

**MVE137**

# **Probability and Statistical Learning Using Python**

## **Supervised Learning**

Alexandre Graell i Amat

`alexandre.graell@chalmers.se`

<https://sites.google.com/site/agraellamat>

September 20 and 21, 2022



**CHALMERS**

# Why statistical learning?

A set of tools for **making sense** of **complex datasets**, i.e., **learning from** and **understanding data**.

- Aims at finding a **predictive function** (**model**) that relates an output to inputs and can be used for prediction.

## Applications:

- computer vision
- speech recognition
- bioinformatics, ...

**Statistical learning** & **machine learning**:

**Statistical learning** provides a **formal framework** for **machine learning**.

# Supervised learning

## Statistical learning

**Goal:** Infer a **predictive function** (**model**), such that it can be used to predict the output for new, yet unseen data.

## Supervised learning

**Goal:** Given a training set  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ , identify an algorithm to predict the outcome  $y$  for a new (yet unseen) data point  $\mathbf{x}$ .

- Learning a model from labeled data
  - Predicting output of new data based on the learned model
- 
- $\mathbf{x}_i$ : free variables (features, predictors, covariates, domain points)
  - $y_i$ : target variables (dependent variables, labels, responses)

# Supervised learning

Three classes of responses:

- **Continuous** (**quantitative**): take on numerical values
- **Discrete** (**qualitative**, **categorical**): a discrete set of categories
- **Order categorical**: the order is important

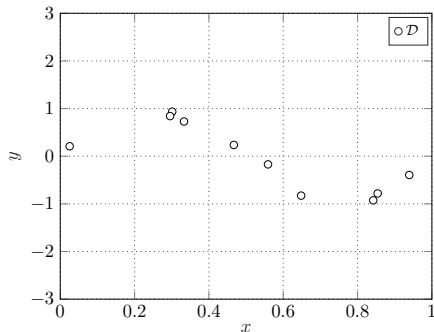
Two learning algorithms:

- **Regression**: predict a quantitative output
- **Classification**: predict a qualitative output

# Statistical inference (learning) and decision theory

**Goal:** Given a training set  $\mathcal{D} = \{(x_1, y_1), \dots, (x_N, y_N)\}$ , predict  $y$  for a new (yet unseen) data point  $x$ .

Impossible if **no information** on the mechanism relating  $x$  and  $y$ !



# Statistical learning and decision theory

**Goal:** Given a training set  $\mathcal{D} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$ , predict  $y$  for a new (yet unseen) data point  $\mathbf{x}$ .

Impossible if no information on the mechanism relating  $\mathbf{x}$  and  $y$ !

We may assume that  $\mathbf{x}$  and  $y$  are related via a function

$$y = \tilde{f}(\mathbf{x})$$

**Goal:** Find best possible approximation of  $\tilde{f}(\mathbf{x})$ ,  $f(\mathbf{x})$  (prediction model).  
Predict the outcome  $y$  for  $\mathbf{x}$  as  $\hat{y} = f(\mathbf{x})$ .

- Treat  $\mathbf{x}$  and  $y$  as random variables, and

$$(\mathbf{x}_i, y_i) \sim_{i.i.d.} p(\mathbf{x}, y), \quad i \in [N].$$

# Statistical learning and decision theory

How should we choose  $f(\mathbf{x})$ ?

- Loss function  $\ell(y, \hat{y}) = \ell(y, f(\mathbf{x}))$ : cost (loss or risk) incurred when the correct value is  $y$  while the estimate is  $\hat{y}$
- Quadratic loss function:

$$\ell(y, \hat{y}) = (y - \hat{y})^2 = (y - f(\mathbf{x}))^2$$

Expected prediction error (expected generalization loss/error):

$$\begin{aligned} L(\hat{y}) &= \mathbb{E}_{\mathbf{x}, y \sim p(\mathbf{x}, y)} [\ell(y, f(\mathbf{x}))] \\ &= \int \int \ell(y, f(\mathbf{x})) p(\mathbf{x}, y) d\mathbf{x} dy \end{aligned}$$

# Statistical learning and decision theory

How should we choose  $f(\mathbf{x})$ ?

Optimal prediction  $\hat{y}(\mathbf{x})$  obtained by minimizing generalization loss:

$$\begin{aligned} f^*(\mathbf{x}) &= \arg \min_f L(\hat{y}) \\ &= \arg \min_f \mathbb{E}_{\mathbf{x}, y \sim p(\mathbf{x}, y)} [\ell(y, f(\mathbf{x}))] \\ &= \arg \min_f \mathbb{E}_{\mathbf{x} \sim p_{\mathbf{x}}} [\mathbb{E}_{y \sim p_{y|\mathbf{x}}} [\ell(y, f(\mathbf{x}))]] \end{aligned}$$

Hence, it suffices to solve

$$f^*(\mathbf{x}) = \arg \min_f \mathbb{E}_{y \sim p_{y|\mathbf{x}}} [\ell(y, f(\mathbf{x}))]$$



