

Machine Learning 1

- Lecture 3 -

Supervised Learning: Linear Regression

-Patrick Forré-



*Slides created by:
Rianne van den Berg*

Overview

- 1. Recap: Probability theory**
- 2. Recap: Statistical learning principles**
- 3. Linear Models for Regression**
- 4. Model selection/supervised learning**

Overview

- 1. Recap: Probability theory**
2. Recap: Statistical learning principles
3. Linear Models for Regression
4. Model selection/supervised learning

The rules of probability theory

For random variables $X \in \mathcal{X}$ and $Y \in \mathcal{Y}$:

	Discrete	Continuous
Additivity	$p(X \in A) = \sum_{x \in A} p(x)$	$p(X \in A) = \int_A p(x) dx$
Positivity	$p(x) \geq 0$	$p(x) \geq 0$
Normalization	$\sum_x p(x) = 1$	$\int_{\mathcal{X}} p(x) dx = 1$
Sum Rule	$p(x) = \sum_{y \in \mathcal{Y}} p(x, y)$	$p(x) = \int_{\mathcal{Y}} p(x, y) dy$
Product Rule	$p(x, y) = p(x y)p(y)$	$p(x, y) = p(x y)p(y)$

↳ generalizations → measure theory

Bayes Rule

$$p(y|x) = \frac{p(x|y)p(y)}{p(x)}$$

- ▶ $p(y)$: the prior probability of $Y = y$
- ▶ $p(y | x)$: the posterior probability of $Y = y$
- ▶ $p(x | y)$: the likelihood of $X = x$ given $Y = y$
- ▶ $p(x)$: the evidence for $X = x$

Multivariate Gaussian Distribution

- D -dimensional vector $\mathbf{x} = (x_1, x_2, \dots, x_D)^T$
 - $\mathcal{N}(\mathbf{x}|\mu, \Sigma) = \frac{1}{(2\pi)^{D/2} |\Sigma|^{\frac{1}{2}}} \exp\left(-\frac{1}{2} \sum (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$
 - $|\Sigma| = \det \Sigma$
 - ($D \times D$ matrix)
 - $\Sigma = \text{cov}(\mathbf{x}, \mathbf{x})$
 - $\mathbb{E}[\mathbf{x}] = \mu$
- quadratic in x
and μ

$$\int \exp\left\{-\frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x}\right\} d^D x = \frac{(2\pi)^{D/2}}{|\mathbf{A}|^{1/2}}$$

Independent Random Variables

Two random variables X and Y are *independent* iff measuring X gives no information on Y , and vice versa.

- Formally: X and Y are called independent if

$$p(x,y) = p(x) \cdot p(y) \text{ for all } x \in \mathcal{X}, y \in \mathcal{Y}$$

- Equivalent to

$$p(x|y) = p(x) \quad \forall x \quad \forall y \quad p(y) > 0$$

Overview

1. Recap: Probability theory
- 2. Recap: Statistical learning principles**
3. Linear Models for Regression
4. Model selection/supervised learning

Statistical Learning Principles

- Given: data set $D = (x_1, x_2, \dots, x_N)$ of N observations
- Choice: (parametric statistical) model class:
 $p(x_1, x_2, \dots | \mathbf{w})$ parameterised by $\mathbf{w} \in \mathcal{W}$ as “proposal distributions”
e.g. $N(x_1..x_N | \underbrace{\mu}_{\mathbf{w}}, \Sigma)$
- Questions: How do we choose a model to make predictions? Which is the best model to explain the data? etc..
- 3 general statistical learning principles to go from data to models:
 1. Maximum Likelihood Estimation (MLE)
 2. Maximum a Posteriori (MAP)
 3. Bayesian Model Averaging

(Frequentist, point estimator) Maximum Likelihood Estimation (MLE)

- Given: data set $D = (x_1, x_2, \dots, x_N)$ of N observations. $(\log p(D|w))$
- Likelihood of the dataset: $p(D|w)$ $\left| \underset{w}{\operatorname{argmax}} \, p(D|w) \right.$
(as a fct of w)
- Maximum likelihood principle: the most likely “explanation” of D is given by \mathbf{w}_{ML} which maximizes the (log-)likelihood function: $\mathbf{w}_{ML} :=$
- Under i.i.d. assumptions: each $x_i \in D$ is independently distributed according to the same underlying true distribution: $x_i \sim \underline{p_{\text{true}}(x)}$
- so may assume the *independence conditioned on w* : $x_i \sim \underline{p(x_i|w)}$ *iid* *unknown*
- For i.i.d. samples: $p(D|w) = p(x_1, x_2, \dots, x_N | w) = \prod_{i=1}^N p(x_i | w)$
- Predictions via: *new data point x*
predict : $p(x | w_{ML})$

(Bayesian point estimator) Maximum A Posteriori Estimates (MAP)

- Given: Data set $D = (x_1, x_2, \dots, x_N)$ of N observations.

