

Machine Learning 1

Lecture 4 - Model Selection - Bias Variance
Decomposition - Gaussian Posteriors -
Sequential Bayesian Learning - Bayesian
Predictive Distributions

Erik Bekkers



Machine Learning 1

Lecture 3.5 - Supervised Learning
Regularized Least Squares

Erik Bekkers

(Bishop 3.1.4)



Parameter estimates with Gaussians

- Given **Likelihood**/Data model:

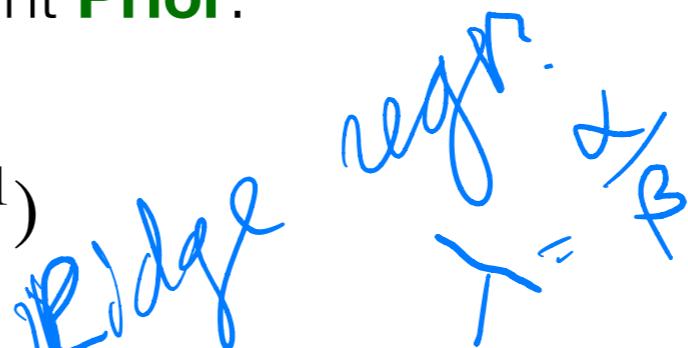
$$p(t | x, \mathbf{w}, \beta) = \mathcal{N}(t | y(x, \mathbf{w}), \beta^{-1})$$

- The **ML parameter estimate** is obtained via **least squares**:

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

- Additionally, given Gaussian weight **Prior**:

$$p(\mathbf{w} | \alpha) = \prod_{i=1}^M \mathcal{N}(w_i | 0, \alpha^{-1})$$



- The **MAP parameter estimate** is obtained via **regularized least squares**:

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

Example: Regularized Polynomial Regression

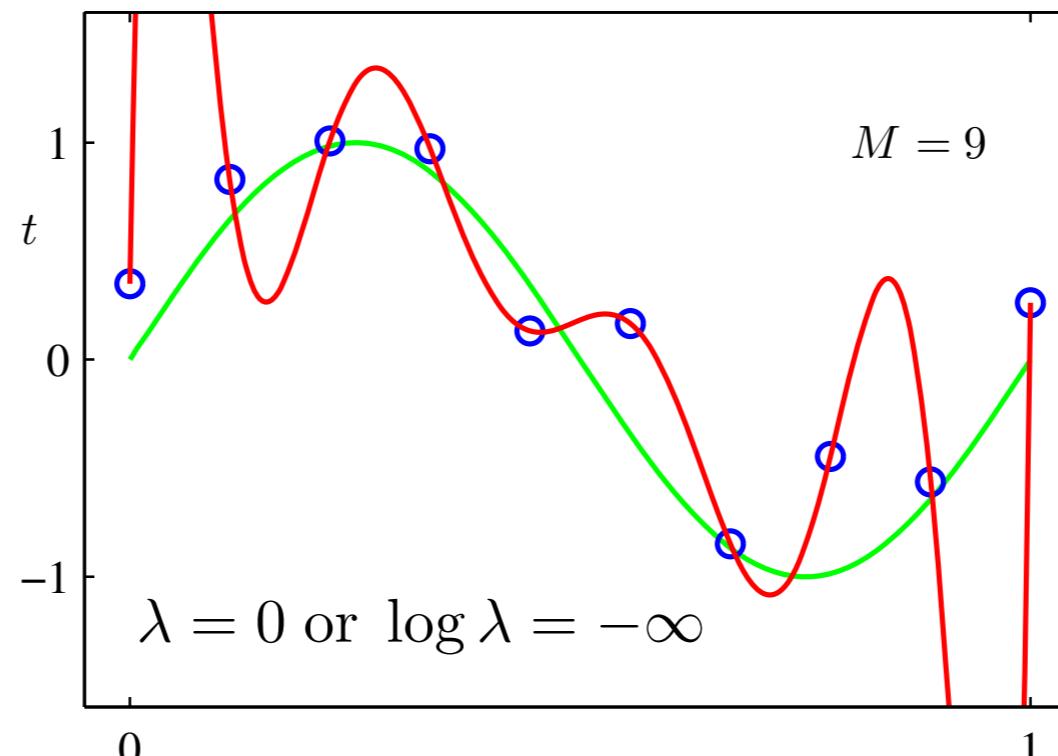


Figure: polynomial regression (Bishop 1.4)