# Statistical learning and decision theory

For the quadratic loss  $\ell(y, \hat{y}) = (y - f(\mathbf{x}))^2$ ,

$$\begin{aligned} f^*(\mathbf{x}) &= \arg \min_f \mathbb{E}_{y \sim p_{y|\mathbf{x}}} [\ell(y, f(\mathbf{x}))] \\ &= \arg \min_f \mathbb{E}_{y|\mathbf{x}} [(y - f(\mathbf{x}))^2] \\ &= \mathbb{E}_{y|\mathbf{x}} [y|\mathbf{x}] \end{aligned}$$

**Problem:** Since  $p(y|\mathbf{x})$  unknown, cannot find optimal prediction via

$$f^*(\mathbf{x}) = \arg \min_f \mathbb{E}_{y|\mathbf{x}} [\ell(y, f(\mathbf{x}))]$$

**Idea:** Approximate  $\mathbb{E}_{y|\mathbf{x}} [y|\mathbf{x}]$ .

# Statistical learning and decision theory

$$f^*(\mathbf{x}) = \arg \min_f \mathbb{E}_{y \sim p_{y|\mathbf{x}}} [\ell(y, f(\mathbf{x}))]$$

**Problem:** The space of all possible functions  $f(\mathbf{x})$  is enormous!

Need to restrict the search space.

- Linear regression
- $k$ -nearest neighbors

# Linear regression

**Linear regression:** Assumes linear model for  $f(\mathbf{x})$ ,

$$y = f(\mathbf{x}) = w_0 + \sum_{j=1}^d x_j w_j,$$

with  $\mathbf{x} = (x_1, x_2, \dots, x_d)^\top$ .

$w_0$ : intercept or bias

- For simplicity, we will write

$$\tilde{\mathbf{x}} = (1, x_1, \dots, x_d)^\top \quad \text{and} \quad \mathbf{w} = (w_0, \dots, w_d)^\top$$

so that

$$f(\mathbf{x}) = \sum_{j=0}^d \tilde{x}_j w_j = \tilde{\mathbf{x}}^\top \mathbf{w}$$

**Linear regression** assumes  $\tilde{f}(\mathbf{x}) \approx \mathbf{x}^\top \mathbf{w}$  or, equivalently,  $\mathbb{E}_{y|\mathbf{x}}[y|\mathbf{x}] \approx \tilde{\mathbf{x}}^\top \mathbf{w}^*$ .

## Linear regression: Optimal $\mathbf{w}$

$$\begin{aligned}\mathbf{w}^* &= \arg \min_{\mathbf{w}} L(\hat{y}) \\ &= \arg \min_{\mathbf{w}} \mathbb{E}_{\mathbf{x}, y \sim p(\mathbf{x}, y)} [\ell(y, f(\mathbf{x}))] \\ &= \arg \min_{\mathbf{w}} \mathbb{E}_{\mathbf{x}, y \sim p(\mathbf{x}, y)} [(y - \mathbf{x}^\top \mathbf{w})^2]\end{aligned}$$

**Idea:** Approximate the expectation by the **empirical average** over the  $N$  training points  $\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N$ ,

$$\begin{aligned}\mathbf{w}^* &= \arg \min_{\mathbf{w}} \frac{1}{N} \sum_{i=1}^N (y_i - \mathbf{x}_i^\top \mathbf{w})^2 \\ &= \arg \min_{\mathbf{w}} \underbrace{\sum_{i=1}^N (y_i - \mathbf{x}_i^\top \mathbf{w})^2}_{\text{RSS}}\end{aligned}$$

**RSS:** residual sum of squares

## Least squares linear regression

$$\begin{aligned}\mathbf{w}^* &= \arg \min_{\mathbf{w}} \sum_{i=1}^N (y_i - \mathbf{x}_i^\top \mathbf{w})^2 \\ &= \arg \min_{\mathbf{w}} \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 \\ &= \arg \min_{\mathbf{w}} (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w})\end{aligned}$$

with

$$\mathbf{X} = \begin{pmatrix} 1 & x_{1,1} & x_{1,2} & \dots & x_{1,d} \\ 1 & x_{2,1} & x_{2,2} & \dots & x_{2,d} \\ \vdots & \vdots & \vdots & \dots & \vdots \\ 1 & x_{N,1} & x_{N,2} & \dots & x_{N,d} \end{pmatrix}$$

and  $\mathbf{y} = (y_1, \dots, y_N)^\top$

# Least squares linear regression

$$\mathbf{w}^* = \arg \min_{\mathbf{w}} (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w})$$

1. Differentiate:

$$\begin{aligned} \frac{\partial (\mathbf{y} - \mathbf{X}\mathbf{w})^\top (\mathbf{y} - \mathbf{X}\mathbf{w})}{\partial \mathbf{w}} &= \frac{\partial \mathbf{y}^\top \mathbf{y} - 2\mathbf{w}^\top \mathbf{X}^\top \mathbf{y} + \mathbf{w}^\top \mathbf{X}^\top \mathbf{X} \mathbf{w}}{\partial \mathbf{w}} \\ &= -2\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) \end{aligned}$$