- Likelihood of data:

$$p(D|w) = \underset{w}{\operatorname{argmax}} p(D|w) \cdot p(w)$$
$$= \underset{w}{\operatorname{argmin}} -\log p(D|w) - \log p(w)$$

- Prior belief in w :

$$p(w)$$

- MAP estimate: choose most probable w given the data, i.e. a posterior mode:

$$w_{\text{MAP}} = \underset{w}{\operatorname{argmax}} p(w|D)$$

- For i.i.d. data:

$$\underset{w}{\operatorname{argmin}} \left[\sum_{n=1}^N -\log p(x_n|w) - \log p(w) \right] \frac{p(D|w) \cdot p(w)}{p(D)}$$

- Predictions via:

$$P(x|w_{\text{MAP}})$$

$$P(D) = \int p(D|w) p(w) dw$$

Bayesian Model Averaging (Full Bayes)

- Given: Data set $D = (x_1, x_2, \dots, x_N)$ of N observations.
- Likelihood of data: $p(D|w)$
- Prior belief in w : $p(w)$ should represent some prior knowledge/belief of the plausibility of w .
- After observing data $D = (x_1, x_2, \dots, x_N)$, posterior distribution:
$$p(w|D) = \frac{p(D|w)}{p(D)} \cdot p(w)$$
- Predictive distribution via posterior model averaging:

$$p(x|\mathcal{D}) := \int p(x|w') p(w'|D) dw'$$

Overview

1. Recap: Probability theory
2. Recap: Statistical learning principles
3. **Linear Models for Regression**
4. Model selection/supervised learning

curve fitting
supervised learning
for continuous RV.

