

# Machine Learning 1

Lecture 3 - Linear Regression - Basis functions  
- Maximum Likelihood - Stochastic Gradient  
Descent - Underfitting and Overfitting -  
Regularized Least Squares

*Erik Bekkers*

*(Bishop 3.1)*



# Machine Learning 1

Lecture 2.3 - Maximum Likelihood

Erik Bekkers

(Bishop 1.2.3 - 1.2.5)



***Three Statistical Learning Principles:***

- Maximum Likelihood**
- Maximum A Posteriori**
- Bayesian Prediction**

# Maximum Likelihood Principle

- **Maximum Likelihood Optimisation:** Consider  $D = \{x_1, x_2, \dots, x_N\} \sim p(D | \mathbf{w})$ , assuming it came from some parametric distribution  $p(D | \mathbf{w})$ . Then solve:

$$\mathbf{w}_{ML} = \underset{\mathbf{w}}{\operatorname{argmax}} p(D | \mathbf{w}) \quad (\text{find parameters } \mathbf{w} \text{ that maximise the likelihood } p(D | \mathbf{w}))$$

- N(x<sub>i</sub>; μ, σ<sup>2</sup>)*
- Assume samples in  $D$  to be **independent, identically distributed** (i.i.d), meaning each  $x_i \sim p(x_i | \mathbf{w})$ :

$$\begin{aligned} p(D | \mathbf{w}) &= p(x_1, x_2, \dots, x_N | \mathbf{w}) \\ &= p(x_1 | \mathbf{w})p(x_2 | \mathbf{w}) \dots p(x_N | \mathbf{w}) \\ &= \prod_{i=1}^N p(x_i | \mathbf{w}) \end{aligned}$$

- Define equivalent problem of minimising **negative log likelihood (NLL)**:

$$\mathbf{w}_{ML} = \underset{\mathbf{w}}{\operatorname{argmin}} -\log p(D | \mathbf{w}) \quad \stackrel{i.i.d.}{\implies} \quad \mathbf{w}_{ML} = \underset{\mathbf{w}}{\operatorname{argmin}} -\sum_{i=1}^N \log p(x_i | \mathbf{w})$$

# Maximum Likelihood Principle

$\mu=2, \sigma=1$   
Likelihood:  $3.8199 \times 10^{-151}$ ,  
NLL: 369.376

$\mu=3, \sigma=1$   
Likelihood:  $2.31499 \times 10^{-94}$ ,  
NLL: 215.604

$\mu=4, \sigma=1$   
Likelihood:  $5.21914 \times 10^{-71}$ ,  
NLL: 161.831

$\mu=5, \sigma=1$   
Likelihood:  $4.37723 \times 10^{-51}$ ,  
NLL: 208.059

$\mu=6, \sigma=1$   
Likelihood:  $1.36569 \times 10^{-154}$ ,  
NLL: 354.286

# Maximum Likelihood Principle

Analytic solutions for Gaussians, i.e., the following can be solved in closed form

$$\mathbf{w}_{ML} = \underset{\mathbf{w}}{\operatorname{argmin}} - \sum_{i=1}^N \log p(x_i | \mathbf{w})$$

The solution being:

$\mu_{ML}$  is Sample mean:

$$\mu_{ML} = \frac{1}{N} \sum_{i=1}^N x_i$$

$\sigma_{ML}^2$  is sample variance:

$$\sigma_{ML}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu_{ML})^2$$

Is the estimator unbiased?

$$\mathbb{E}_{D \sim p(D|\mu, \sigma)} [\mu_{ML}] = \mu$$

$\neq$  biased!

$$\mathbb{E}_{D \sim p(D|\mu, \sigma)} [\underline{\sigma_{ML}^2}] = \frac{N-1}{N} \sigma^2$$

# Machine Learning 1

Lecture 2.4 - Maximum Likelihood: An Example

Erik Bekkers

(Bishop 1.2.3 - 1.2.5)



# Curve Fitting: Maximum Likelihood Estimates

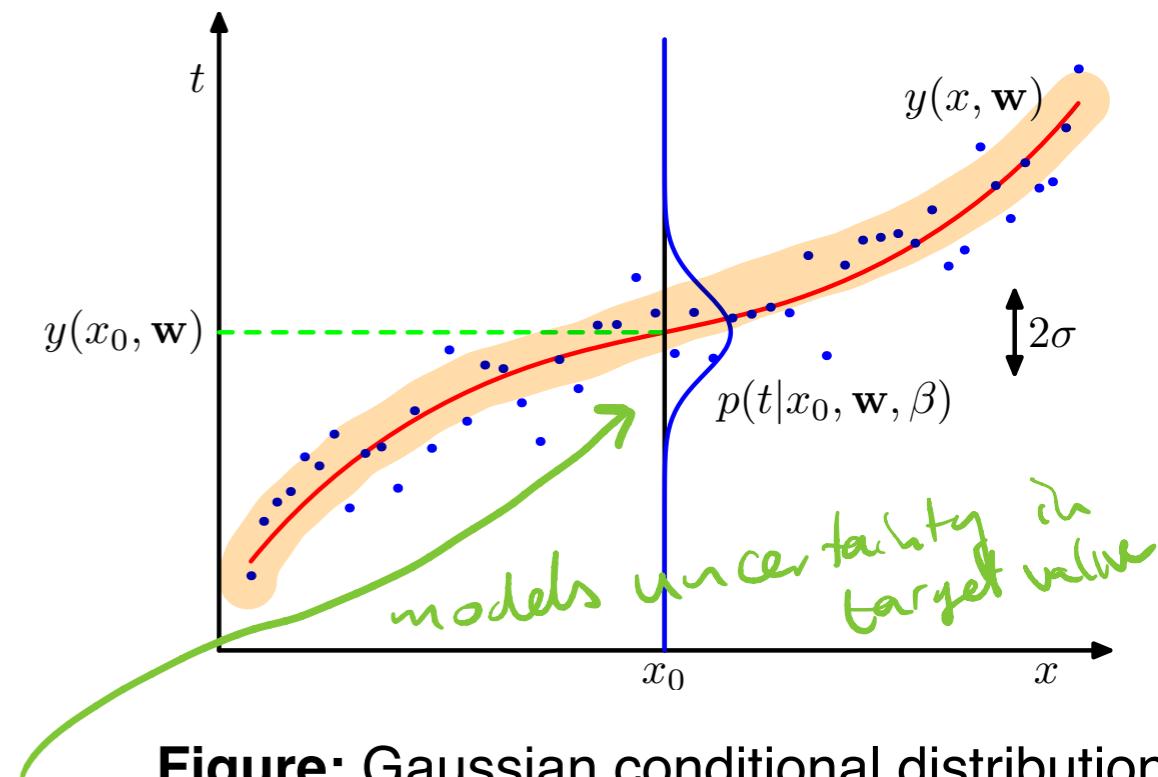
- Data  $D = \{(x_1, t_1), \dots, (x_N, t_N)\} = \{\mathbf{x}, \mathbf{t}\}$
- Assume data is generated through underlying model

$$t = y(x, \mathbf{w}) + \sigma \epsilon$$

with Gaussian noise  $\epsilon \sim \mathcal{N}(\mathbf{0}, 1)$

- Define precision  $\beta$  such that  $\beta^{-1} = \sigma^2$

- Under these model assumptions we have target distribution:



**Figure:** Gaussian conditional distribution  
(Bishop 1.16)

$$p(t | x, \mathbf{w}, \beta) = \mathcal{N}(t | y(x, \underline{\mathbf{w}}), \beta^{-1})$$

$y(x, \underline{\mathbf{w}})$  mean per position

- Log likelihood:

$$\begin{aligned} \log p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) &= \log \prod_i^N \mathcal{N}(t_i | y(x_i, \mathbf{w}), \beta^{-1}) \\ &= \frac{N}{2} \log \beta - \frac{N}{2} \log 2\pi - \frac{\beta}{2} \sum_{i=1}^N (t_i - y(x_i, \mathbf{w}))^2 \end{aligned}$$

# Curve Fitting: Maximum Likelihood Estimates

- ML: minimize  $E(\mathbf{x}, \mathbf{t}, \mathbf{w}, \beta) = -\log p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta)$  w.r.t.  $\mathbf{w}$  and  $\beta$