2. Equate to zero:

$$\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{w}) = 0$$

Optimal  $\mathbf{w}^*$ :

$$\mathbf{w}^* = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

# Least squares linear regression for classification

Binary classification:  $y \in \{a, b\}$

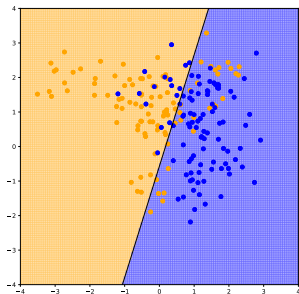
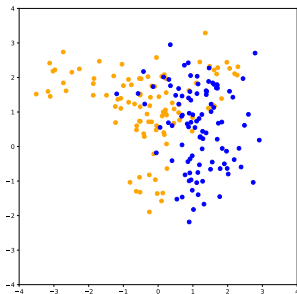
**Idea:** Encode  $y = a \implies y = 0$  and  $y = b \implies y = 1$  and apply plain linear regression plus **thresholding**:

$$\hat{y}(\mathbf{x}) = \begin{cases} 0 & \text{if } f^*(\mathbf{x}) = \mathbb{E}_{y|\mathbf{x}}[y|\mathbf{x}] = \mathbf{x}^\top \mathbf{w}^* \leq 0.5 \\ 1 & \text{otherwise} \end{cases}$$

**Observations:**

- The linear model determines the decision boundary  $\{\mathbf{x} : \mathbf{x}^\top \mathbf{w}^* = 0.5\}$
- We can interpret  $f^*(\mathbf{x}) = \mathbb{E}_{y|\mathbf{x}}[y|\mathbf{x}] = \mathbf{x}^\top \mathbf{w}$  as the probability  $p(y = 1|\mathbf{x})$

# Least squares linear regression for classification



- 100 samples from two bivariate Gaussian distributions
- Blue:  $\mathcal{N}((1, 0)^T, \mathbf{I})$
- Orange:  $\mathcal{N}((0, 1)^T, \mathbf{I})$

**Goal:** For a new coordinate  $\mathbf{x} = (x_1, x_2)$ , determine to which of the Gaussians corresponds to.



## $k$ -Nearest neighbor regression

$k$ -Nearest neighbor regression: Given  $\mathbf{x}$ , it predicts  $y$  as

$$\hat{y}(\mathbf{x}) = \frac{1}{k} \sum_{\mathbf{x}_i \in \mathcal{N}_k(\mathbf{x})} y_i(\mathbf{x}_i)$$

$\mathcal{N}_k(\mathbf{x})$ : set of  $k$  nearest neighbors to  $\mathbf{x}$  in the training set.

- For  $\mathbf{x} \in \mathbb{R}^d$ , we consider the Euclidean distance  $d(\mathbf{x}, \mathbf{x}_i) = \|\mathbf{x} - \mathbf{x}_i\|^2 \longrightarrow \mathcal{N}_k(\mathbf{x})$  is the set  $\{\mathbf{x}_i\}$  closest (in Euclidean distance) to  $\mathbf{x}$ .

# $k$ -Nearest neighbor regression

$k$ -Nearest neighbor regression:

$$\hat{y}(\mathbf{x}) = \frac{1}{k} \sum_{\mathbf{x}_i \in \mathcal{N}_k(\mathbf{x})} y_i(\mathbf{x}_i)$$

Optimal solution:

$$\hat{y}(x) = f^*(\mathbf{x}) = \mathbb{E}_{y|\mathbf{x}}[y|\mathbf{x}]$$

$k$ -nearest neighbors tries to accomplish this directly!

1. Replacing  $\mathbf{x} = x$  by neighborhood of  $x$  in the training data  $\mathcal{N}_k(x)$
2. Replacing expectation by average over the  $k$  training neighbors

We make **two approximations**:

- Expectation approximated by sample average
- Conditioning at point  $x$  relaxed to conditioning on region close  $x$

## $k$ -Nearest neighbor regression

What's the **role** of  $k$ ?

For a larger  $k$  ...

- Average is more accurate and stable (reduced variance)
- Neighborhood is bigger and less representative of  $\mathbf{x} = \mathbf{x}$  (increased bias)

Under mild regularity conditions on  $p(\mathbf{x}, y)$ , as  $N \rightarrow \infty$ ,  $k(N) \rightarrow \infty$  and  $k(N)/N \rightarrow 0$ ,  $\hat{y}(\mathbf{x}) \rightarrow \tilde{f}(\mathbf{x})$ ,  $\forall \mathbf{x}$ .

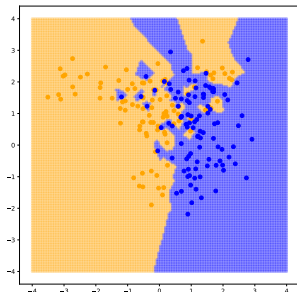
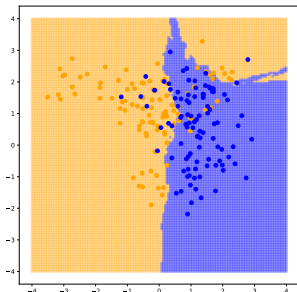
## $k$ -Nearest neighbors for classification

**Idea:** Replace averaging by a **majority vote**,

$$\hat{y}(\mathbf{x}) = \mathbb{1} \left\{ \frac{1}{k} \sum_{\mathbf{x}_i \in \mathcal{N}_k(\mathbf{x})} \mathbb{1}\{y_i(\mathbf{x}_i) = 1\} > 0.5 \right\}$$

1. Find  $\mathcal{N}_k(\mathbf{x})$
2. **Majority vote:** Assign  $\mathbf{x}$  to the class that most predictors in  $\mathcal{N}_k(\mathbf{x})$  belong

## $k$ -Nearest neighbors for classification



- Left: 15-nearest neighbors
- Right: 1-nearest neighbor

Error on training data **decreases** with **decreasing**  $k$ , and is zero for  $k = 1$ .