Linear Regression

- Regression: $D = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N)\}$

$\mathbf{x}_n \in \mathbb{R}^D$

$$\mathbf{x}_{n,d} \in \mathbb{R}^D$$
$$(\mathbf{x}_{n,d})_{d=1 \dots D}^T$$

- Input variables $\mathbf{x}_i \in \mathbb{R}^D$

- Target variables $t_i \in \mathbb{R}$

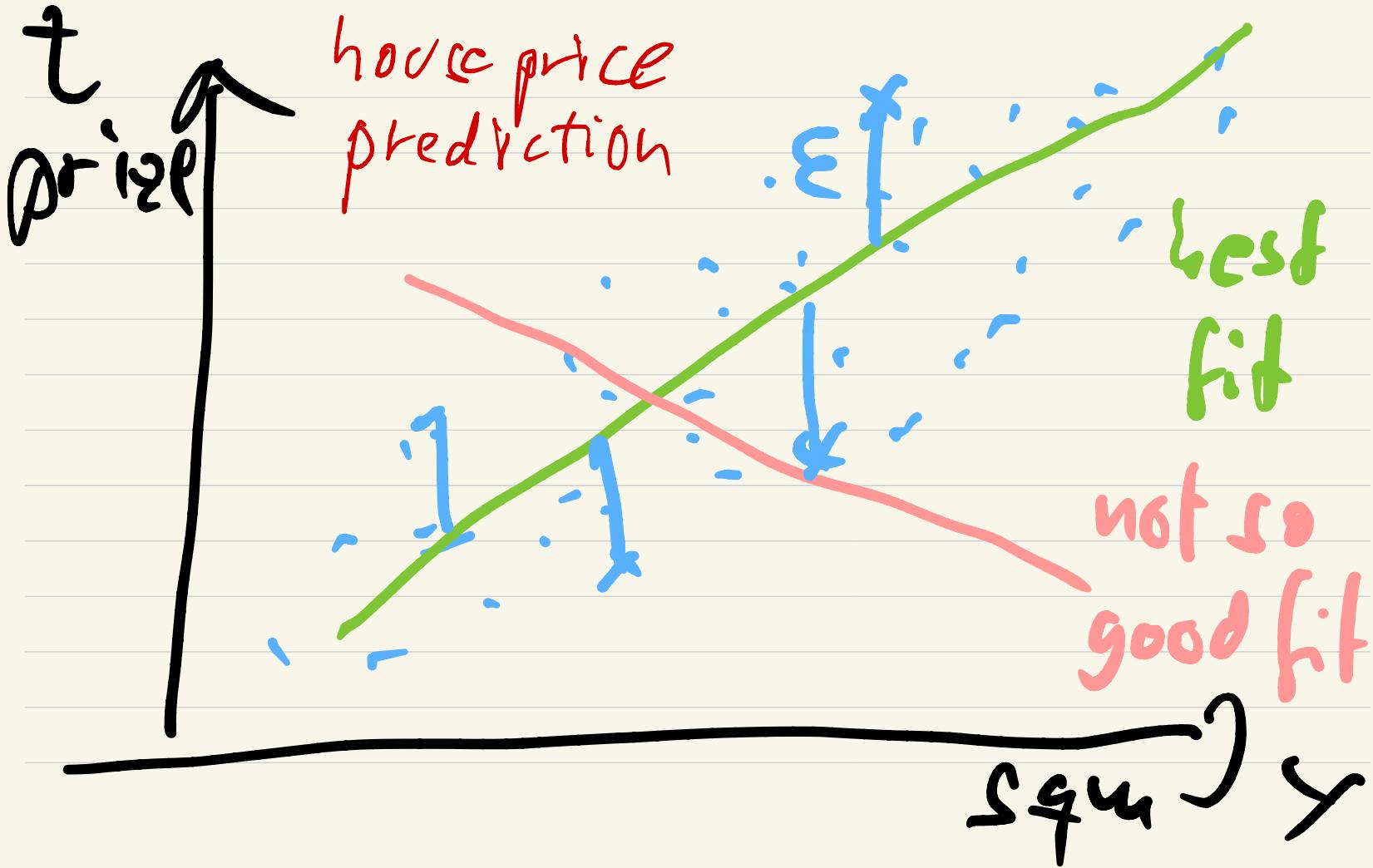
- Simplest linear model:

$$\mathbf{w} = \begin{pmatrix} w_0 \\ \vdots \\ w_{M-1} \end{pmatrix}$$

$$\mathbf{x} = \begin{pmatrix} x_0 \\ \vdots \\ x_D \end{pmatrix}$$

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + w_2 x_2 + w_3 x_3 + \dots + w_M x_M$$
$$= \mathbf{w}^T \mathbf{x}$$

$\xrightarrow{\text{components}}$



Linear Basis Models

- Fix number of parameters M s.t.

$$w \in \mathbb{R}^M$$

- Choose $M - 1$ basis functions/features of x :

- Approximation:

$$x \in \mathbb{R}^D$$

$$\phi(x) \in \mathbb{R}^M$$

$$\phi_1(x), \dots, \phi_{M-1}(x)$$

$$w = \begin{pmatrix} w_0 \\ \vdots \\ w_{M-1} \end{pmatrix} y(x, w) = w_0 + w_1 \phi_1(x) + w_2 \phi_2(x) + \dots + w_{M-1} \phi_{M-1}(x)$$

w_0 : bias/offset ↳ not linear in x
Set $\phi_0(x) = 1$ such that $\phi = [\phi_0, \dots, \phi_{M-1}]^\top$ ↳ linear in w !

$$y(x, w) = w^\top \phi(x)$$

↳ linear in w

Example: Basis Functions (I)

- Projection on input components : $\phi_i(\mathbf{x}) = x_i$

$$\mathbf{x} = (x_1, \dots, x_D)$$

for $\mathbf{x} = (x_1, x_2, \dots, x_D)^T$: $y(\mathbf{x}, \mathbf{w}) =$

$$w_0 + w_1 x_1 + w_2 x_2 + \dots$$

↳ multidim linear regression

- i -power map for $x \in \mathbb{R}^n$: $\phi_i(x) = x^i$ (monomials)

$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + w_3 x^3 + \dots$$

↳ one-dim polynomial regression

Example: Basis Functions (III)