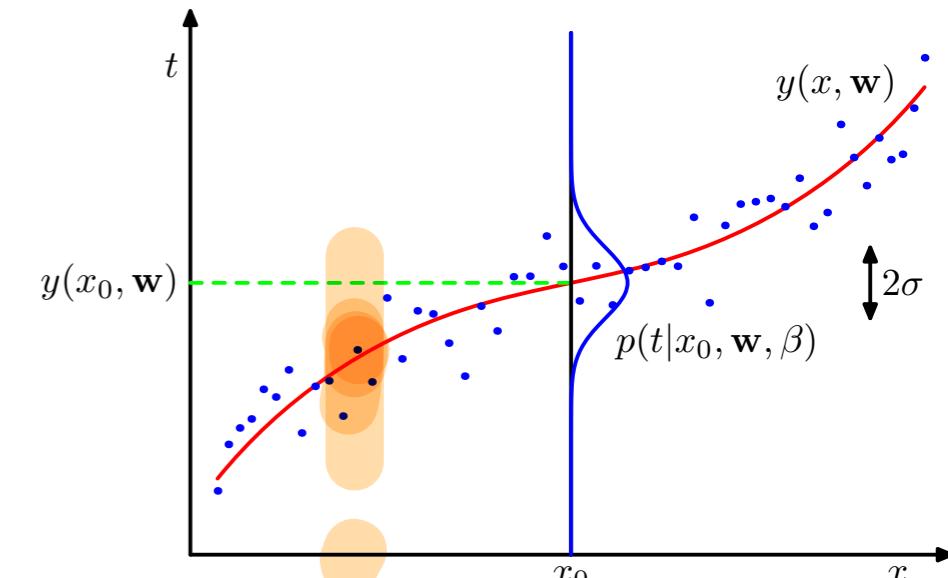
$$E(\mathbf{x}, \mathbf{t}, \mathbf{w}, \beta) = \frac{\beta}{2} \sum_{i=1}^N (y(x_i, \mathbf{w}) - t_i)^2 - \frac{N}{2} \log \beta + \frac{N}{2} \log 2\pi$$

~~does not depend on  $\mathbf{w}$~~

- Maximum likelihood solution (least squares fit!):

$$\mathbf{w}_{ML} = \operatorname{argmin}_{\mathbf{w}} \frac{\beta}{2} \sum_{i=1}^N (y(x_i, \mathbf{w}) - t_i)^2$$

$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{i=1}^N (y(x_i, \mathbf{w}_{ML}) - t_i)^2$$



**Figure:** Gaussian conditional distribution (Bishop 1.16)

- Predictive distribution:

best point estimate  $\hat{t}$

$$p(t' | x', \mathbf{w}_{ML}, \beta_{ML}) = \mathcal{N}(t' | y(x', \mathbf{w}_{ML}), \beta_{ML}^{-1})$$

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

# Machine Learning 1

Lecture 2.5 - Maximum A Posteriori

Erik Bekkers

(Bishop 1.2.5 - 1.2.6)



## **Three Statistical Learning Principles:**

Maximum Likelihood

*Data-driven*

Maximum A Posteriori

*Model-driven*

Bayesian Prediction



# The posterior

- › Dataset  $D = (x_1, x_2, \dots, x_N)$  of  $N$  independent observations.
- › Model the data distribution with some  $p(D | \mathbf{w})$  parametrized by  $\mathbf{w}$
- › MAP estimate: choose most probable  $\mathbf{w}$  given the data.

$$\mathbf{w}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmax}} p(\mathbf{w} | D)$$

- › With posterior  $p(\mathbf{w} | D)$  obtained via Bayes' rule

$$p(\mathbf{w} | D) = \frac{p(D | \underline{\mathbf{w}}) p(\underline{\mathbf{w}})}{p(D)}$$

# Maximum A Posteriori Estimates

- › Dataset  $D = (x_1, x_2, \dots, x_N)$  of  $N$  independent observations.
- › Model the data distribution with some  $p(D | \mathbf{w})$  parametrized by  $\mathbf{w}$
- › MAP estimate: choose most probable  $\mathbf{w}$  given the data.

$$\mathbf{w}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmax}} p(\mathbf{w} | D)$$



# Maximum A Posteriori Estimates

- › Dataset  $D = (x_1, x_2, \dots, x_N)$  of  $N$  independent observations.
- › Model the data distribution with some  $p(D | \mathbf{w})$  parametrized by  $\mathbf{w}$
- › MAP estimate: choose most probable  $\mathbf{w}$  given the data.

$$\mathbf{w}_{MAP} = \operatorname{argmax}_{\mathbf{w}} p(\mathbf{w} | D)$$

$$\mathbf{w}_{MAP} = \operatorname{argmax}_{\mathbf{w}} \frac{p(D | \mathbf{w}) p(\mathbf{w})}{p(D)}$$



$$\mathbf{w}_{MAP} = \operatorname{argmax}_{\mathbf{w}} p(D | \mathbf{w}) p(\mathbf{w})$$

# Maximum A Posteriori Estimates

- Dataset  $D = (x_1, x_2, \dots, x_N)$  of  $N$  independent observations.
- Model the data distribution with some  $p(D | \mathbf{w})$  parametrized by  $\mathbf{w}$
- MAP estimate: choose most probable  $\mathbf{w}$  given the data.

$$\mathbf{w}_{MAP} = \operatorname{argmax}_{\mathbf{w}} p(\mathbf{w} | D)$$

$$\mathbf{w}_{MAP} = \operatorname{argmax}_{\mathbf{w}} \frac{p(D | \mathbf{w}) p(\mathbf{w})}{p(D)}$$

$$\mathbf{w}_{MAP} = \operatorname{argmax}_{\mathbf{w}} p(D | \mathbf{w}) p(\mathbf{w})$$

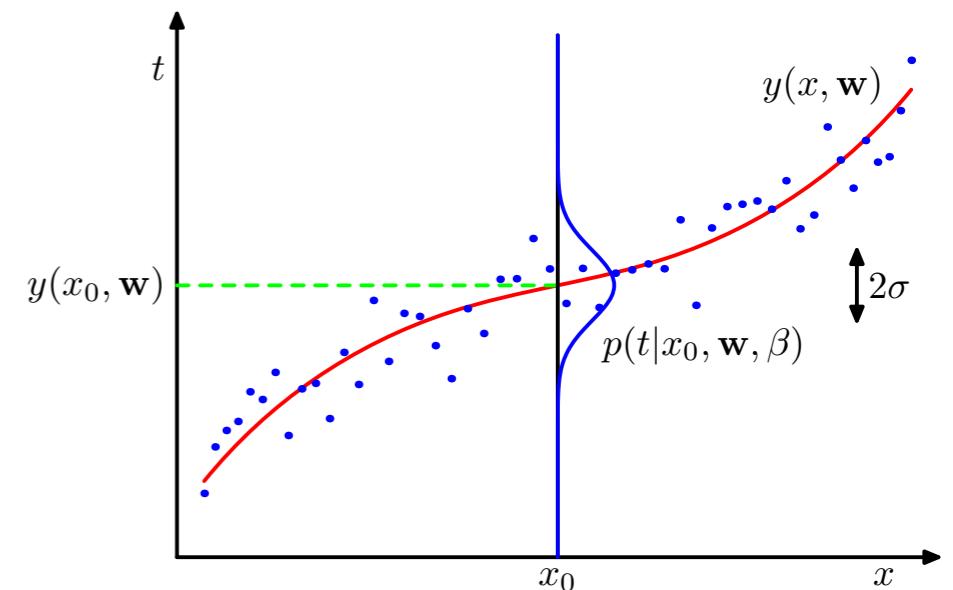
$$\mathbf{w}_{MAP} = \operatorname{argmin}_{\mathbf{w}} -\log p(D | \mathbf{w}) - \log p(\mathbf{w})$$

log brick  
neg log likelihood  
log prior

# MAP - Example Gaussian distribution

- **Likelihood**/Data model:

$$\begin{aligned} p(t|x, \mathbf{w}, \beta) &= \mathcal{N}(t|y(x, \mathbf{w}), \beta^{-1}) \\ &= \sqrt{\frac{\beta}{2\pi}} \exp \left[ -\frac{\beta}{2}(t - y(x, \mathbf{w}))^2 \right] \end{aligned}$$





$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} & \Sigma_{13} \\ \Sigma_{21} & \Sigma_{22} & \Sigma_{23} \\ \vdots & \vdots & \vdots \end{pmatrix}$$

Covariance  $[w_1, w_2] = 0$

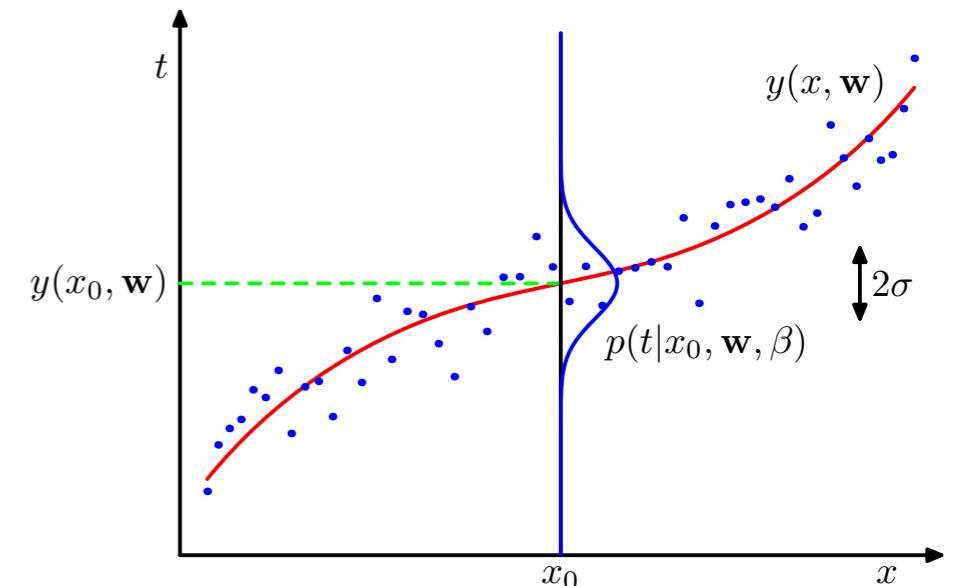
$$\Sigma = \begin{pmatrix} \Sigma_{11} & 0 & 0 \\ 0 & \Sigma_{22} & 0 \\ 0 & 0 & \Sigma_{33} \end{pmatrix}$$

$$= \begin{pmatrix} \alpha^{-1} & 0 & 0 \\ 0 & \alpha^{-1} & 0 \\ 0 & 0 & \alpha^{-1} \end{pmatrix} = \alpha^{-1} \text{Id}$$