Is  $k = 1$  **optimal**?

# Least squares regression vs $k$ -nearest neighbors regression

## Least squares regression:

- Assumes  $\tilde{f}(x)$  well approximated by a globally linear function
- Very smooth boundary
- Linear decision boundary (strong assumption!)
- Low variance and (potentially) high bias

## $k$ -nearest neighbors regression:

- No stringent assumptions about underlying data
- Can adapt to any shape of the data
- Not smooth boundary (for small  $k$ )
- High variance and low bias

# Parametric vs nonparametric models

**Parametric models:** Build  $f(x)$  as a parametric model that applies to the whole space.

1. Select parametric model (hypothesis class), with fixed number of parameters
2. Learn parameters to fit the training data  $\mathcal{D}$

**Nonparametric models:** Don't make explicit assumptions about form  $f(x)$ , but describe it in terms of local behavior of the training data in the region near  $x$ .

1. Seek an estimate of  $f$  that gets as close to the data points as possible without being too wiggly
2. Advantage: accurately fit a wider range of possible shapes for  $f$
3. Disadvantage: number parameters grows with amount of training data

# Classification and statistical learning

- $y \in \mathcal{C} = \{C_1, \dots, C_K\}$ ,  $|\mathcal{C}| = K$
- We assume that  $\mathbf{x}$  and  $y$  are related via function  $\tilde{f}(\mathbf{x})$  (rule)

**Goal:** Learn a rule  $f(\mathbf{x})$  which maps  $\mathbf{x}$  to one of the classes  $\{C_1, \dots, C_K\}$ .

How should we choose  $f(\mathbf{x})$ ?

- 0–1 loss function:

$$\ell(y, f(\mathbf{x})) = \begin{cases} 0 & y = f(\mathbf{x}) \\ 1 & y \neq f(\mathbf{x}) \end{cases}$$

Equivalently,

$$\ell(y, f(\mathbf{x})) = \mathbb{1}\{y \neq f(\mathbf{x})\}$$



# Classification and statistical learning

Optimal rule  $f^*(\mathbf{x})$ :

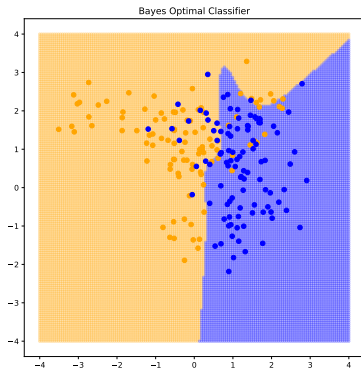
$$\begin{aligned} f^*(\mathbf{x}) &= \arg \min_f \mathbb{E}_{\mathbf{x}, y \sim p(\mathbf{x}, y)} [\ell(y, f(\mathbf{x}))] \\ &= \arg \min_f \mathbb{E}_{\mathbf{x} \sim p_{\mathbf{x}}} [\mathbb{E}_{y \sim p_{y|\mathbf{x}}} [\ell(y, f(\mathbf{x}))]] \\ &= \arg \min_f \mathbb{E}_{\mathbf{x} \sim p_{\mathbf{x}}} [\mathbb{E}_{y \sim p_{y|\mathbf{x}}} [\mathbb{1}\{y \neq f(\mathbf{x})\}]] \\ &= \arg \min_f \mathbb{E}_{y \sim p_{y|\mathbf{x}}} [\mathbb{1}\{y \neq f(\mathbf{x})\}] \end{aligned}$$

# Classification and statistical learning

$$\begin{aligned} f^*(\mathbf{x}) &= \arg \min_f \mathbb{E}_{y \sim p_{y|\mathbf{x}}} [\mathbb{1}\{y \neq f(\mathbf{x})\}] \\ &= \arg \min_f \sum_{i=1}^K \mathbb{1}\{i \neq f(\mathbf{x})\} p(y = i|\mathbf{x}) \\ &= \arg \min_f \sum_{i \neq f(\mathbf{x})} p(y = i|\mathbf{x}) \\ &= \arg \min_f \left[ 1 - p(y = f(\mathbf{x})|\mathbf{x}) \right] \\ &= \operatorname{argmax}_f p(y = f(\mathbf{x})|\mathbf{x}) \\ &= \operatorname{argmax}_{j \in [K]} p(y = j|\mathbf{x}) \end{aligned}$$

Optimal classification: Bayes' classifier

# Classification and statistical learning



# Classification and statistical learning

Optimal classification: Bayes' classifier

$$f^*(\mathbf{x}) = \operatorname{argmax}_{i \in [K]} p(y = i | \mathbf{x})$$

$k$ -nearest neighbors: Approximates the optimal solution as

$$f^*(\mathbf{x}) = \operatorname{argmax}_{j \in [K]} \frac{1}{k} \sum_{\mathbf{x}_i \in \mathcal{N}_k(\mathbf{x})} \mathbb{1}\{y_i = j\}$$

# Curse of dimensionality (regression)

Nonparametric (local) methods often perform poorly for large number of dimensions.