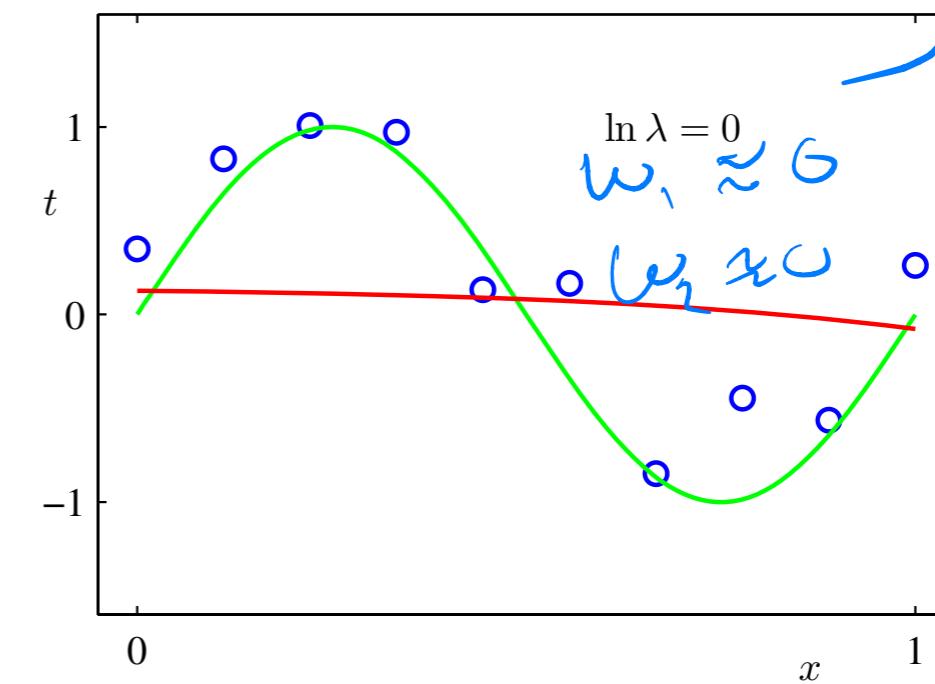
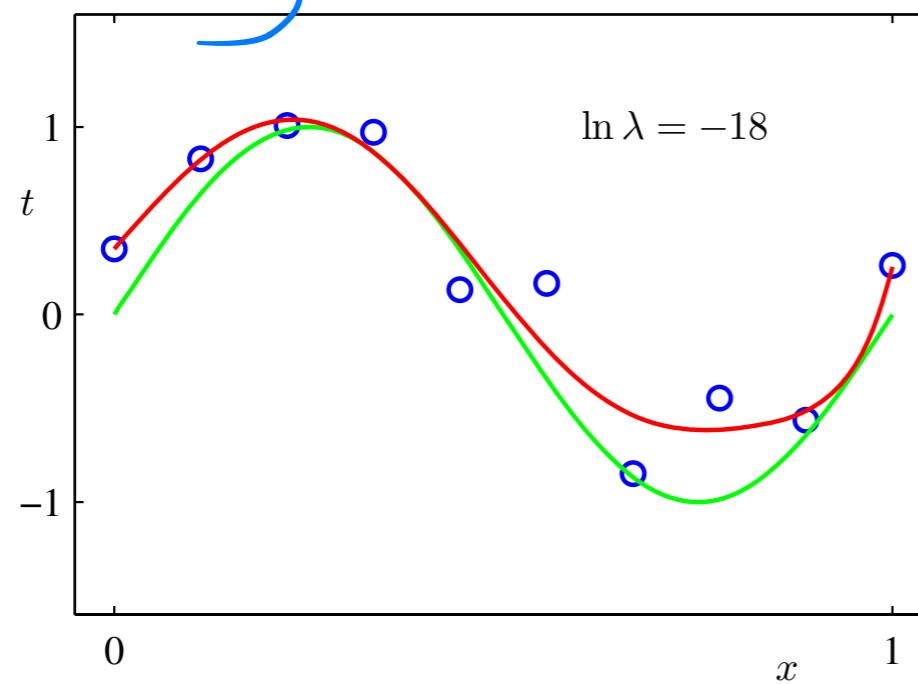


Figure: Regularized polynomial regression (Bishop 1.7)