- ▶ Gaussian basis functions: $\phi_i(\mathbf{x}) = \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_i)^T \boldsymbol{\Sigma}_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i)\right)$

$$\mathbf{x} \in \mathbb{R}^D$$

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_1)^T \boldsymbol{\Sigma}_1^{-1} (\mathbf{x}-\boldsymbol{\mu}_1)} + w_2 e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}_2^{-1} (\mathbf{x}-\boldsymbol{\mu}_2)} + \dots$$

- ▶ Logistic sigmoid functions: $\phi_i(x) = \sigma\left(\frac{x - \mu_i}{s_i}\right) \quad x \in \mathbb{R}$

with $\sigma(x) := \frac{1}{1 + \exp(-x)}$

$$y(x, \mathbf{w}) = w_0 + w_1 G\left(\frac{x - \mu_1}{s_1}\right) + w_2 G\left(\frac{x - \mu_2}{s_2}\right) + \dots$$

μ_i, s_i : fixed.

Example: Basis Functions (IV)

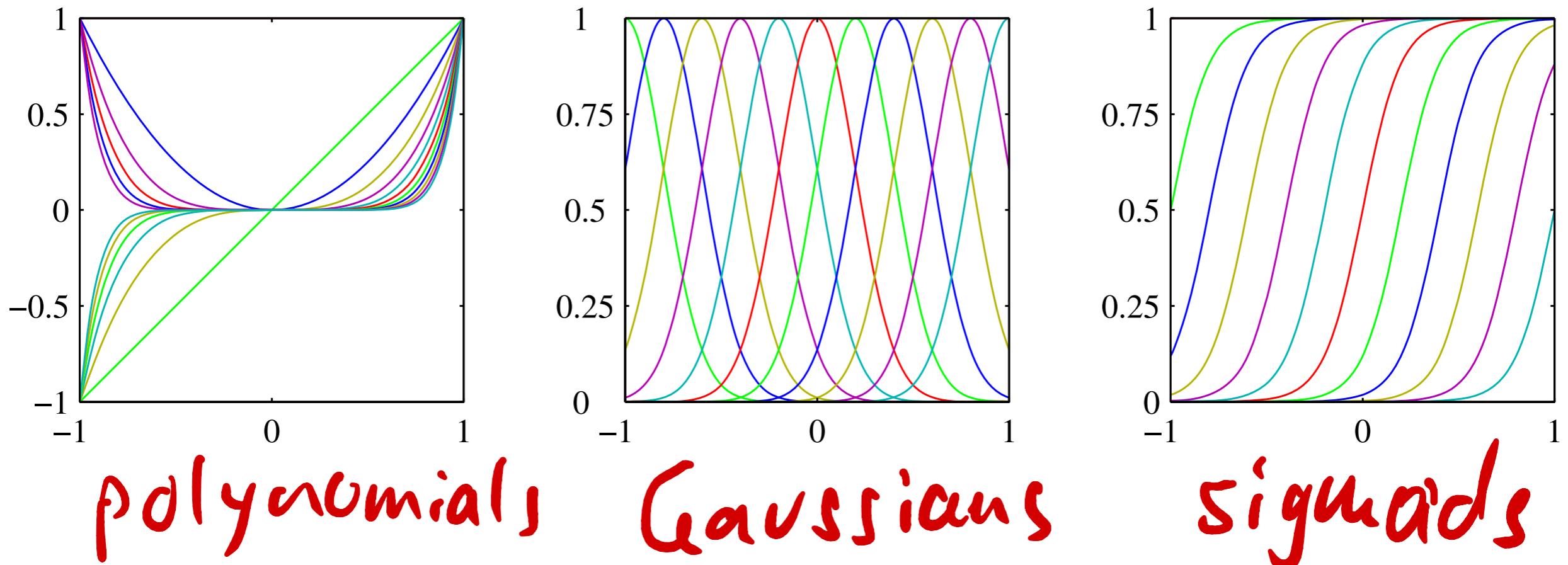


Figure: Example of basis functions (Bishop 3.1)