Optimal solution:

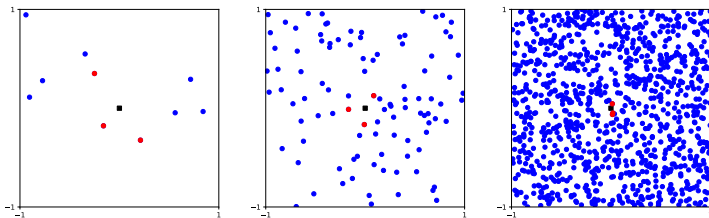
$$\hat{y}(x) = f^*(x) = \mathbb{E}_{y|x}[y|x]$$

$k$ -nearest neighbors: Approximates the optimal solution as

$$\hat{y}(x) = \frac{1}{k} \sum_{x_i \in \mathcal{N}_k(x)} y_i(x_i)$$

Can we accurately approximate the optimal solution by considering a very large training set? How large?

# Curse of dimensionality



## 3-nearest neighbors

- Blue and red: Training data points  $\mathbf{x}_i \in \mathcal{X} = [-1, 1]^2$  (10, 100, 1000)
- Black: New data point
- Red: 3 nearest neighbors

As  $N$  increases:

$$\hat{y}(\mathbf{x}) = \frac{1}{k} \sum_{\mathbf{x}_i \in \mathcal{N}_k(\mathbf{x})} y_i(\mathbf{x}_i) \xrightarrow{N \rightarrow \infty} \mathbb{E}[y|\mathbf{x}]$$

# Curse of dimensionality

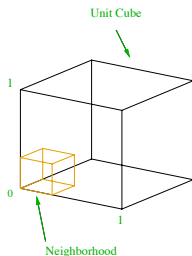
... but  $k$ -nearest neighbors (and other local methods) do not work well with **high-dimensional** inputs!

**Curse of dimensionality:** number of points exponential in number of dimensions!

1. Nearest neighbors not so close to  $x$
2.  $k$ -NNs of  $x$  closer to the boundary of  $\mathcal{X}$
3. Need a prohibitive number of training samples to densely sample  $\mathcal{X} \in \mathbb{R}^d$

(see “The elements of statistical learning,” Section 2.5, for 2 and 3)

For large  $d$ , the nearest neighbors are not so close



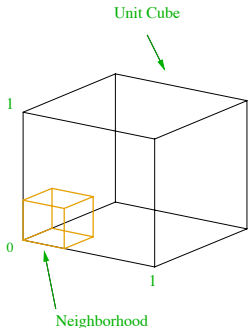
$k$ -nearest neighbors to data with training points uniformly distributed in a  $d$ -dimensional unit hypercube,  $\mathcal{X} = [0, 1]^d$ .

Want to estimate the **density** of class labels around a test point  $\mathbf{x}$  by growing a **hypercube** around  $\mathbf{x}$  until we capture a fraction  $\rho$  of the data points.

**Expected length** of the **side** of the smallest hypercube containing a fraction  $\rho$  of the data points:



For large  $d$ , the nearest neighbors are not so close



- Estimate based on 10% of the data ( $\rho = 1/10$ ):  $r_{10}(1/10) = 0.8$
- Estimate based on 1% of the data ( $\rho = 1/100$ ):  $r_{10}(1/100) = 0.63$

$k$ -nearest neighbors is **not local** in **higher dimensions**!  $\rightarrow$  Far-away data points may not be good predictors for the behavior of the function at  $x$ .

# Probabilistic models for learning

Up to now we assumed

$$y = \tilde{f}(\mathbf{x})$$

Typically assume a **probabilistic model** of the form

$$y = \tilde{f}(\mathbf{x}) + \varepsilon$$

with  $\varepsilon$  **independent** of  $\mathbf{x}$  and  $\mathbb{E}[\varepsilon] = 0$ .

Hence,

$$\mathbb{E}[y|\mathbf{x} = \mathbf{x}] = \tilde{f}(\mathbf{x})$$

How do we effectively use the training data  $\mathcal{D}$  to guide the learning of  $f(\mathbf{x})$ ?

# Probabilistic models for learning

The space of all possible regression functions  $f(x)$  is enormous!

**Idea:** Consider a **parametric form** of  $f(x)$ ,  $f_{\theta}(x)$ , with parameters  $\theta$ .

Example: **Linear regression**,

$$\begin{aligned}\hat{y} &= \sum_{j=0}^p x_j w_j + \varepsilon \\ &= \mathbf{x}^T \mathbf{w} + \varepsilon.\end{aligned}$$

Linear function of the parameters  $w_0, \dots, w_p$  and of the input variables  $x_1, \dots, x_p$ .

# Probabilistic models for learning

More in general we may consider

$$f_{\theta}(\mathbf{x}, \mathbf{w}) = \sum_{i=1}^M \theta_i \phi_i(\mathbf{x})$$

$\phi_j$ : basis functions or basis expansion

$f_{\theta}(\mathbf{x}, \mathbf{w})$ : linear basis expansion

Resulting model is much richer (nonlinear in  $\mathbf{x}$ ), but still a linear function in  $\theta$ !

