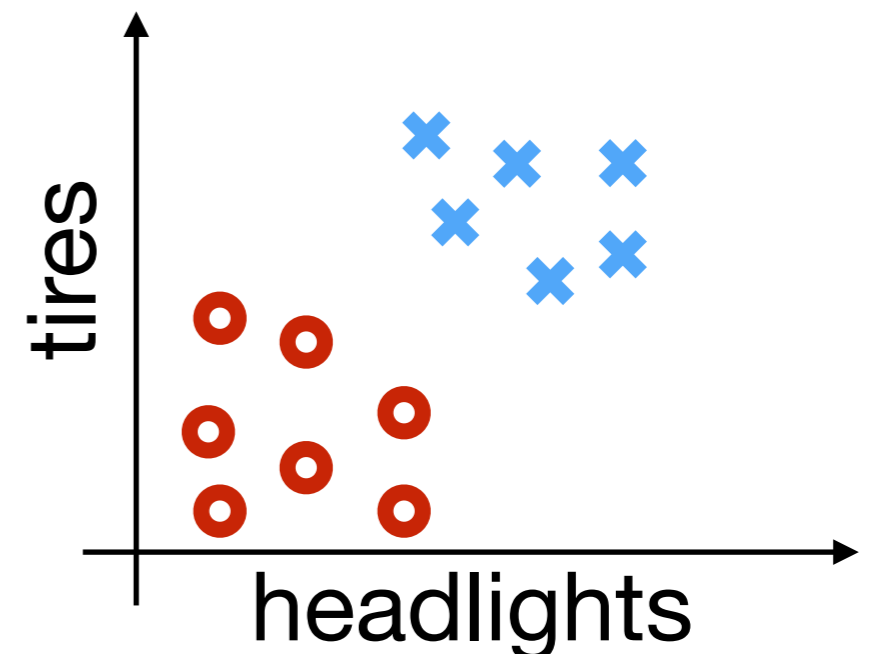
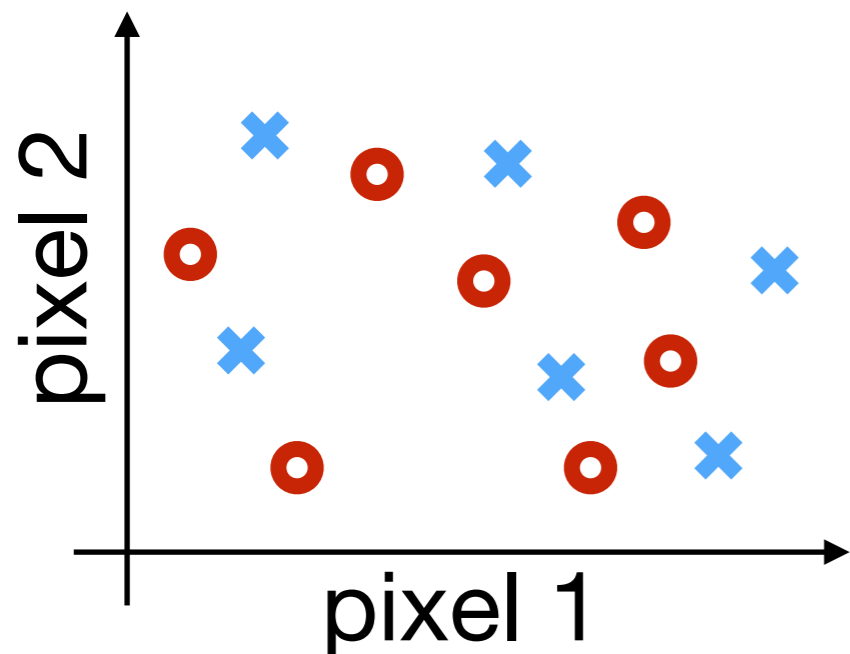
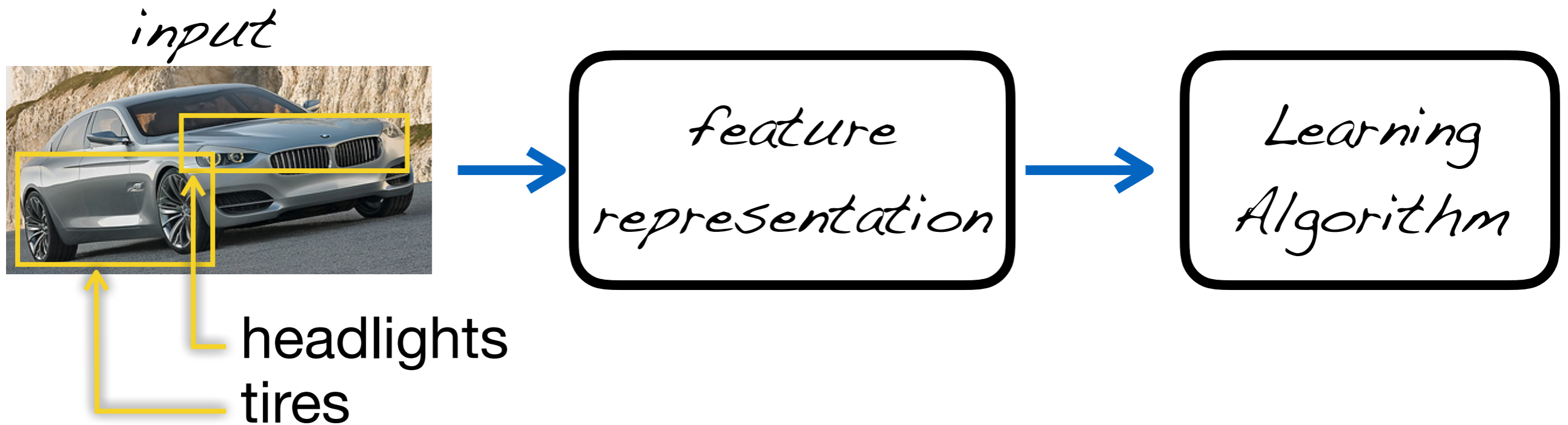
$$p(y = 1|x; \theta) > p(y = 0|x; \theta)$$