Example: Regularized Polynomial Regression

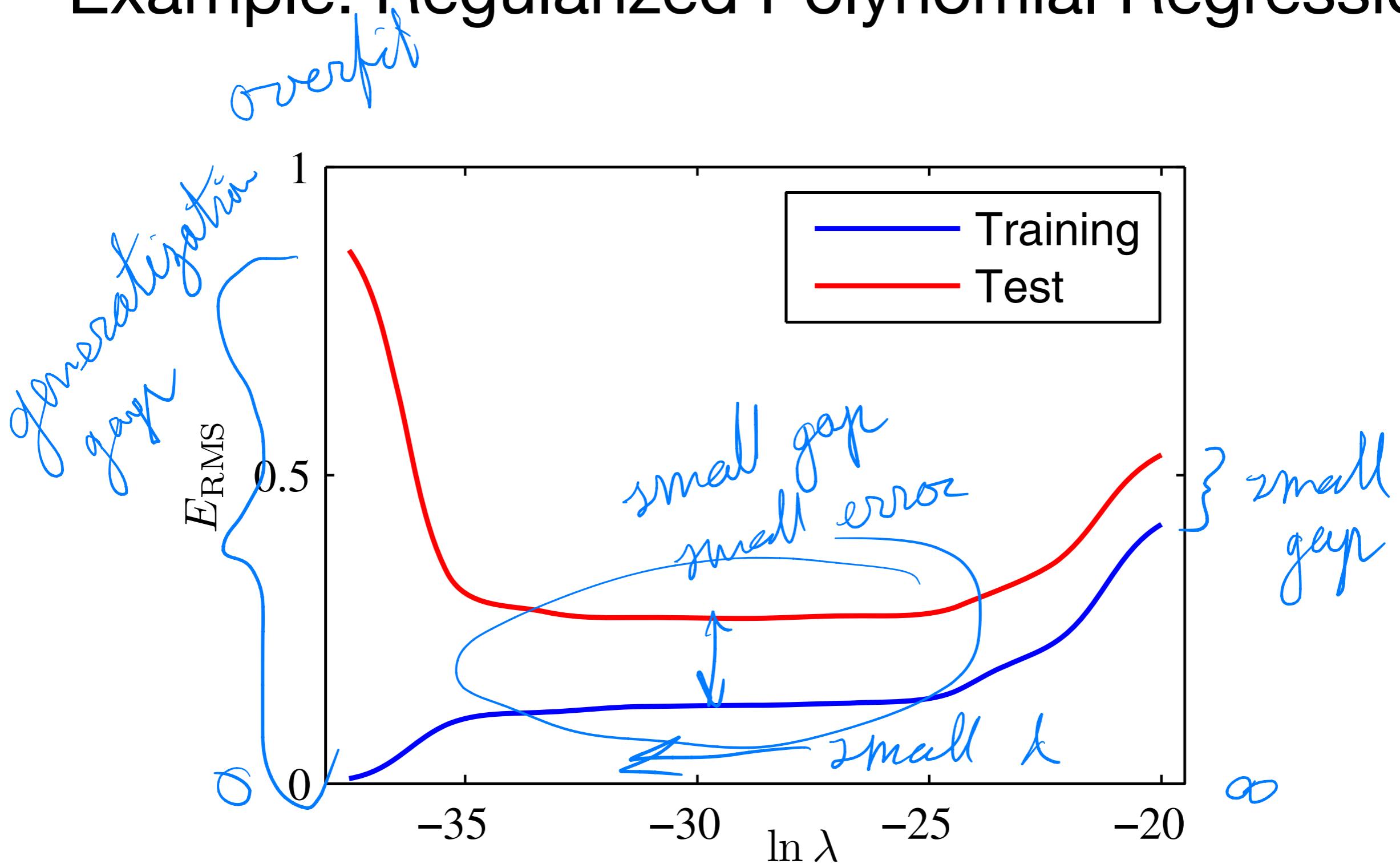


Figure: train and test errors for regularized M=9 polynomial regression (Bishop 1.8)

Regularized Least Squares: sparse weights

$$\text{minimize} \quad \frac{1}{2} \sum_{i=1}^N \{t_i - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i)\}^2 \quad \text{subject to constraint}$$