# Maximum likelihood learning

How do we learn the model parameters  $\theta$  (for given  $f_{\theta}(x)$ )?

Before (least-squares regression):

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N (y_i - f_{\theta}(x_i))^2$$

**Maximum likelihood:** select  $\theta$  for which the training set  $\mathcal{D}$  has the maximum probability of being observed.

# Maximum likelihood learning

- $\mathbf{x}_{\mathcal{D}} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
- $y_{\mathcal{D}} = \{y_1, \dots, y_N\}$

Choose  $\theta$  that maximizes

$$p(y_{\mathcal{D}}|\mathbf{x}_{\mathcal{D}}, \theta) = \prod_{i=1}^N p(y_i|\mathbf{x}_i, \theta)$$

or, equivalently, the **log-likelihood (LL) function**

$$\ln p(y_{\mathcal{D}}|\mathbf{x}_{\mathcal{D}}, \theta) = \sum_{i=1}^N \ln p(y_i|\mathbf{x}_i, \theta)$$

# Maximum likelihood learning

$$\begin{aligned}\theta_{\text{ML}} &= \arg \max_{\theta} \ln p(y_{\mathcal{D}} | \mathbf{x}_{\mathcal{D}}, \theta) \\ &= \arg \min_{\theta} - \ln p(y_{\mathcal{D}} | \mathbf{x}_{\mathcal{D}}, \theta) \\ &= \arg \min_{\theta} - \sum_{i=1}^N \ln p(y_i | \mathbf{x}_i, \theta) \\ &= \arg \min_{\theta} - \frac{1}{N} \sum_{i=1}^N \ln p(y_i | \mathbf{x}_i, \theta)\end{aligned}$$

If  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ , from  $y = f_{\theta}(\mathbf{x}) + \varepsilon$ :

$$y_i | \mathbf{x}_i, \theta \sim \mathcal{N}(y_i; f_{\theta}(\mathbf{x}_i), \sigma^2)$$

$$p(y_i | \mathbf{x}_i, \theta) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - f_{\theta}(\mathbf{x}_i))^2}{2\sigma^2}}$$

## Maximum likelihood learning

$$\begin{aligned}\sum_{i=1}^N \ln p(y_i | \mathbf{x}_i, \boldsymbol{\theta}) &= \sum_{i=1}^N \ln \left( \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - f_{\boldsymbol{\theta}}(\mathbf{x}_i))^2}{2\sigma^2}} \right) \\&= -\frac{1}{2} \sum_{i=1}^N \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - f_{\boldsymbol{\theta}}(\mathbf{x}_i))^2 \\&= -\frac{N}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - f_{\boldsymbol{\theta}}(\mathbf{x}_i))^2\end{aligned}$$



# Maximum likelihood learning

$$\sum_{i=1}^N \ln p(y_i | \mathbf{x}_i, \boldsymbol{\theta}) = -\frac{N}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - f_{\boldsymbol{\theta}}(\mathbf{x}_i))^2$$

$$\boldsymbol{\theta}_{\text{ML}} = \arg \min_{\boldsymbol{\theta}} -\frac{1}{N} \sum_{i=1}^N \ln p(y_i | \mathbf{x}_i, \boldsymbol{\theta})$$

$$= \arg \min_{\boldsymbol{\theta}} \frac{1}{2} \ln(2\pi\sigma^2) + \frac{1}{2\sigma^2 N} \sum_{i=1}^N (y_i - f_{\boldsymbol{\theta}}(\mathbf{x}_i))^2$$

$$= \arg \min_{\boldsymbol{\theta}} \frac{1}{N} \sum_{i=1}^N (y_i - f_{\boldsymbol{\theta}}(\mathbf{x}_i))^2$$

**ML** learning coincides with the **least squares regression** criterion (for  $f_{\boldsymbol{\theta}}(\mathbf{x}) = \mathbf{x}^T \mathbf{w}$  and Gaussian noise)!

# The need of structured regression models

**Goal:** Choose a function  $f \in \mathcal{F}$  that minimizes a given loss function  $L(\hat{y}) = L(f; \mathcal{D})$  based on training set  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}, i \in [N]$ .

Example:

$$\theta^* = \arg \min_{\theta} \text{RSS} \triangleq \sum_{i=1}^N (y_i - f_{\theta}(\mathbf{x}_i))^2$$

**Observation:** If  $\mathcal{F}$  set of all possible functions, can make  $\text{RSS}(f(\mathbf{x})) = 0$  (any  $f(\mathbf{x})$  that passes through the training points)

... but not all will **generalize** well to new data.

Need to impose some **constraints** on  $f(\mathbf{x})$ !

# The need of structured regression models

Which **constraints** should we impose?

- Restrict to parametric functions  $f_\theta$  (**linear regression**:  $\mathcal{F}$  family of all linear functions)
- Smooth functions
- ...