# MAP - Example Gaussian distribution

- **Likelihood**/Data model:

$$\begin{aligned} p(t | x, \mathbf{w}, \beta) &= \mathcal{N}(t | y(x, \mathbf{w}), \beta^{-1}) \\ &= \sqrt{\frac{\beta}{2\pi}} \exp \left[ -\frac{\beta}{2}(t - y(x, \mathbf{w}))^2 \right] \end{aligned}$$



- **Prior** (assumption:  $\mathbf{w}$  takes on small values with certainty  $\alpha$ ):

$$\begin{aligned} p(\mathbf{w} | \alpha) &= \prod_{i=1}^M \mathcal{N}(w_i | 0, \alpha^{-1}) \\ &= \left( \frac{\alpha}{2\pi} \right)^{M/2} \prod_{i=1}^M e^{-\frac{\alpha}{2} w_i^2} \\ &= \left( \frac{\alpha}{2\pi} \right)^{M/2} e^{-\frac{\alpha}{2} \mathbf{w}^T \mathbf{w}} \end{aligned}$$

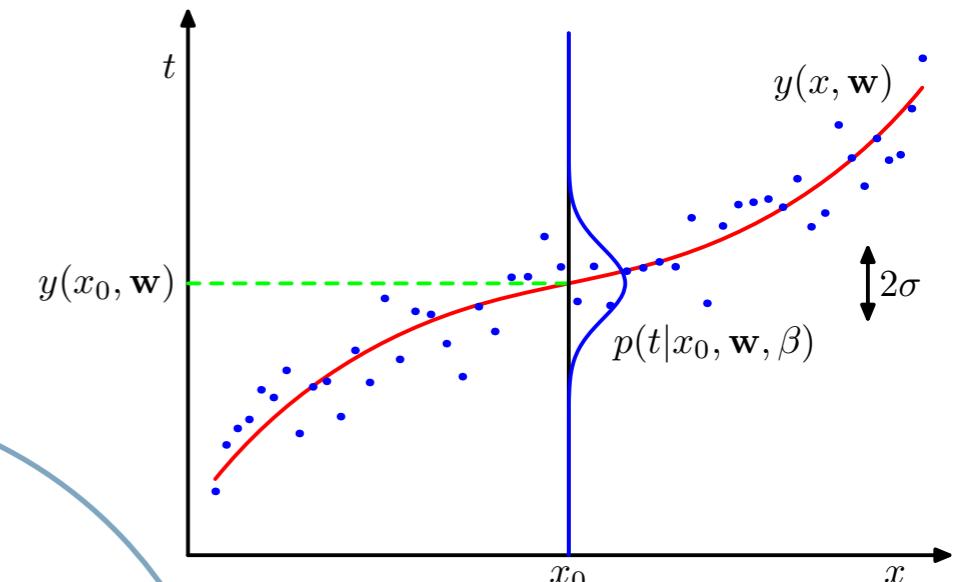
- **MAP solution** obtained by solving:

$$\mathbf{w}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmin}} -\log p(D | \mathbf{w}, \beta) - \log p(\mathbf{w}, \alpha)$$

# MAP - Example Gaussian distribution

- › **Likelihood**/Data model:

$$\begin{aligned} p(t | x, \mathbf{w}, \beta) &= \mathcal{N}(t | y(x, \mathbf{w}), \beta^{-1}) \\ &= \sqrt{\frac{\beta}{2\pi}} \exp \left[ -\frac{\beta}{2}(t - y(x, \mathbf{w}))^2 \right] \end{aligned}$$



- › **Prior** (assumption:  $\mathbf{w}$  takes on small values with certainty  $\alpha$ ):

$$\begin{aligned} p(\mathbf{w} | \alpha) &= \prod_{i=1}^M \mathcal{N}(w_i | 0, \alpha^{-1}) \\ &= \left( \frac{\alpha}{2\pi} \right)^{M/2} \prod_{i=1}^M e^{-\frac{\alpha}{2} w_i^2} \\ &= \left( \frac{\alpha}{2\pi} \right)^{M/2} e^{-\frac{\alpha}{2} \mathbf{w}^T \mathbf{w}} \end{aligned}$$

$$\begin{aligned} \log p(D | \mathbf{w}, \beta) &= -\frac{\beta}{2} \sum_{i=1}^N (y(x_i, \mathbf{w}) - t_i)^2 \\ &\quad + \frac{N}{2} \log \beta - \frac{N}{2} \log 2\pi \end{aligned}$$

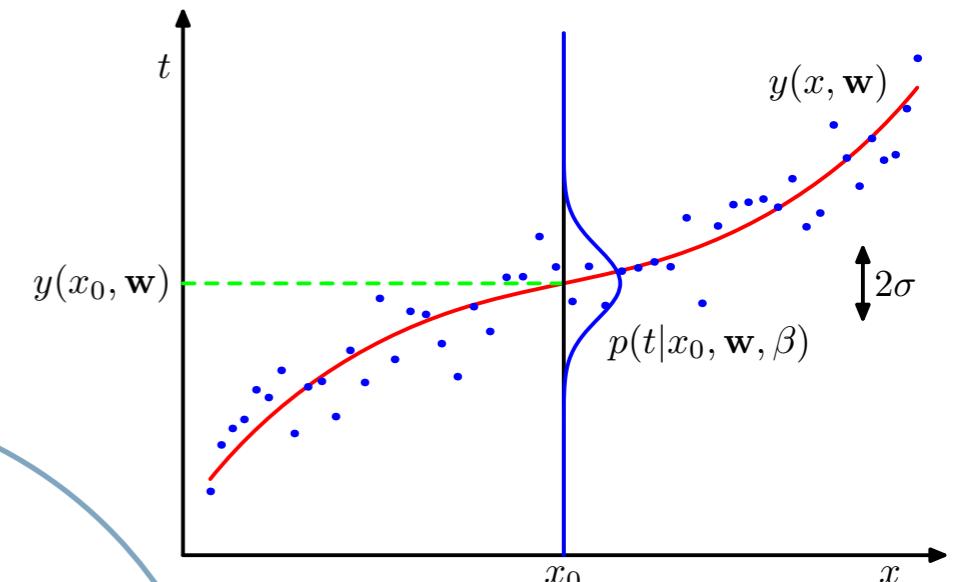
- › **MAP solution** obtained by solving:

$$\mathbf{w}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{\beta}{2} \sum_{i=1}^N (y(x_i, \mathbf{w}) - t_i)^2 - \log p(\mathbf{w}, \alpha)$$

# MAP - Example Gaussian distribution

- › **Likelihood**/Data model:

$$\begin{aligned} p(t | x, \mathbf{w}, \beta) &= \mathcal{N}(t | y(x, \mathbf{w}), \beta^{-1}) \\ &= \sqrt{\frac{\beta}{2\pi}} \exp \left[ -\frac{\beta}{2}(t - y(x, \mathbf{w}))^2 \right] \end{aligned}$$



- › **Prior** (assumption:  $\mathbf{w}$  takes on small values with certainty  $\alpha$ ):

$$\begin{aligned} p(\mathbf{w} | \alpha) &= \prod_{i=1}^M \mathcal{N}(w_i | 0, \alpha^{-1}) \\ &= \left( \frac{\alpha}{2\pi} \right)^{M/2} \prod_{i=1}^M e^{-\frac{\alpha}{2} w_i^2} \\ &= \left( \frac{\alpha}{2\pi} \right)^{M/2} e^{-\frac{\alpha}{2} \mathbf{w}^T \mathbf{w}} \end{aligned}$$

$\log p(D | \mathbf{w}, \beta) = -\frac{\beta}{2} \sum_{i=1}^N (y(x_i, \mathbf{w}) - t_i)^2 + \frac{N}{2} \log \beta - \frac{N}{2} \log 2\pi$