$$\sum_{i=1}^M |w_i|^q \leq \eta$$

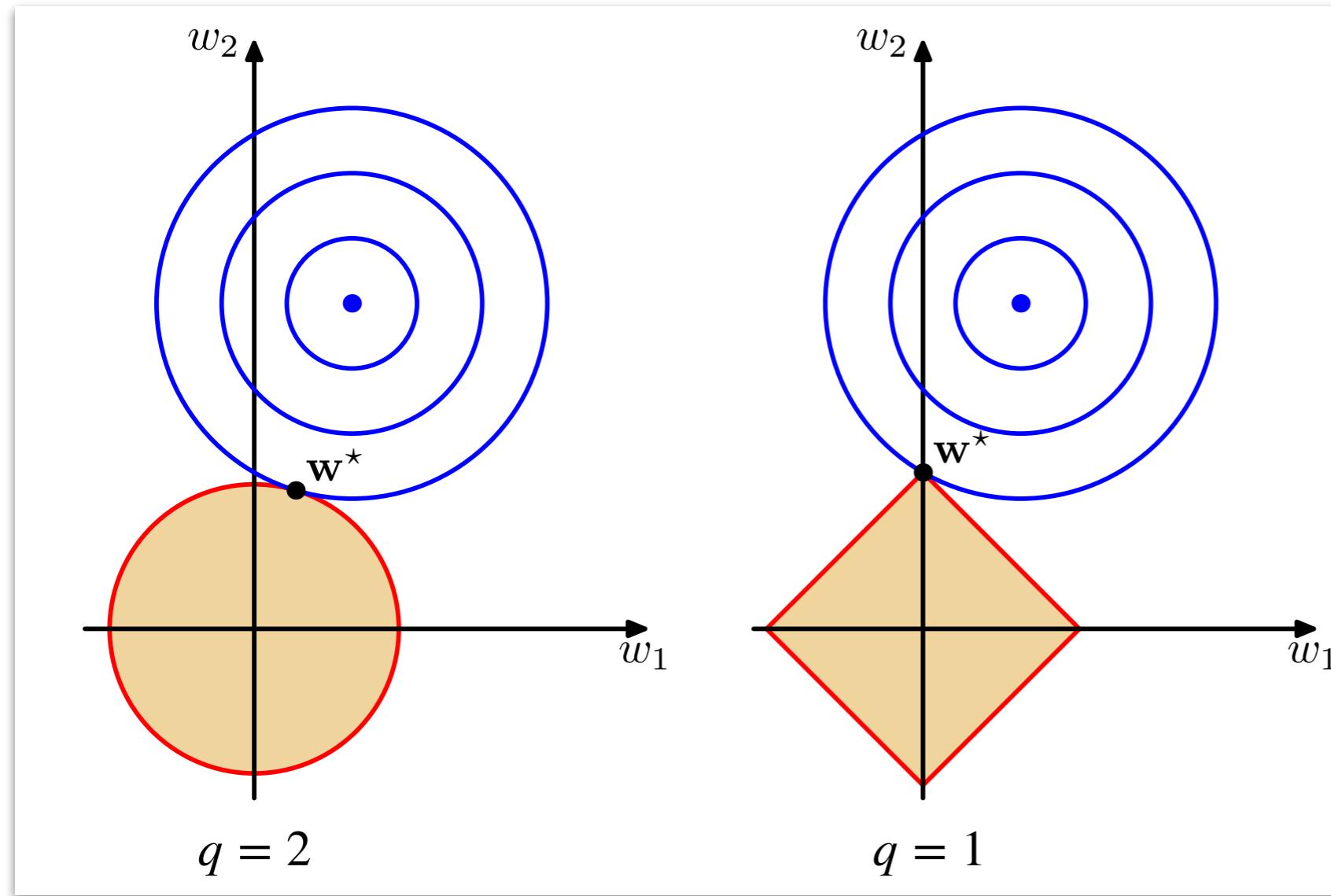


Figure: regularization as constrained optimization (Bishop 3.4)

Machine Learning 1

Lecture 4.1 - Supervised Learning
Model Selection

Erik Bekkers

(Bishop 1.3)

λ is a hyperparameter



Supervised Learning: Evaluating Errors

Q: How can we reliably estimate the model performance properly for unknown data?

Q: How can we choose the optimal hyperparameters?

Supervised Learning: Small Datasets

- Small dataset \rightarrow small validation and test set
- Noisy model selection and estimate of generalization error
- Cross-validation:
 - Split data $D = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N)\}$ in to K folds
 - Train model y on $K - 1$ folds (fold k left out) $\rightarrow \hat{y}^{-k}$

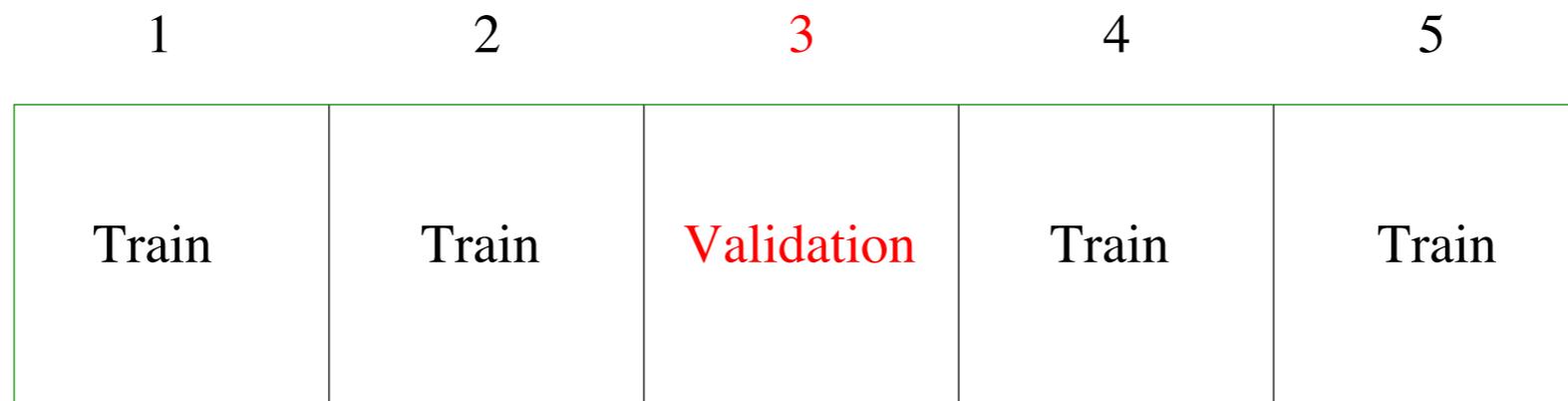


Figure: K-fold splitting of dataset (ESL 7.10)