Three classes of structured regression models:

- **Roughness penalty**
- **Kernel methods**
- **Basis functions and dictionary methods**

# Class 1: Roughness penalty

Assuming a measure of “niceness” (e.g., **smoothness**)  $J(f)$ ,

$$f(\mathbf{x}) = \arg \min_{f \in \mathcal{F}: L(f; \mathcal{D})=0} J(f)$$

**Example:** smoothing splines

For one-dimensional data  $x \in [0, 1]$ ,  $\mathcal{F}$  is the family of all twice-differentiable functions, and we choose  $J(f)$  as

$$J(f) = \int_0^1 (f''(x))^2 dx$$

Can **relax requirement** that  $L(f; \mathcal{D}) = 0$  via

$$f(\mathbf{x}) = \arg \min_{f \in \mathcal{F}} L(f; \mathcal{D}) + \lambda J(f)$$

**Regularization methods:** Trade-off between loss and smoothness

## Class 2: Kernel methods

Estimate regression (or classification) function in a **local neighborhood**

- Need to specify nature of local neighborhood and class of functions used for local fit

**Simplest form:** **Nadaraya-Watson** weighted average,

$$f(\mathbf{x}) = \frac{\sum_{i=1}^N K_{\lambda}(\mathbf{x}, \mathbf{x}_i) y_i}{\sum_{i=1}^N K_{\lambda}(\mathbf{x}, \mathbf{x}_i)}$$

$K_{\lambda}(\mathbf{x}, \mathbf{x}_i)$ : **Kernel function**; assigns weights to  $\mathbf{x}_i$  depending on closeness to  $\mathbf{x}$

**Example 1:** **k-nearest neighbors**

$$K_k(\mathbf{x}, \mathbf{x}_i) = \mathbb{1}\{\|\mathbf{x}_i - \mathbf{x}\| \leq \|\mathbf{x}_{(k)} - \mathbf{x}\|\}$$

$\mathbf{x}_{(k)}$ :  $k$ -th closest input in data set to  $\mathbf{x}$

**Example 2:** **Gaussian kernel**

$$K_{\lambda}(\mathbf{x}, \mathbf{x}_i) = \frac{1}{\lambda} \exp\left(-\frac{\|\mathbf{x} - \mathbf{x}_i\|^2}{2\lambda}\right)$$

## Class 2: Kernel methods

More in general,

$$\theta^* = \arg \min_{\theta} \sum_{i=1}^N K_{\lambda}(\mathbf{x}, \mathbf{x}_i) (y_i - f_{\theta}(\mathbf{x}_i))^2$$

## Class 3: Basis functions and dictionary methods

$f$  modeled as a linear expansion of basis functions:

$$f_{\theta}(\mathbf{x}, \mathbf{w}) = \sum_{i=1}^M \theta_i \phi_i(\mathbf{x})$$

# The bias–variance trade-off

Model parameters control the **capacity** to fit available data.

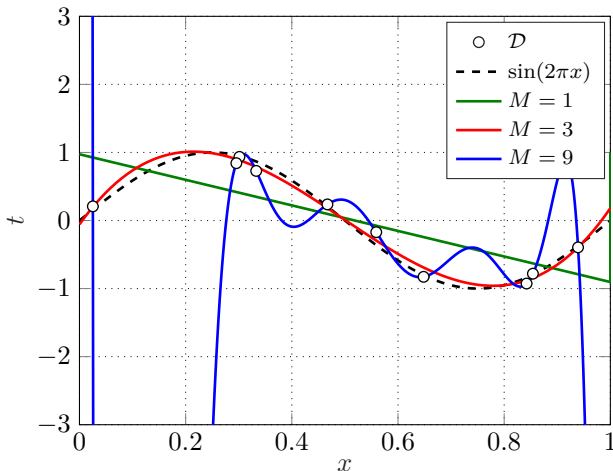
- **$k$ -NN method**:  $k$ , the number of nearest neighbors
- **Linear models**:  $M$ , the model order (the number of basis functions)
- **Roughness penalty methods**:  $\lambda$ , the weight of the penalty term

**Higher capacity** (higher complexity)  $\longrightarrow$  fit training data more accurately  
... but unlikely to generalize well

**Low capacity**  $\longrightarrow$  cannot capture all variations present on data and may generalize poorly



## Model capacity



- $p(x, y) = p(x)p(t|y)$ ,  $x \sim \mathcal{U}(0, 1)$ ,  $y|x \sim \mathcal{N}(\sin(2\pi x), 0.1)$
- Optimal predictor under quadratic loss:  $f^*(x) = \mathbb{E}_{y|x} [y|x] = \sin(2\pi x)$
- Polynomial regression

# The bias–variance trade-off

**Idea:** Decompose expected prediction error into components (bias and variance for squared error loss; approximation and estimation error for general case).