$\log p(\mathbf{w} | \alpha) = -\frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + \frac{N}{2} \log \alpha - \frac{N}{2} \log 2\pi$

- › **MAP solution** obtained by solving:

$$\mathbf{w}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{\beta}{2} \sum_{i=1}^N (y(x_i, \mathbf{w}) - t_i)^2 + \frac{\alpha}{2} \mathbf{w}^T \mathbf{w}$$



# Machine Learning 1

Lecture 2.6 - Bayesian Prediction

Erik Bekkers

(Bishop 1.2.6)



***Three Statistical Learning Principles:***

- Maximum Likelihood
- Maximum A Posteriori
- Bayesian Prediction

# Bayesian Approach

- › Dataset  $D = (x_1, x_2, \dots, x_N)$  of  $N$  independent observations.
- › Frequentist approach: search for **one** optimal **estimate** of  $\mathbf{w}$

$$\mathbf{w}_{ML} = \underset{\mathbf{w}}{\operatorname{argmax}} p(D | \mathbf{w})$$
$$\mathbf{w}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmax}} p(\mathbf{w} | D)$$

point estimate

- › Bayesian approach: Given a prior belief over  $\mathbf{w}$ ,  $p(\mathbf{w})$ , and our data  $D$ , we are interested in the entire posterior distribution

$$p(\mathbf{w} | D) = \frac{p(D | \mathbf{w}) p(\mathbf{w})}{p(D)}$$

Here we do not pick one value for  $\mathbf{w}$  but consider all of them. We will weigh each model parameterized with  $\mathbf{w}$  according to their posterior probability

- ›  $p(\mathbf{w} | D)$  reflects the **plausibility** of different  $\mathbf{w}$ , given our prior knowledge and how likely our data is generated using  $\mathbf{w}$ .

- we consider both  $x$  and  $\underline{w}$  random variables according to joint prob  $p(x, \underline{w}|D)$ . i.e.  $p(x, \underline{w}|D)$  gives the prob. of observing a new  $x'$  and model params  $\underline{w}$  given the already observed dataset  $D$ .
- In the predictive setting we are only interested in  $x'$ , and thus we compute the marginal  $p(x'|D)$  by integrating out  $\underline{w}$ .
- Following the product rule  

$$p(x', \underline{w}|D) = p(x'|\underline{w}, D) p(\underline{w}|D)$$
the marginal gets the interpretation of a weighted average sum/integral.  
Weighting each predictive model  $p(x'|\underline{w})$  with the posterior  $p(\underline{w}|D)$  for that  $\underline{w}$ .

conditional



# Three Statistical Learning Principles

Three general statistical learning principles to go from data to models (parametric predictive/proposal distributions):

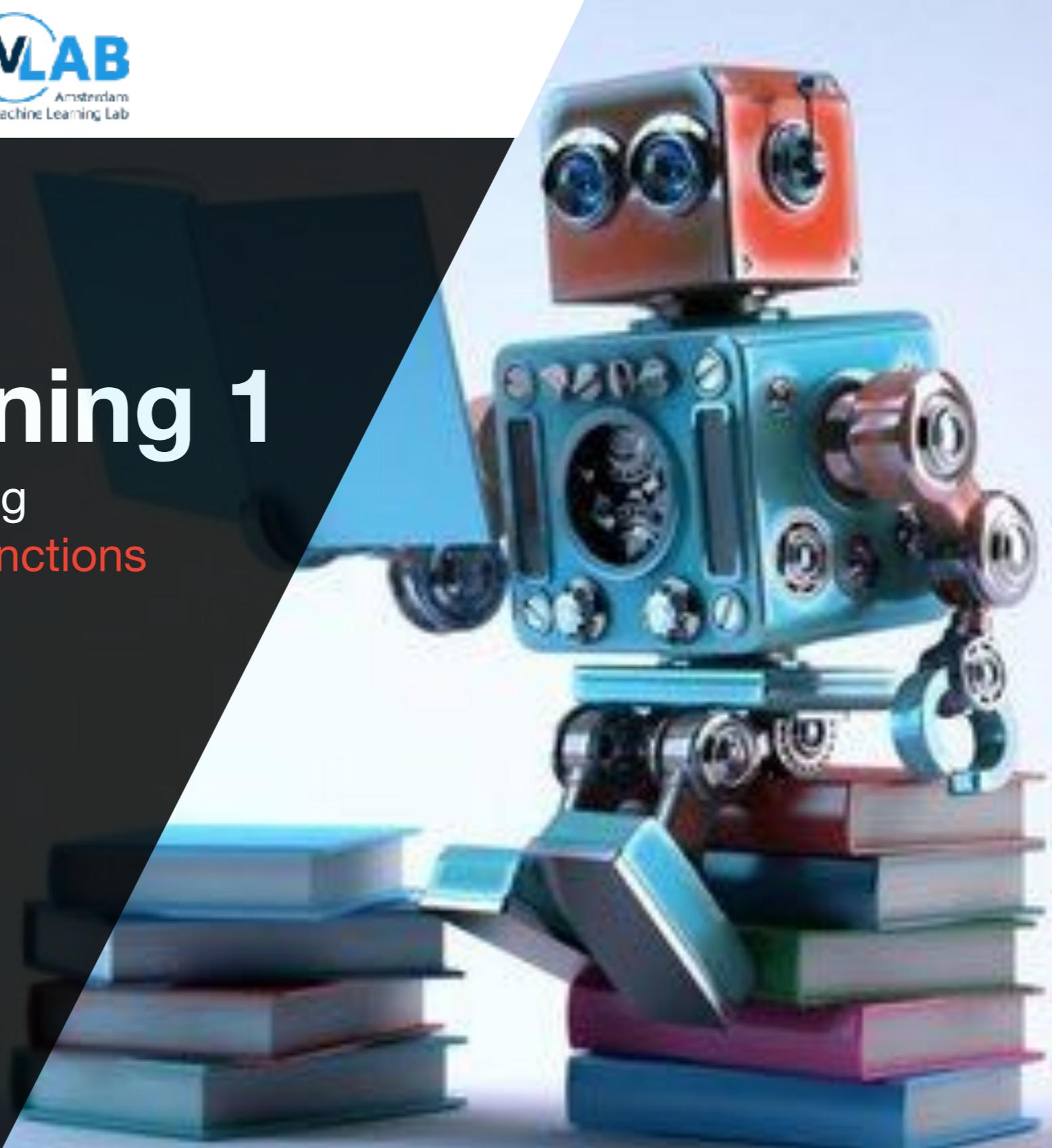
- I. Maximum likelihood
- II. Maximum a posteriori
- III. Bayesian prediction

# Machine Learning 1

Lecture 3.1 - Supervised Learning  
**Linear Regression With Basis Functions**

*Erik Bekkers*

*(Bishop 3.1)*





# Machine Learning

Lecture  
Linear F

Erik Bek  
(Bishop)

# Machine Learning 1

*The probabilistic  
view*



W7

Kernel methods

W6

Continuous Latent variable models

W5

Discrete Latent variable models

W4

Neural networks

W3

Classification

W2

Regression

W1

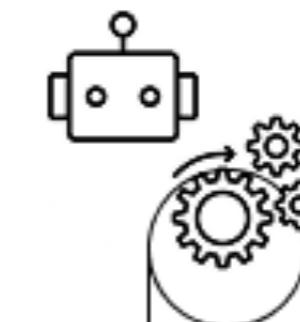
Probabilistic modeling

Random variables,  
distributions, maximum  
likelihood, MAP, ...

Optimization

Exact/analytic, (stochastic)  
gradient descent, constraint  
optimization, 2nd order