- Leave-one-out cross validation: $K = N$

Cross-Validation: Model Selection

- Hyperparameter selection: α

$$CV(\hat{y}_\alpha) = \frac{1}{N} \sum_{i=1}^N E(\hat{y}_\alpha^{-k(i)}(\mathbf{x}_i), t)$$

$$\alpha^* = \underset{\alpha}{\operatorname{arg\,min}} CV(\hat{y}_\alpha)$$

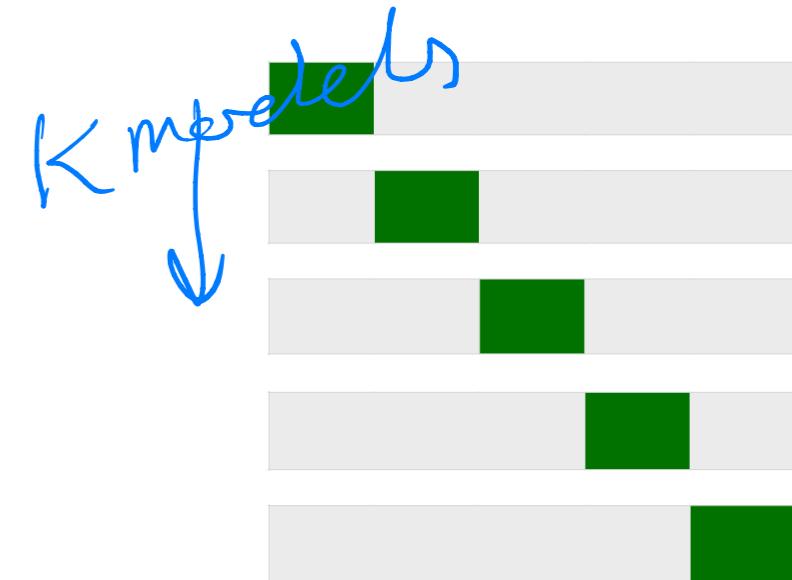
- Multiple hyperparameters: $\alpha \in \{\alpha_1, \alpha_2\}$, $\beta \in \{\beta_1, \beta_2, \beta_3\}$

- How many times should CV be performed?

$$2 \times 3 = 6$$

- Total number of training runs?

$$6 \times K$$



Nested Cross-Validation

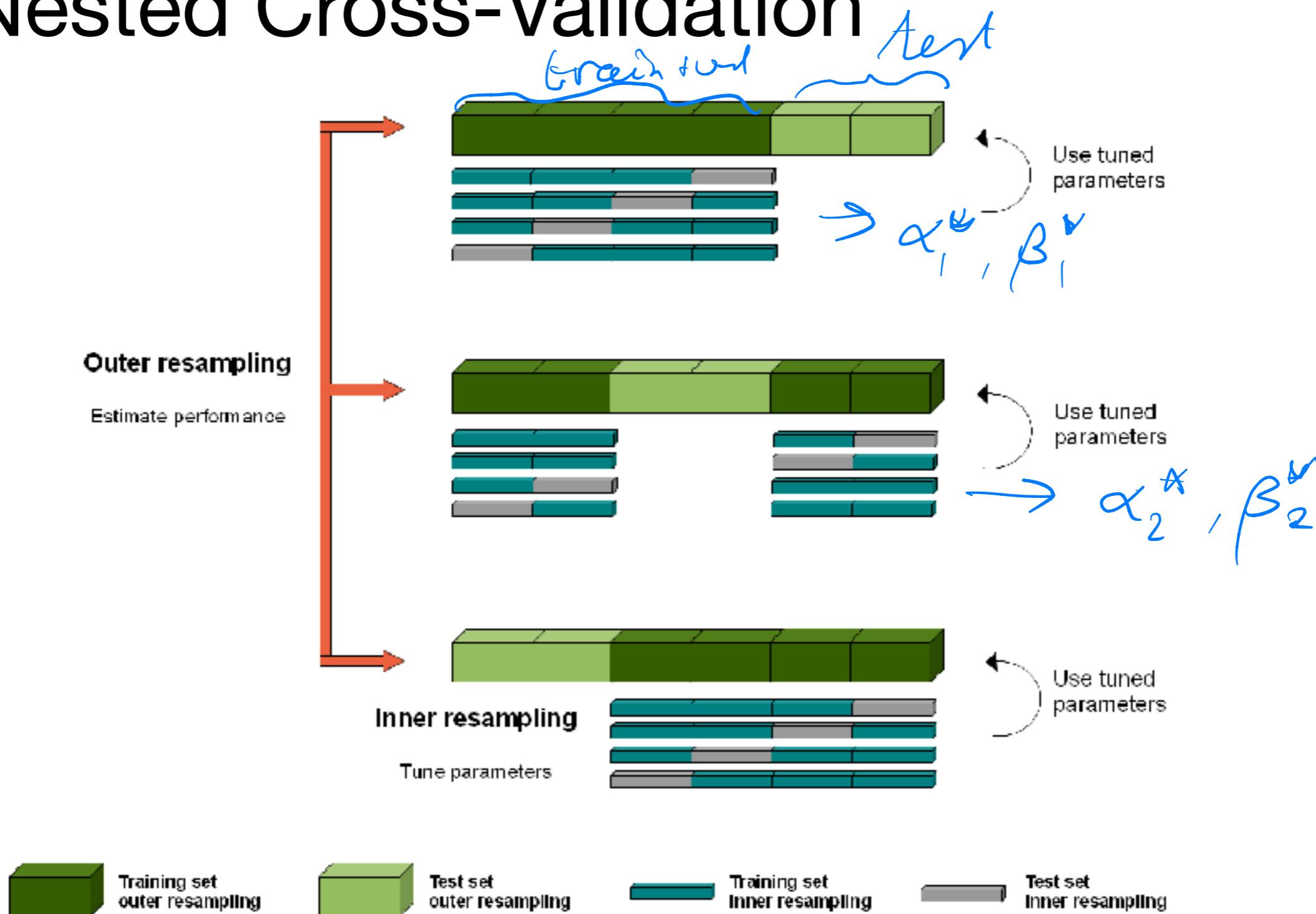


Figure: Nested cross-validation

https://mlr-org.github.io/mlr-tutorial-devel/html/nested_resampling/index.html
(site is offline unfortunately)