Maximum Likelihood

$$y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \phi(\mathbf{x})$$

Hyperparameters
 β^{-1} precision
 β^{-1} variance

- Assume gaussian noise around the target

$$t = \underbrace{\mathbf{w}^T \phi(\mathbf{x})}_{\text{determin.}} + \underbrace{\varepsilon}_{\text{random}}$$

$$\varepsilon \sim N(0, \beta^{-1})$$

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = N(t | \mathbf{w}^T \phi(\mathbf{x}), \beta^{-1})$$

$$\varepsilon \perp \!\!\! \perp \mathbf{x} \text{ (indep)}$$

- Dataset: $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ and $t = (t_1, \dots, t_N)^T$

- Likelihood function (iid)

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N N(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1})$$

ML: Sum-of-Squares Error

- › Likelihood: $p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{i=1}^N \mathcal{N}(t_i | \mathbf{w}^T \phi(\mathbf{x}_i), \beta^{-1})$
$$\frac{1}{\sqrt{2\pi}} \exp\left(-\frac{\beta}{2} (t_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2\right)$$
- › Log likelihood $\log p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) =$
$$N \log \frac{\beta}{\sqrt{2\pi}} - \frac{\beta}{2} \sum_{n=1}^N (t_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2$$
- › Sum-of-squares error: $E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (t_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2$
- › For comparison of different dataset sizes N

$$E_D^{\text{RMSE}}(\mathbf{w}) = \sqrt{\frac{1}{N} \sum (t_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2}$$

Example: Sum-of-Squares Error

(sensitive to outliers)