*The algorithmic  
view*



Mathematical tools

Multi-variate calculus,  
linear algebra

# Linear Regression

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

- Input variables  $\mathbf{x}_n \in \mathbb{R}^d$

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

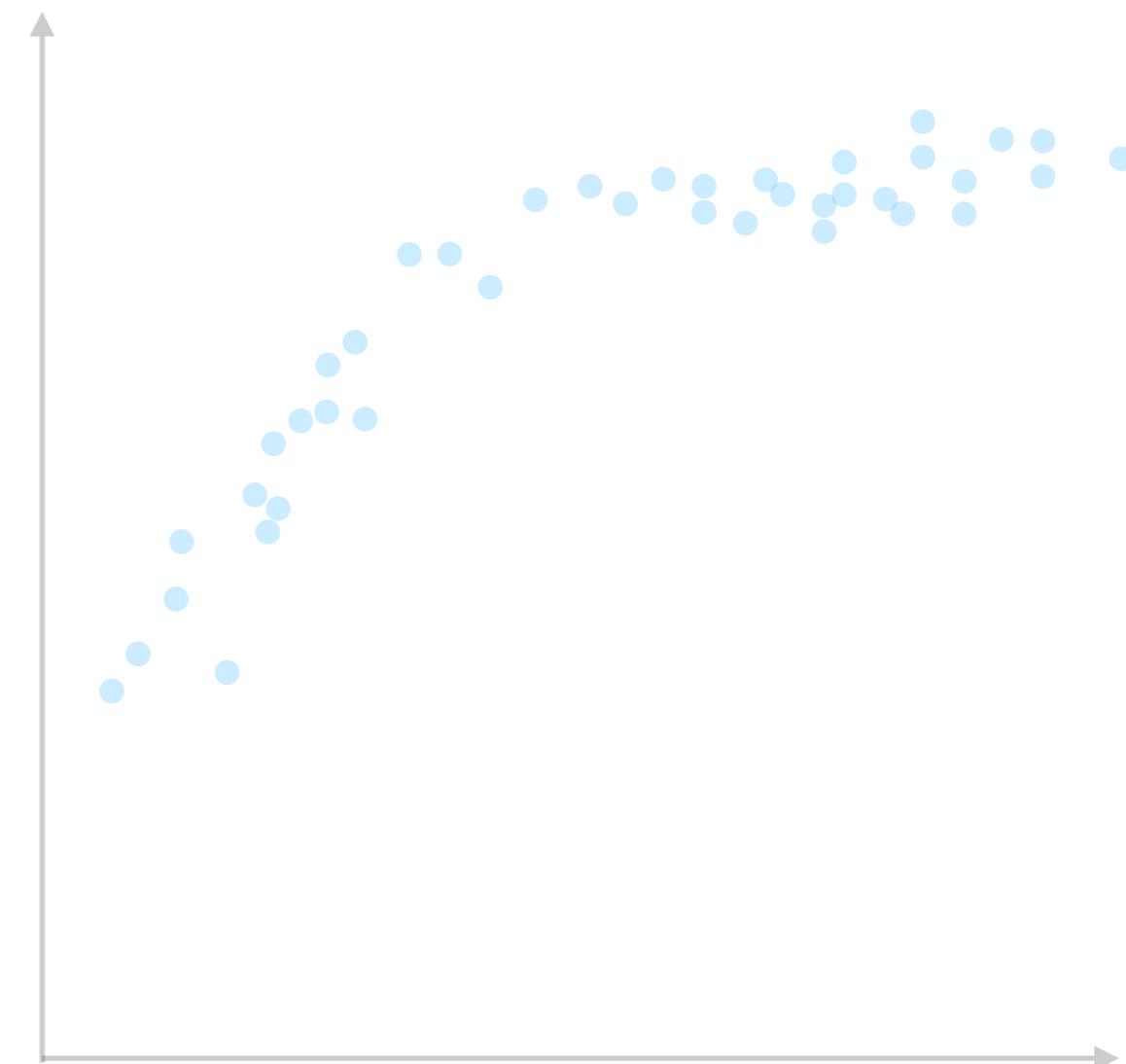
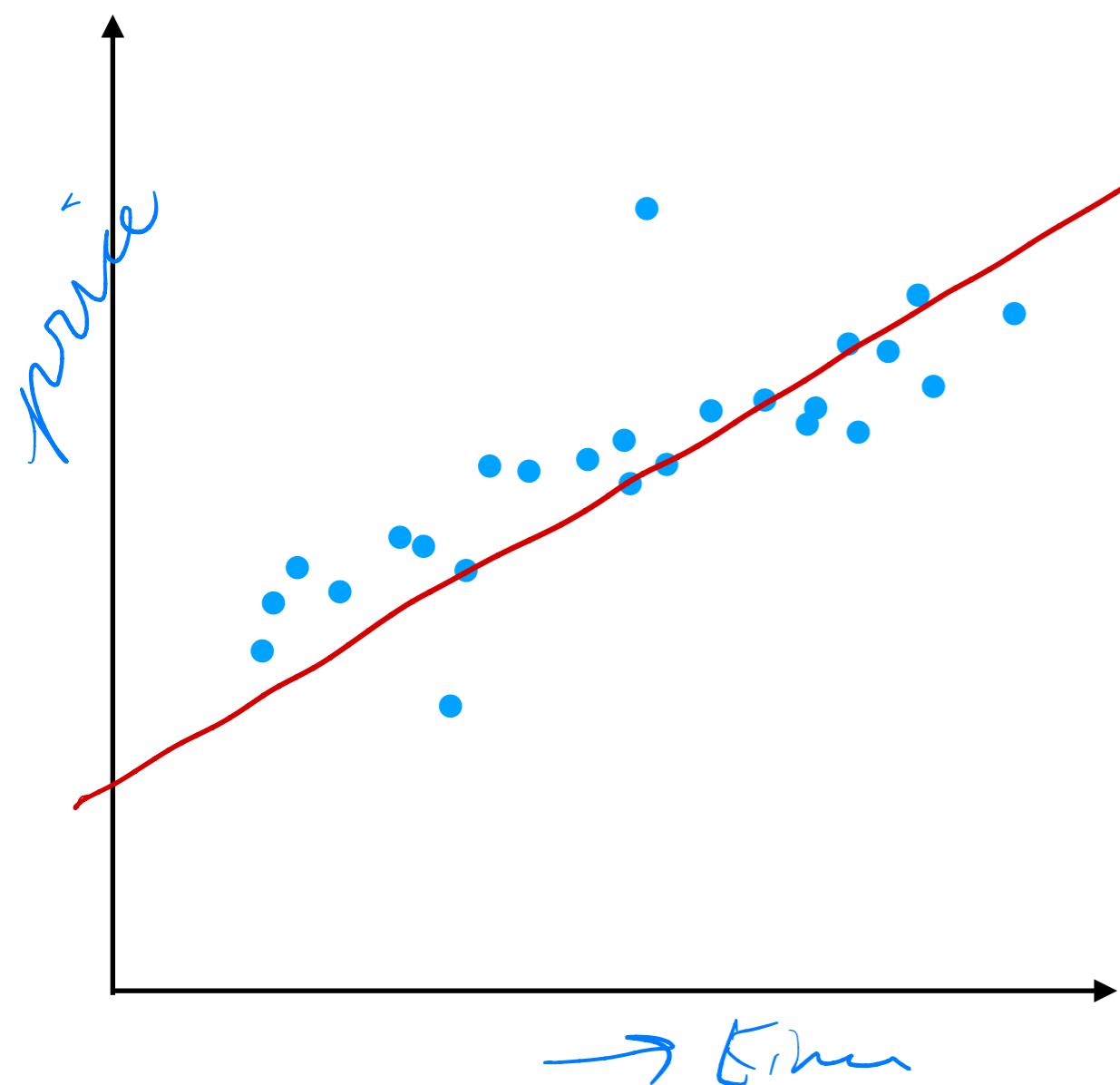
- Simplest linear model:

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + w_2 \cdot x_2 + \dots + w_d x_d$$