which is equivalent to  $\theta^\top x > 0$

# Features



# Features



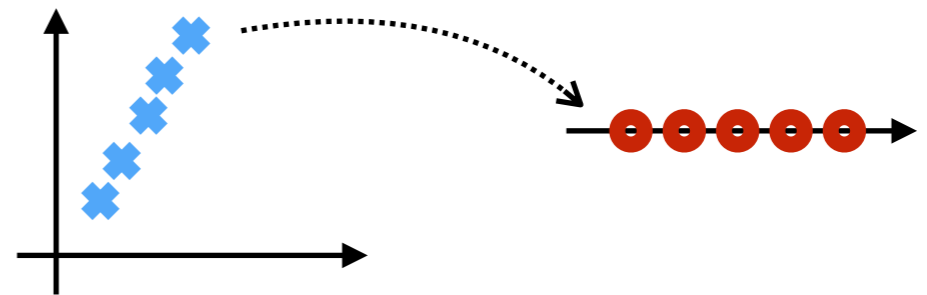
# Unsupervised Learning

- Aim is to find a suitable **data representation**
  - Probability density estimator
  - Sampling procedure
  - Data denoising
  - Manifold learning
  - Clustering

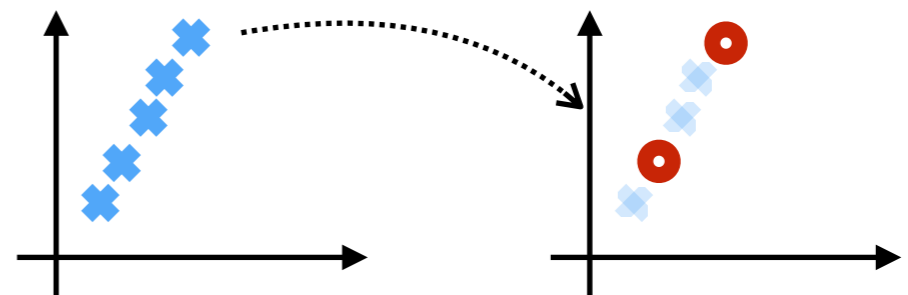
# Data Representation

- The ideal data representation should:
  1. **Preserve** all task-relevant information
  2. Be **simpler** than the original data and **easier** to use