Machine Learning 1

Lecture 4.2 - Supervised Learning
Bias Variance Decomposition

Erik Bekkers

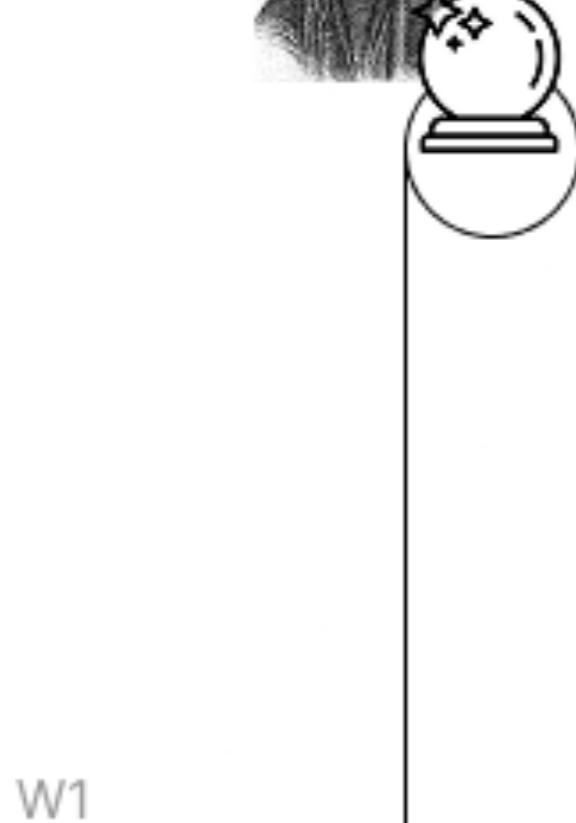
(Bishop 1.5.5, 3.2)