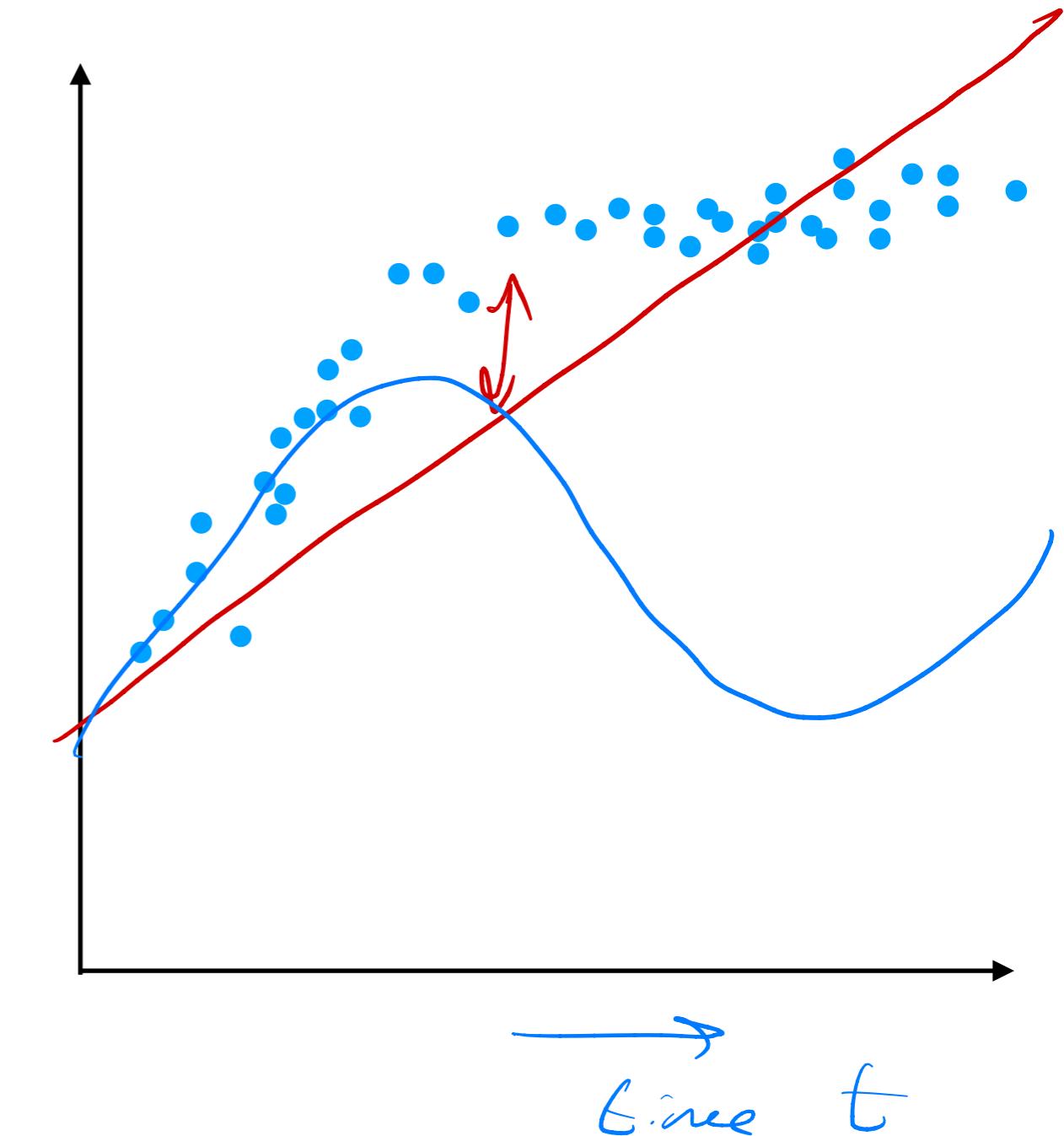
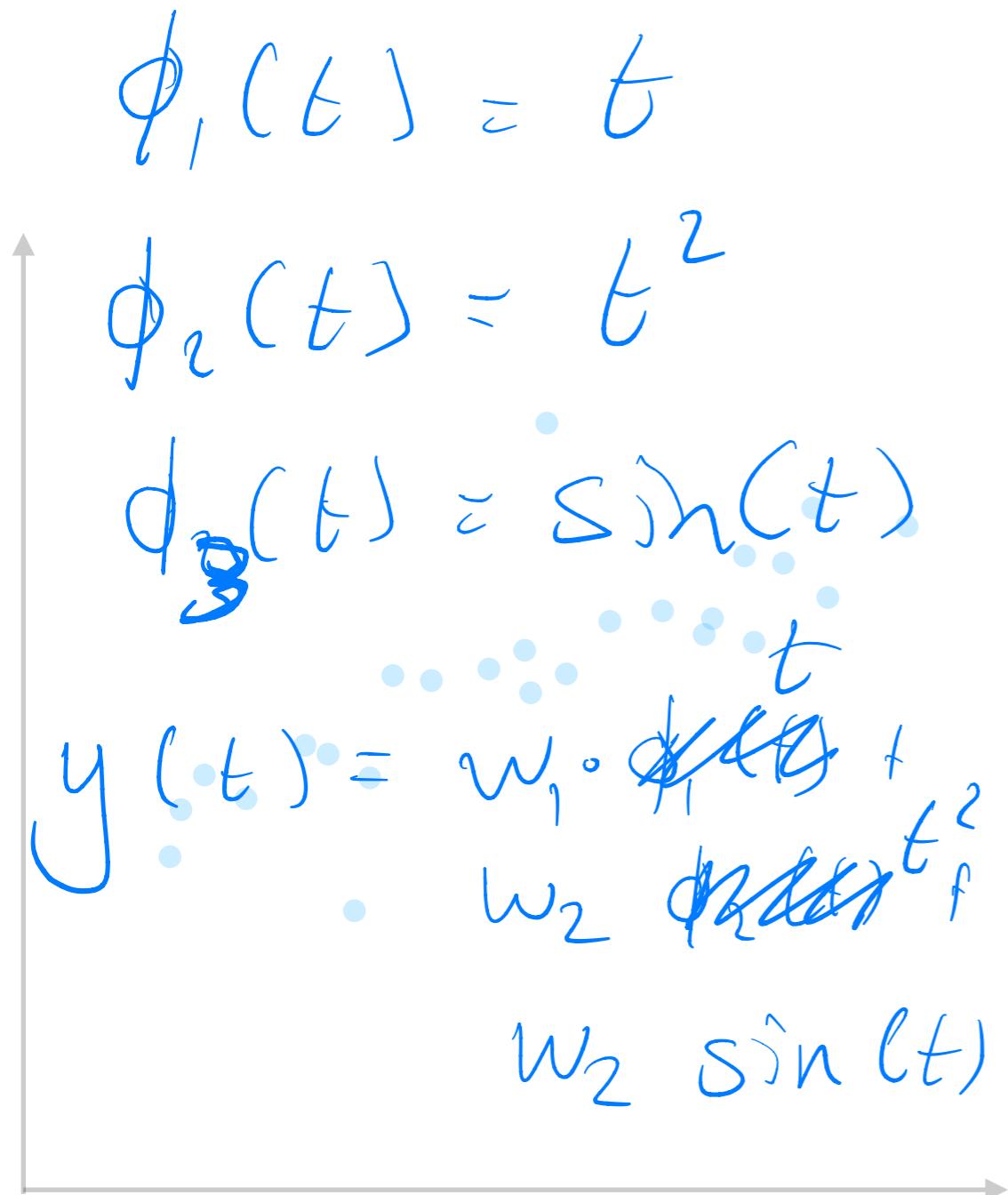


$$\underline{\mathbf{X}} = \begin{pmatrix} x_{n1} \\ x_{n2} \\ x_{n3} \\ \vdots \\ x_{nd} \end{pmatrix} \in \mathbb{R}^d, \quad \underline{\mathbf{w}} = \begin{pmatrix} w_0 \\ w_1 \\ w_2 \\ \vdots \\ w_d \end{pmatrix}$$

# Linear Regression



# Linear Regression

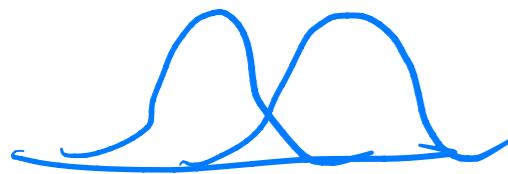






# Example: Basis Functions (II)

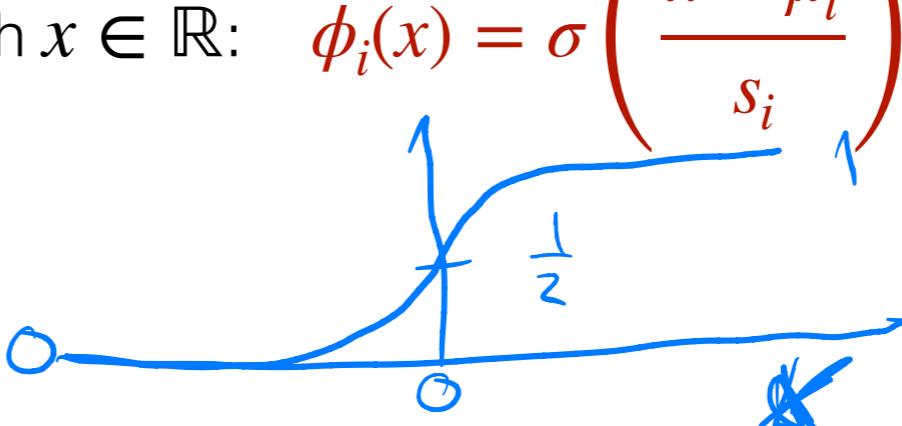
- Gaussian basis functions with  $\mathbf{x} \in \mathbb{R}^D$ :  $\phi_i(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i)\right)$



$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{i=1}^{M-1} w_i \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_i)^T \Sigma_i^{-1} (\mathbf{x} - \boldsymbol{\mu}_i)\right)$$

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

with  $\sigma(x) = \frac{1}{1 + e^{-x}}$



$$y(\mathbf{x}, \mathbf{w}) = w_0 + \sum_{i=1}^{M-1} w_i \sigma\left(\frac{x - \mu_i}{s_i}\right)$$