We consider:

$$y = \tilde{f}(\mathbf{x}) + \varepsilon,$$

with

- $\varepsilon \sim \mathcal{N}(0, \sigma^2) \longrightarrow \mathbb{E}[y|\mathbf{x}] = \tilde{f}(\mathbf{x})$
- $\mathbf{x} = \mathbf{x}$  fixed
- square loss function  $\ell(y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2$
- Large number of datasets drawn from  $p(y, \mathbf{x})$

Expected prediction (generalization) error:

$$\text{Err}(y, f(\mathbf{x})) = \mathbb{E}_{\mathcal{D}, \mathbf{x}, y}[\ell(y, f(\mathbf{x}))]$$

# The bias–variance trade-off

$$\text{Err}(y, f(\mathbf{x})) = \mathbb{E}_{\mathcal{D}, \mathbf{x}, y}[\ell(y, f(\mathbf{x}))]$$

For  $\mathbf{x} = \mathbf{x}$  fixed:

$$\begin{aligned}\text{Err}(\mathbf{x}) &= \mathbb{E}_{y|\mathbf{x}, \mathcal{D}}[\ell(y, f(\mathbf{x})) | \mathbf{x} = \mathbf{x}] \\ &= \mathbb{E}_{y|\mathbf{x}} \mathbb{E}_{\mathcal{D}}[\ell(y, f(\mathbf{x})) | \mathbf{x} = \mathbf{x}] \\ &= \mathbb{E}_{y|\mathbf{x}} \mathbb{E}_{\mathcal{D}}[(y - f(\mathbf{x}))^2 | \mathbf{x}] \\ &= \mathbb{E}_{y|\mathbf{x}} \mathbb{E}_{\mathcal{D}}[(y - \mathbb{E}_{y|\mathbf{x}}[y|\mathbf{x}] + \mathbb{E}_{y|\mathbf{x}}[y|\mathbf{x}] - \mathbb{E}_{\mathcal{D}}[f(\mathbf{x})] + \mathbb{E}_{\mathcal{D}}[f(\mathbf{x})] - f(\mathbf{x}))^2] \\ &= \mathbb{E}_{y|\mathbf{x}}[(y - \mathbb{E}_{y|\mathbf{x}}[y|\mathbf{x}])^2] + (\mathbb{E}_{y|\mathbf{x}}[y|\mathbf{x}] - \mathbb{E}_{\mathcal{D}}[f(\mathbf{x})])^2 + \mathbb{E}_{y|\mathbf{x}} \mathbb{E}_{\mathcal{D}}[(\mathbb{E}_{\mathcal{D}}[f(\mathbf{x})] - f(\mathbf{x}))^2] \\ &= \mathbb{E}_{y|\mathbf{x}}[(y - \tilde{f}(\mathbf{x}))^2] + (\tilde{f}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[f(\mathbf{x})])^2 + \mathbb{E}_{y|\mathbf{x}} \mathbb{E}_{\mathcal{D}}[(\mathbb{E}_{\mathcal{D}}[f(\mathbf{x})] - f(\mathbf{x}))^2] \\ &= \text{Var}[y|\mathbf{x}] + (\tilde{f}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[f(\mathbf{x})])^2 + \mathbb{E}_{\mathcal{D}}[(\mathbb{E}_{\mathcal{D}}[f(\mathbf{x})] - f(\mathbf{x}))^2]\end{aligned}$$

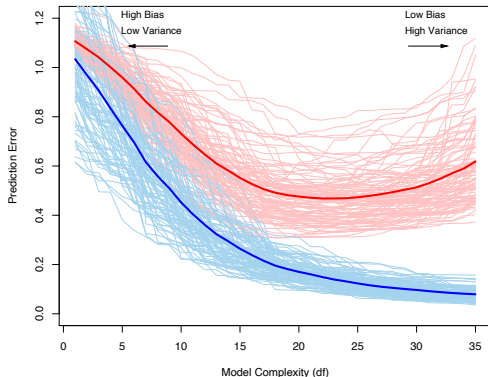
# The bias–variance trade-off

$$\begin{aligned}\text{Err}(\mathbf{x}) &= \text{Var}[y|\mathbf{x}] + (\tilde{f}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[f(\mathbf{x})])^2 + \mathbb{E}_{\mathcal{D}} [(\mathbb{E}_{\mathcal{D}}[f(\mathbf{x})] - f(\mathbf{x}))^2] \\ &= \sigma^2 + \underbrace{(\tilde{f}(\mathbf{x}) - \mathbb{E}_{\mathcal{D}}[f(\mathbf{x})])^2}_{\text{bias}} + \underbrace{\text{Var}_{\mathcal{D}} [f(\mathbf{x})]}_{\text{variance}}\end{aligned}$$

## Bias-variance decomposition

- $\sigma^2$ : Irreducible error due to data randomness
- Bias: Approximation error due to the limited flexibility model
- Variance: Estimation error; sensitivity/variability of model due to randomness in  $\mathcal{D}$

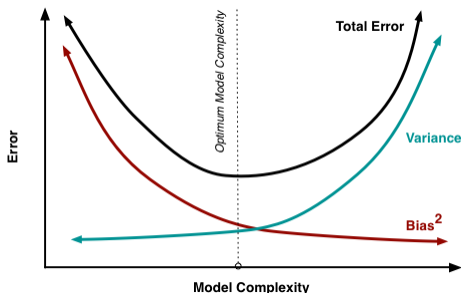
# The bias-variance trade-off



Bias-variance trade-off:

- High capacity models: low bias, high variance (overfitting)
- Low capacity models: high bias, low variance (underfitting)

# The bias–variance trade-off



Minimize prediction error → find right balance between bias and variance

**Sweet spot:** Level of complexity at which increase in bias equivalent to reduction in variance.

- If sweet spot exceeded → overfitting
- If falling short of the sweet spot → underfitting