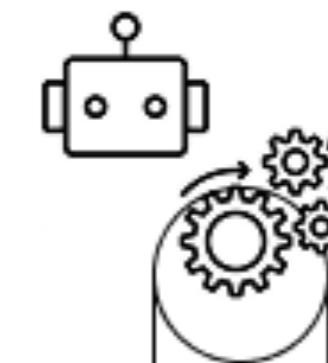


Machine Learning 1

The probabilistic view



The algorithmic view



Probabilistic modeling

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

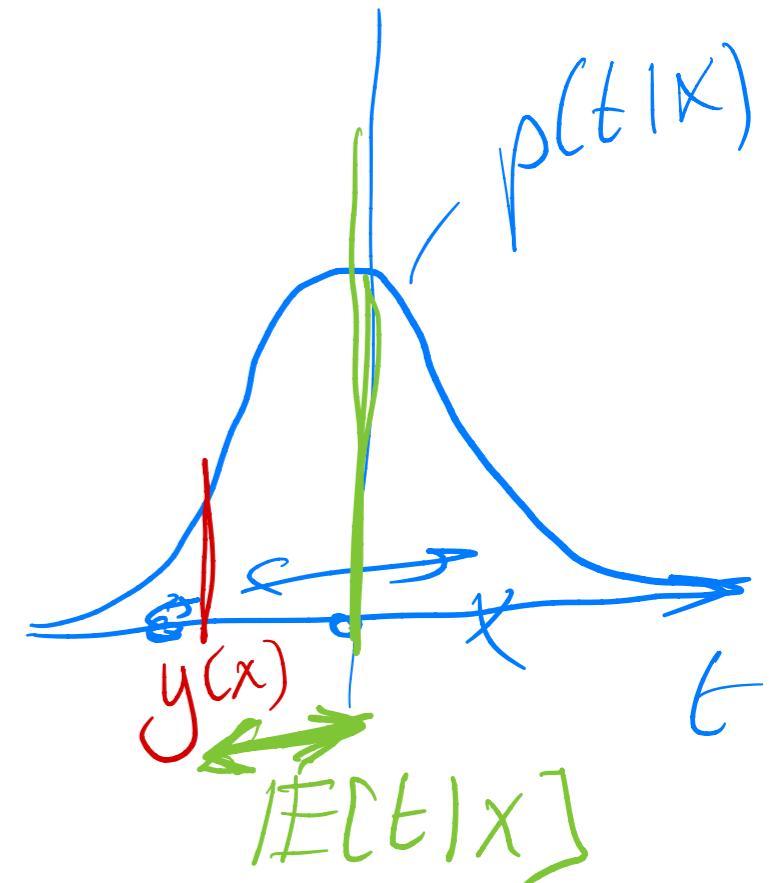
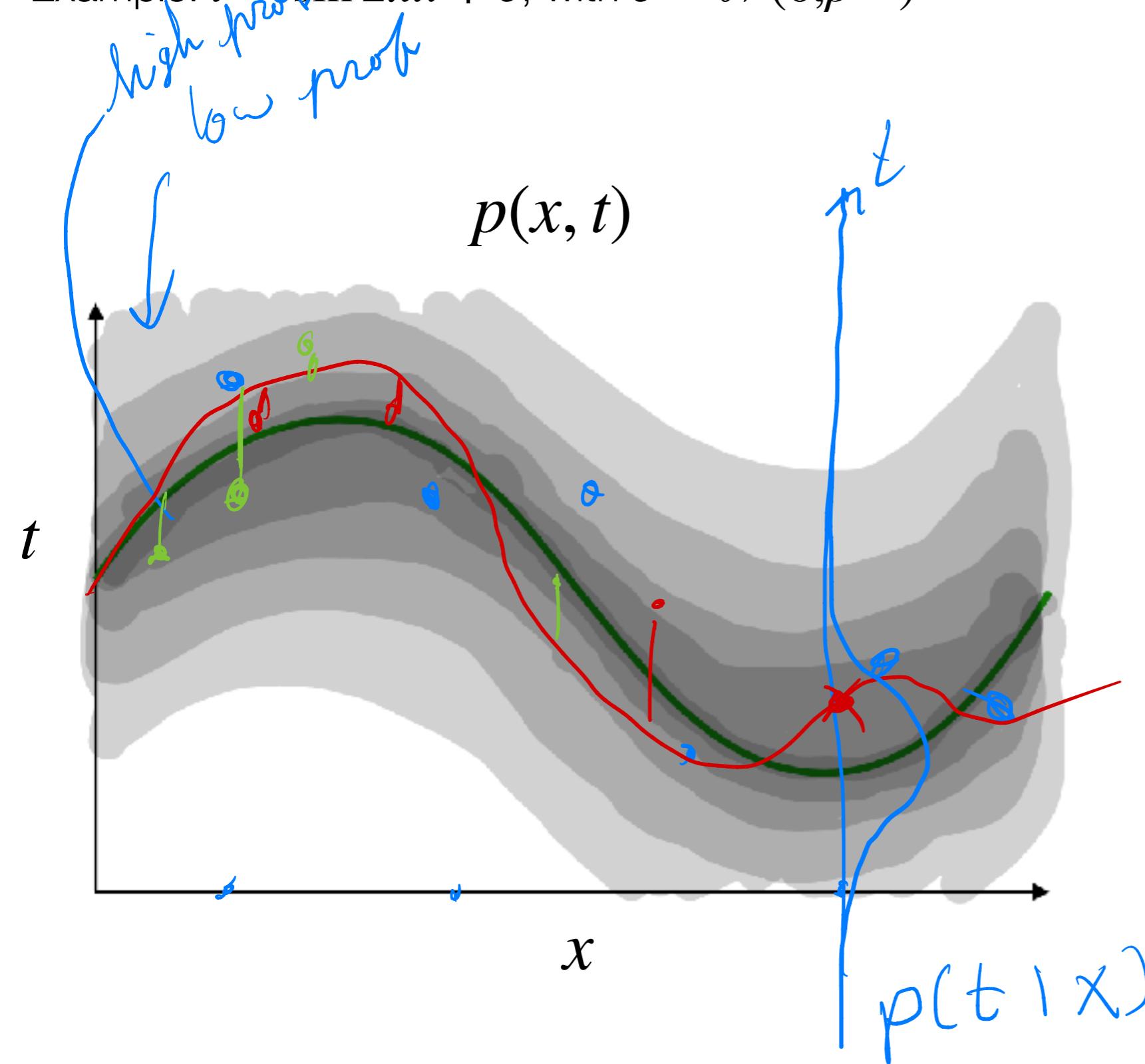
Optimization

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

Mathematical tools

Multi-variate calculus,
linear algebra

Example: $t = \sin 2\pi x + \epsilon$, with $\epsilon \sim \mathcal{N}(0, \beta^{-1})$



Expected Loss for Regression $\mathbb{E}_{(\mathbf{x},t) \sim p(\mathbf{x},t)}[L(t, y(\mathbf{x}))]$

- Consider the expected loss:

$$\mathbb{E}[L] = \int \int (y(\mathbf{x}) - t)^2 p(\mathbf{x}, t) dt d\mathbf{x}$$

- Let's analyze it relative to the regression function $\mathbb{E}[t | \mathbf{x}] := \mathbb{E}_{t \sim p(t|\mathbf{x})}[t | \mathbf{x}]$:

$$\mathbb{E}[L] = \int \int (y(\mathbf{x}) - \underbrace{\mathbb{E}[t | \mathbf{x}]}_a + \underbrace{\mathbb{E}[t | \mathbf{x}] - t}_b)^2 p(\mathbf{x}, t) dt d\mathbf{x}$$

Summary so far...