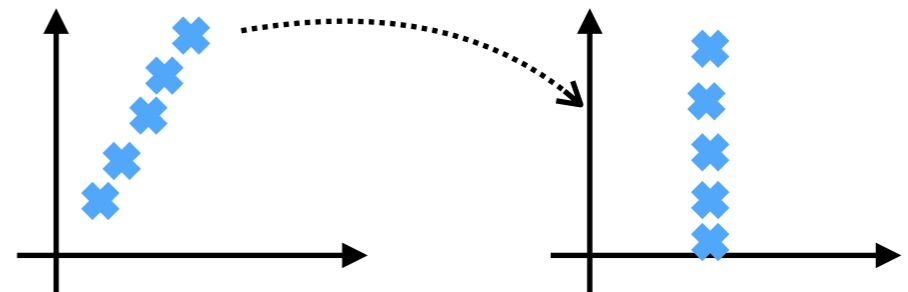
(i) low-dimensional



(ii) sparse



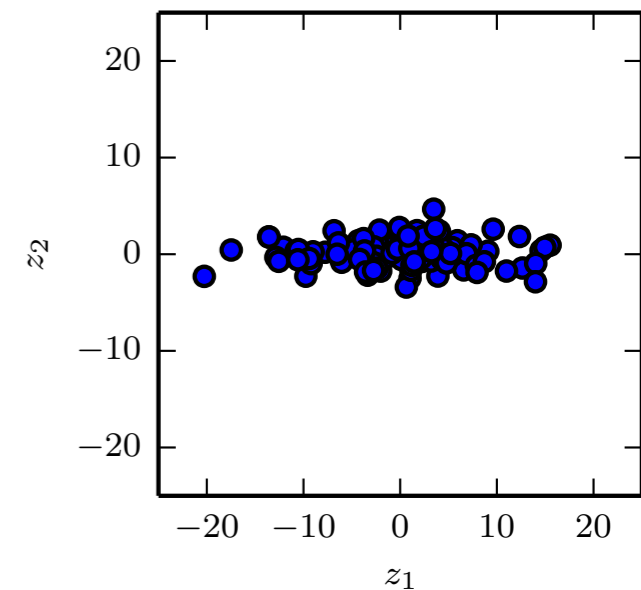
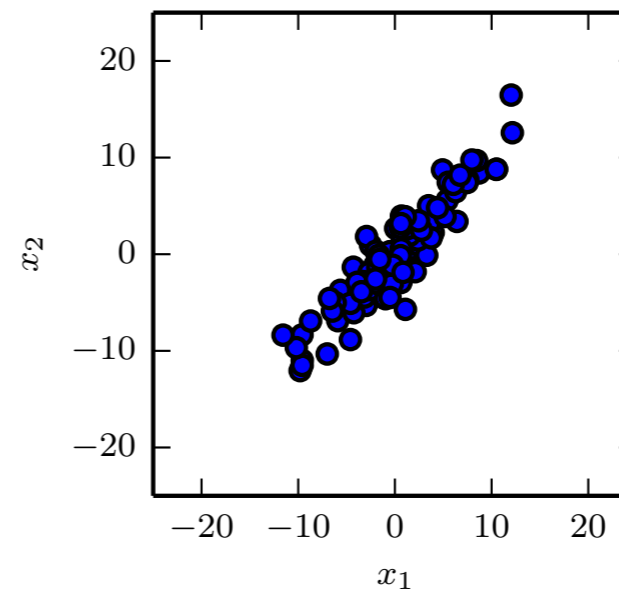
(iii) independent