# Machine Learning 1

Lecture 3.2 - Supervised Learning

**Linear Regression via Maximum Likelihood  
Optimization**

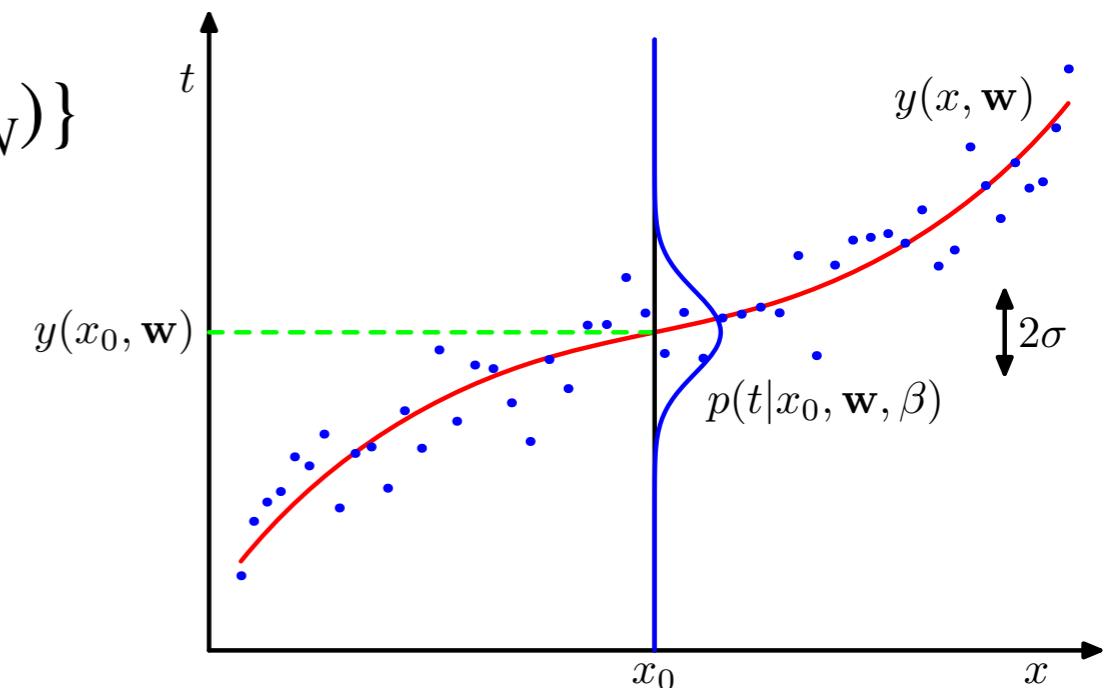
*Erik Bekkers*

*(Bishop 3.1.1)*



# Linear Regression

- ▶ Regression  $D = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N)\}$
- ▶ Input variables  $\mathbf{x}_i \in \mathbb{R}^D$
- ▶ Target variables  $t_i$



**Figure:** Gaussian conditional distribution  
(Bishop 1.16)

- ▶ Linear model with basis functions

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

$\cup$   
 $\mathbb{R}^M$

# Maximum Likelihood

- Model assumptions/choices:

- Data is generated through model  $y(\mathbf{x}, \mathbf{w}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$
- Observations  $t_i$  are subject to Gaussian noise

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon , \quad \epsilon \sim \mathcal{N}(0, \beta^{-1})$$

- Then data distribution (likelihood):

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

$\left( \begin{array}{c} \mathbf{x}_1 \\ \vdots \\ \mathbf{x}_N \end{array} \right) \in \mathbb{R}^{N \times D}$  and  $\mathbf{t} = (t_1, \dots, t_N)^T \in \mathbb{R}^N$

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

- Likelihood function:

$$p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = \prod_{i=1}^N \sqrt{\frac{\beta}{2\pi}} e^{-\frac{\beta}{2}(t_i - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i))^2}$$

# ML: Sum-of-Squares Error

- ▶ Likelihood:

$$p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = \prod_{i=1}^N \mathcal{N}(t_i | \mathbf{w}^T \phi(\mathbf{x}_i), \beta^{-1})$$

- ▶ Log likelihood

$$\log p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = \frac{N}{2} \log \beta - \frac{N}{2} \log 2\pi - \frac{\beta}{2} \sum_{i=1}^N (t_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2$$

- ▶ Sum-of-squares error:

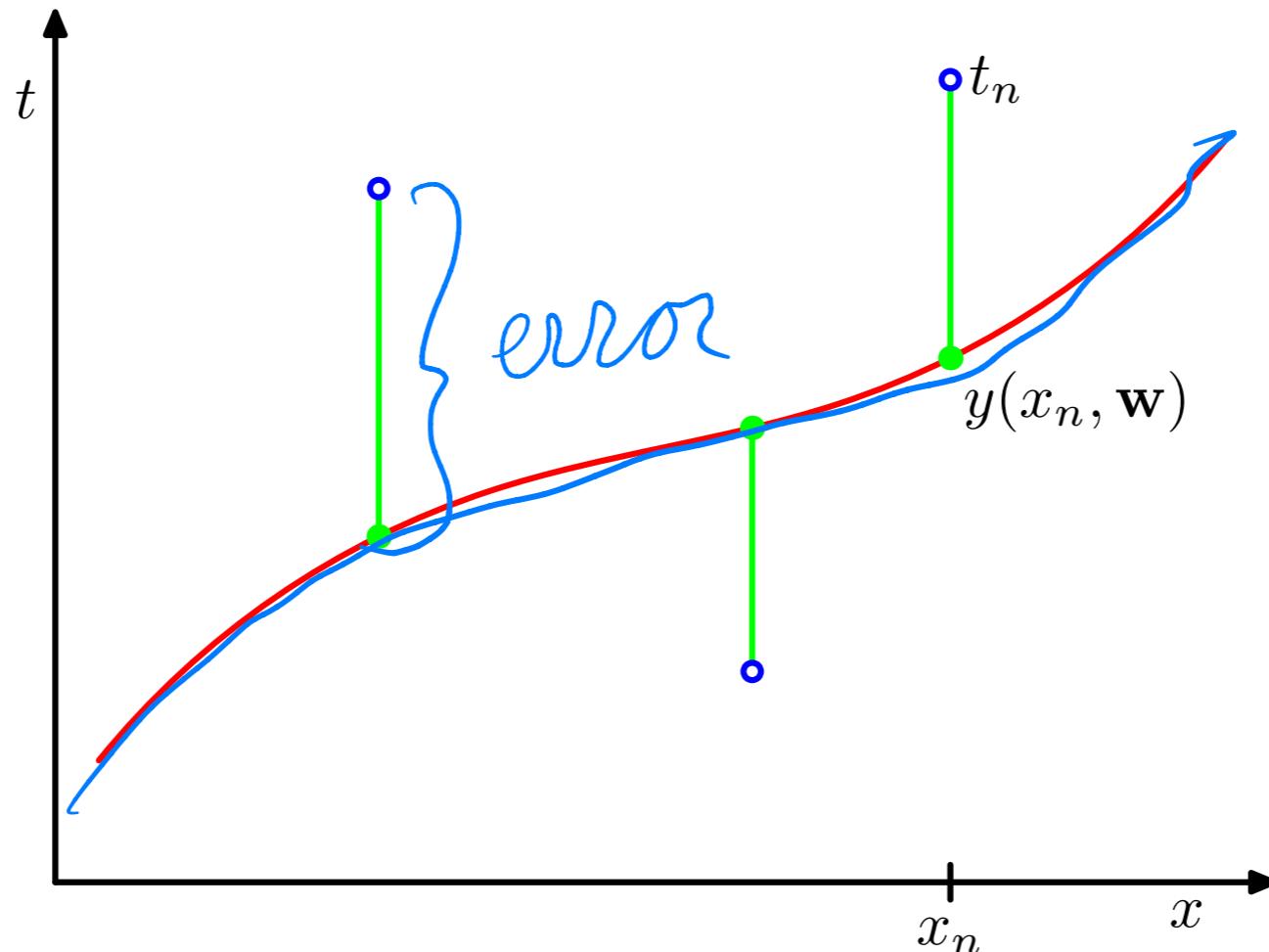
$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N (t_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2$$

- ▶ For comparison of different dataset sizes  $N$

$$E_D^{RMSE}(\mathbf{w}) = \sqrt{\frac{1}{N} \sum_{i=1}^N (\hat{t}_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2}$$

X no dependency on dataset size  
X sensible error in units, e.g. euros instead of squared euros