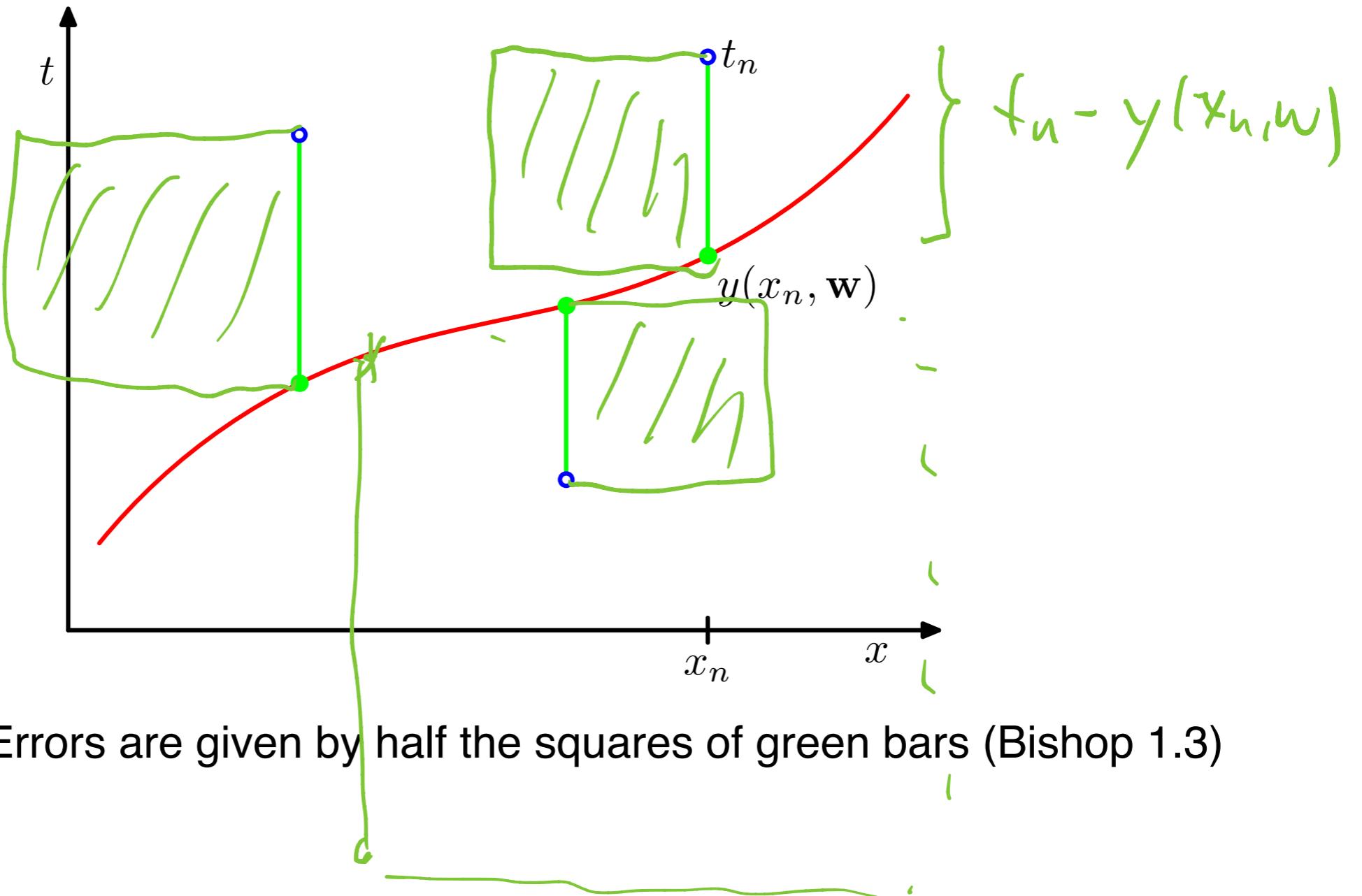


Figure: Errors are given by half the squares of green bars (Bishop 1.3)

$f: U \subseteq \mathbb{R}^d \rightarrow \mathbb{R}$ differentiable
(2nd line cont.)

Thm

1) If x is local minimum of f
 $\Rightarrow \nabla f(x) = 0$

2) If $x \in U$, $\nabla f(x) = 0$

and $\nabla^2 f(x)$ pos. def $\Rightarrow x$ local
minimum

3) If f is convex

(e.g. D^2f pos. def every where)

$\Rightarrow Df(x) = 0$ implies

x global minimum

$\Rightarrow x^* \in \arg\min_x f$

$\Rightarrow x^*$ is global min.

$\Rightarrow x^*$ is local of f

$\Rightarrow \nabla f(x^*) = 0$

so check these points
to find global min.

Maximum Likelihood Estimates

- Maximize the log likelihood / Minimize the sum-of-squares error:

$$\underset{\mathbf{w}}{\operatorname{argmax}} \log p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = \underset{\mathbf{w}}{\operatorname{argmin}} E(\mathbf{w})$$

$$\begin{aligned}\frac{\partial}{\partial \mathbf{w}} \log p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) &= -\beta \frac{\partial}{\partial \mathbf{w}} E_D(\mathbf{x}) = -\beta \frac{\partial}{\partial \mathbf{w}} \frac{1}{2} \sum_{i=1}^N \{t_i - \mathbf{w}^T \phi(\mathbf{x}_i)\}^2 \\ &= -\frac{\beta}{2} \sum_{i=1}^N \frac{\partial}{\partial \mathbf{w}} (t_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2 \\ &= +\frac{\beta}{2} \sum_{i=1}^N 2 \cdot (t_i - \mathbf{w}^T \phi(\mathbf{x}_i)) \cdot \phi(\mathbf{x}_i)^T \stackrel{?}{=} 0\end{aligned}$$

$$\stackrel{T}{\Rightarrow} \sum_{i=1}^N \phi(\mathbf{x}_i) \cdot t_i = \sum_{i=1}^N \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^T \mathbf{w}$$

$$\frac{\partial a}{\partial \mathbf{x}} = \left(\frac{\partial a}{\partial x_1}, \frac{\partial a}{\partial x_2}, \dots \right)$$


ROW vector

Maximum Likelihood Estimates

- Optimal \mathbf{w}^* satisfies

$$\sum_{i=1}^N \phi(\mathbf{x}_i) \phi(\mathbf{x}_i)^T \mathbf{w} = \sum_{i=1}^N \phi(\mathbf{x}_i) t_i$$

$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \dots & \phi_{M-1}(\mathbf{x}_1) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \dots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}$$

$N \times M$

$$\underline{\Phi}^T \underline{\Phi} w = \underline{\Phi}^T t$$

$$t = \begin{pmatrix} t_1 \\ \vdots \\ t_N \end{pmatrix}$$

if invertible $\Rightarrow w_{ML}^* = (\underline{\Phi}^T \underline{\Phi})^{-1} \underline{\Phi}^T t$

else $w^* = \underline{\Phi}^+ t$ $\underline{\Phi}^+ \leftarrow$ Moore-Penrose pseudoinverse

$$\mathbb{E}[t' | \mathbf{x}', \mathbf{w}_{ML}] = y(\mathbf{x}', \mathbf{w}_{ML}) = \mathbf{w}_{ML}^T \phi(\mathbf{x}')$$