# Principal Components Analysis

- **Definition:** Project data  $X$  so that the largest variation of the projected data  $Z = U^\top X$  is axis-aligned



$$X = U \Sigma V^\top$$

$$U^\top U = I$$

$$\Sigma =$$

$$\begin{bmatrix} \sigma_0 & 0 & \dots & 0 \\ 0 & \sigma_1 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma_n \end{bmatrix}$$

$$V^\top V = I$$

*singular values*  $\longrightarrow \sigma_0 \geq \sigma_1 \geq \dots \geq \sigma_n \geq 0$

# Principal Components Analysis

- Unsupervised learning method for **linearly** transformed data
- A low-dimensional representation (by thresholding the singular values)
- Yields independent (uncorrelated) components

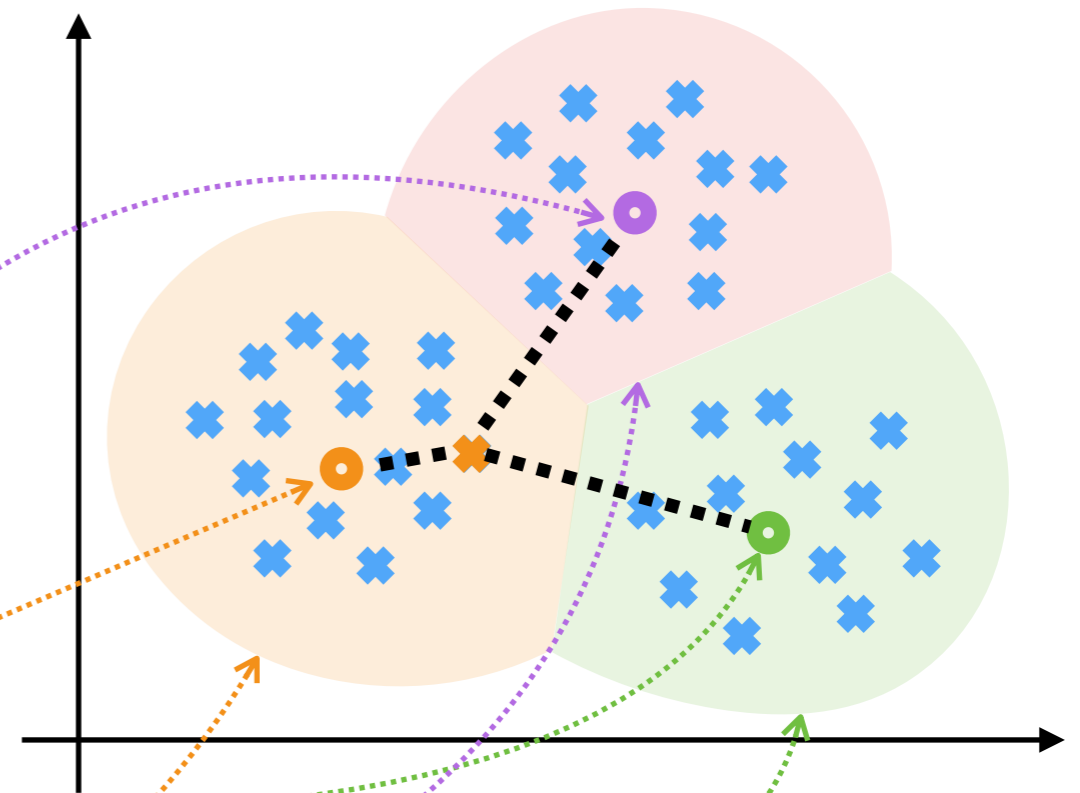
# K-Means Clustering

- **Definition:** Find  $k$  clusters of data samples similar to each other

Alternate between:

$$c_j = \frac{\sum_i \delta[w_i = j] x_i}{\sum_i \delta[w_i = j]}$$

$$w_i = \arg \min_j |x_i - c_j|^2$$



# K-Means Clustering

- Unsupervised learning method (handles nonlinearly transformed data)
- A sparse representation (assignments  $w_i$  encode one sample with one of the cluster centers  $c_j$ )
- Depends on initialization
- Ill-posed (multiple solutions can be valid)
- Number of clusters is usually unknown