# Example: Sum-of-Squares Error



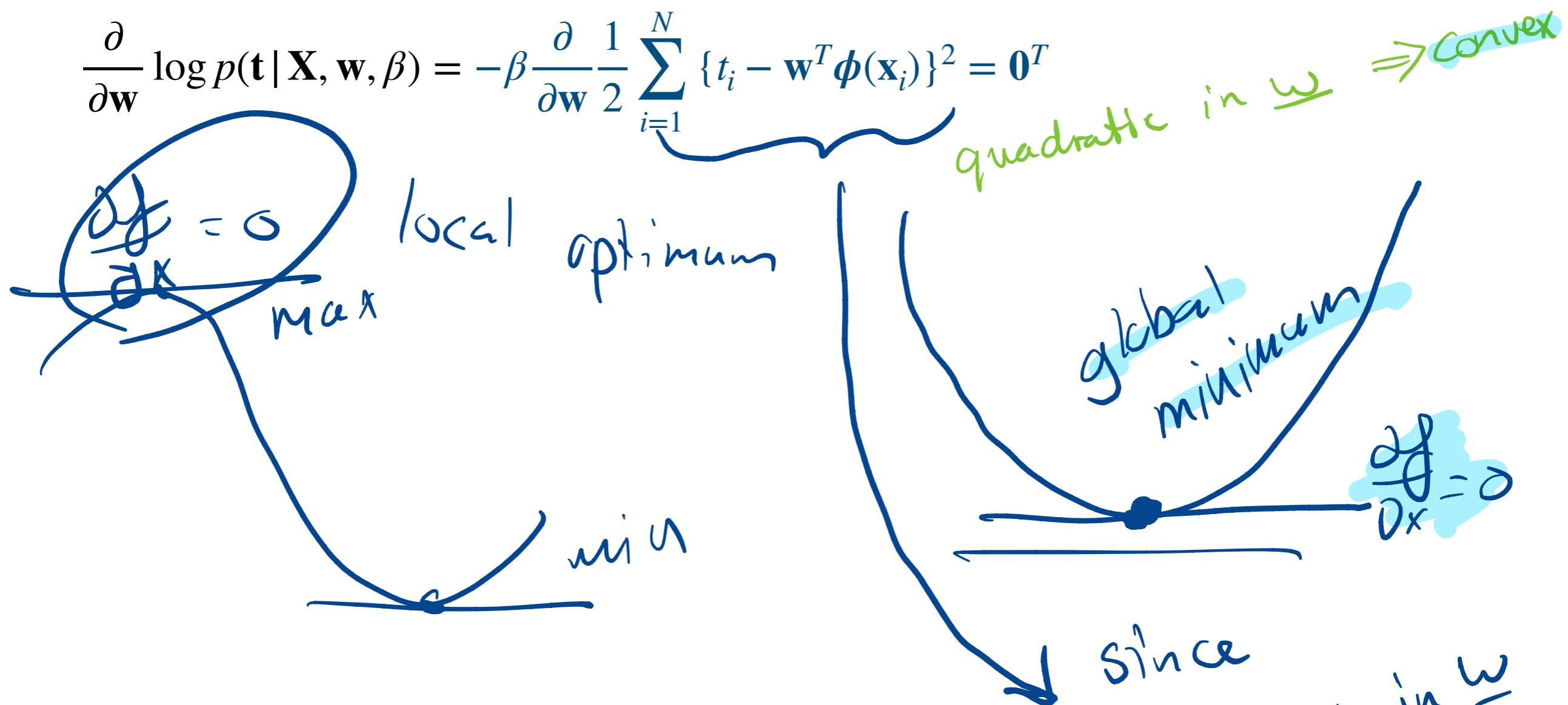
**Figure:** Errors are given by half the squares of green bars (Bishop 1.3)

# Maximum Likelihood Estimates

- Maximize the log likelihood  $\log p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta)$

- Analytic solution by solving:  $\frac{\partial}{\partial \mathbf{w}} \log p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = 0 \quad \text{for } \mathbf{w} \in \mathbb{R}^M$

$$\frac{\partial}{\partial \mathbf{w}} \log p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = -\beta \frac{\partial}{\partial \mathbf{w}} \frac{1}{2} \sum_{i=1}^N \{t_i - \mathbf{w}^T \phi(\mathbf{x}_i)\}^2 = \mathbf{0}^T$$



Convention  $\nabla_{\mathbf{x}} f(\mathbf{x}) := \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} := \left( \frac{\partial f(\mathbf{x})}{\partial x_1}, \frac{\partial f(\mathbf{x})}{\partial x_2}, \dots \right)$  is a row vector



# Maximum Likelihood Estimates

- Maximize the log likelihood  $\log p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta)$
- Analytic solution by solving:  $\frac{\partial}{\partial \mathbf{w}} \log p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = 0 \quad \text{for } \mathbf{w}$

$$\frac{\partial}{\partial \mathbf{w}} \log p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = -\beta \frac{\partial}{\partial \mathbf{w}} \frac{1}{2} \sum_{i=1}^N \{t_i - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i)\}^2 = \mathbf{0}^T$$

$(\beta > 0, \text{ and apply derivative})$   
 $\Leftrightarrow$

$$\sum_{i=1}^N \{t_i - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i)\} \boldsymbol{\phi}(\mathbf{x}_i)^T = \mathbf{0}^T$$

$(\text{rearrange terms})$   
 $\Leftrightarrow$

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

Convention  $\nabla_{\mathbf{x}} f(\mathbf{x}) := \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}} := \left( \frac{\partial f(\mathbf{x})}{\partial x_1}, \frac{\partial f(\mathbf{x})}{\partial x_2}, \dots \right)$  is a row vector





# Machine Learning 1

Lecture 3.3 - Supervised Learning  
**Stochastic Gradient Descent**

*Erik Bekkers*

*(Bishop 3.1.3)*





# Stochastic gradient descent

- For  $N \gg 1$  solution  $\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$  is very costly to compute!

- Needs to process all data  $(\mathbf{x}_1, \dots, \mathbf{x}_N)$  at once.

- Matrix inversion of  $M \times M$  matrix:  $O(M^3)$

Expensive

- Loss is a sum of error terms for each datapoint:

- $$E_D(\mathbf{w}) = \sum_{i=1}^N E(\mathbf{x}_i, t_i, \mathbf{w})$$

- $$E(\mathbf{x}_i, t_i, \mathbf{w}) = \frac{1}{2}(t_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2$$

- Approach for large dataset: stochastic gradient descent

Idea: approximate  $E_D$  with few data points

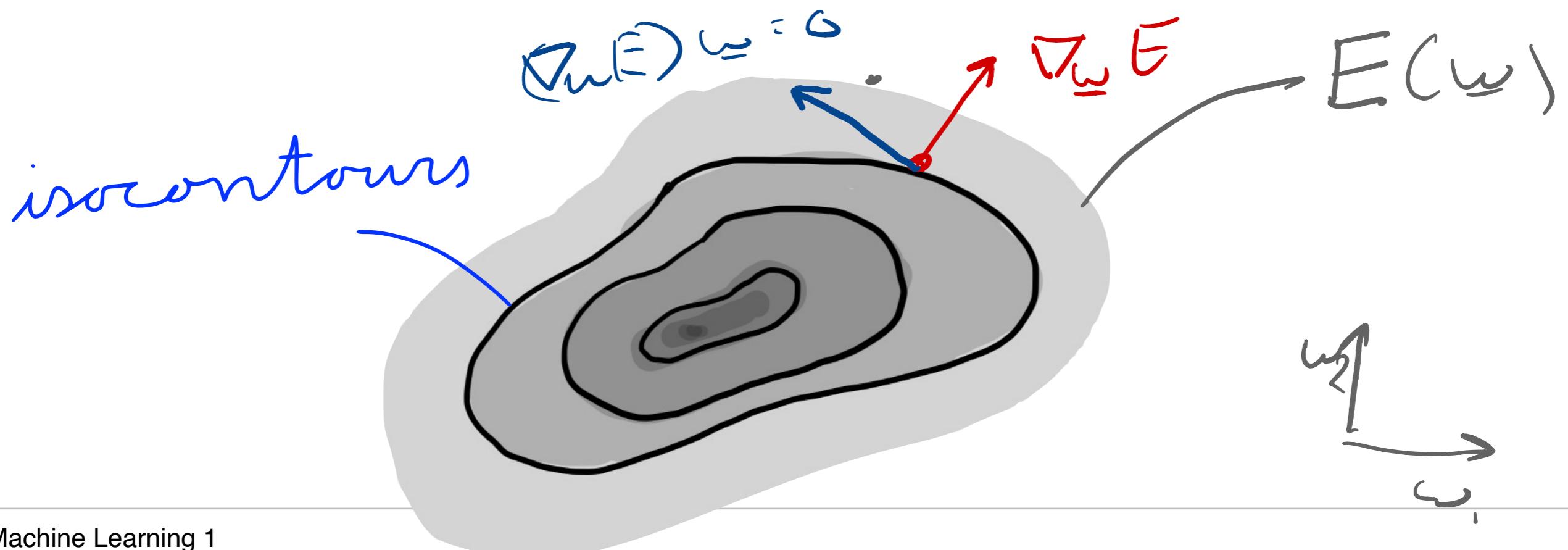
# Recap: The Gradient (see e.g. Bishop App. E)

- The gradient encodes all directional derivatives via scalar product

$$\nabla_{\underline{w}} E := \left( \frac{\partial E}{\partial w_1}, \frac{\partial E}{\partial w_2}, \dots \right) \quad \text{a row vector}$$

directional derivative  $\downarrow$  in direction  $\vec{v}$  is given by  $(\nabla_{\underline{w}} E) \cdot \vec{v}$

- The gradient always points in the direction of steepest ascent  
directional derivative is maximized when  $\vec{v} = \nabla E$
- The gradient is always perpendicular to the contours of a function



# Stochastic gradient descent

- Objective: minimize  $E_D(\mathbf{w}) = \sum_{i=1}^N E(\mathbf{x}_i, t_i, \mathbf{w})$ , with

$$E(\mathbf{x}_i, t_i, \mathbf{w}) = \frac{1}{2}(t_i - \mathbf{w}^T \phi(\mathbf{x}_i))^2$$