- Any model y will make errors, this expected loss decomposes into

$$\mathbb{E}_{(\mathbf{x}, t) \sim p(\mathbf{x}, t)}[L(t, y(\mathbf{x}))] = \int (y(\mathbf{x}) - \mathbb{E}[t | \mathbf{x}])^2 p(\mathbf{x}) d\mathbf{x} + \int \text{var}[t | \mathbf{x}] p(\mathbf{x}) d\mathbf{x}$$

zoom-in

- The best possible model is the regression function $y(\mathbf{x}) = \mathbb{E}[t | \mathbf{x}]$

$$D_1 = \{(x_1, t_1), \dots, (x_N, t_N)\} \rightarrow y_{D_1}$$

$$D_2 = \{(x_1, t_1), \dots, (x_N, t_N)\} \rightarrow y_{D_2}$$

- In practice we approximate it with fits $\hat{y}_D = \operatorname{argmin}_y \sum_{(\mathbf{x}, t) \in D} L(t, y(\mathbf{x}))$

- What can we say about the expected loss of \hat{y}_D ?

The Average Expected Loss

- Let's analyze performance of a learning algorithm by averaging the expected loss over learned y_D for different datasets D

$$\mathbb{E}_D[\mathbb{E}[L]] = \int \mathbb{E}_D[(y_D(\mathbf{x}) - \mathbb{E}[t | \mathbf{x}])^2] p(\mathbf{x}) d\mathbf{x} + \int \text{var}[t | \mathbf{x}] p(\mathbf{x}) d\mathbf{x}$$

The Average Expected Loss

- Let's analyze performance of a learning algorithm by averaging the expected loss over learned y_D for different datasets D

$$\mathbb{E}_D[\mathbb{E}[L]] = \int \mathbb{E}_D[(y_D(\mathbf{x}) - \mathbb{E}[t | \mathbf{x}])^2] p(\mathbf{x}) d\mathbf{x} + \int \text{var}[t | \mathbf{x}] p(\mathbf{x}) d\mathbf{x}$$

- Analyze it relative to the average model $\mathbb{E}_D[y_D(\mathbf{x})]$

$$\mathbb{E}_D[(y_D(\mathbf{x}) - \mathbb{E}[t | \mathbf{x}])^2] = \mathbb{E}_D[(y_D(\mathbf{x}) - \mathbb{E}_D[y_D(\mathbf{x})] + \mathbb{E}_D[y_D(\mathbf{x})] - \mathbb{E}[t | \mathbf{x}])^2]$$

The Average Expected Loss

- Let's analyze performance of a learning algorithm by averaging the expected loss over learned y_D for different datasets D

$$\mathbb{E}_D[\mathbb{E}[L]] = \int \mathbb{E}_D[(y_D(\mathbf{x}) - \mathbb{E}[t | \mathbf{x}])^2] p(\mathbf{x}) d\mathbf{x} + \int \text{var}[t | \mathbf{x}] p(\mathbf{x}) d\mathbf{x}$$

- Analyze it relative to the average model $\mathbb{E}_D[y_D(\mathbf{x})]$

$$\mathbb{E}_D[(y_D(\mathbf{x}) - \mathbb{E}[t | \mathbf{x}])^2] = \mathbb{E}_D[(y_D(\mathbf{x}) - \mathbb{E}_D[y_D(\mathbf{x})] + \mathbb{E}_D[y_D(\mathbf{x})] - \mathbb{E}[t | \mathbf{x}])^2]$$

(Expand square $(a + b)^2 = a^2 + 2ab + b^2$)

$$= \mathbb{E}_D[(y_D(\mathbf{x}) - \mathbb{E}_D[y_D(\mathbf{x})])^2] + \mathbb{E}_D[(\mathbb{E}_D[y_D(\mathbf{x})] - \mathbb{E}[t | \mathbf{x}])^2]$$

$$+ 2 \mathbb{E}_D[(y_D(\mathbf{x}) - \mathbb{E}_D[y_D(\mathbf{x})])(\mathbb{E}_D[y_D(\mathbf{x})] - \mathbb{E}[t | \mathbf{x}])]$$

Bias-Variance Decomposition

- ▶ What can we say about the expected loss of y_D ?
- ▶ On average (over datasets D) our model y_D will make three types of errors:

$$\mathbb{E}_D[\mathbb{E}[L]] = \int \mathbb{E}_D[(y_D(x) - \mathbb{E}_D[y_D(x)])^2] p(x) dx \quad \text{Variance}$$

$$+ \int (\mathbb{E}_D[y_D(x)] - \mathbb{E}[t|x])^2 p(x) dx \quad \text{Bias}^2$$

$$+ \int \text{var}[t|x] p(x) dx \quad \text{Noise}$$

Bias-Variance Decomposition

- ▶ What can we say about the expected loss of y_D ?
- ▶ On average (over datasets D) our model y_D will make three types of errors:

$$\mathbb{E}_D[\mathbb{E}[L]] \approx \frac{1}{N} \sum_{n=1}^N \frac{1}{L} \sum_{l=1}^L (y^{(l)}(x_n) - \bar{y}(x_n))^2 \quad \text{Variance}$$

$$+ \frac{1}{N} \sum_{i=1}^N \bar{y}(x_n) - \sin 2\pi x_n)^2 \quad \text{Bias}^2$$

see lecture
for finite dataset
approximation

Bias-Variance Decomposition: Example

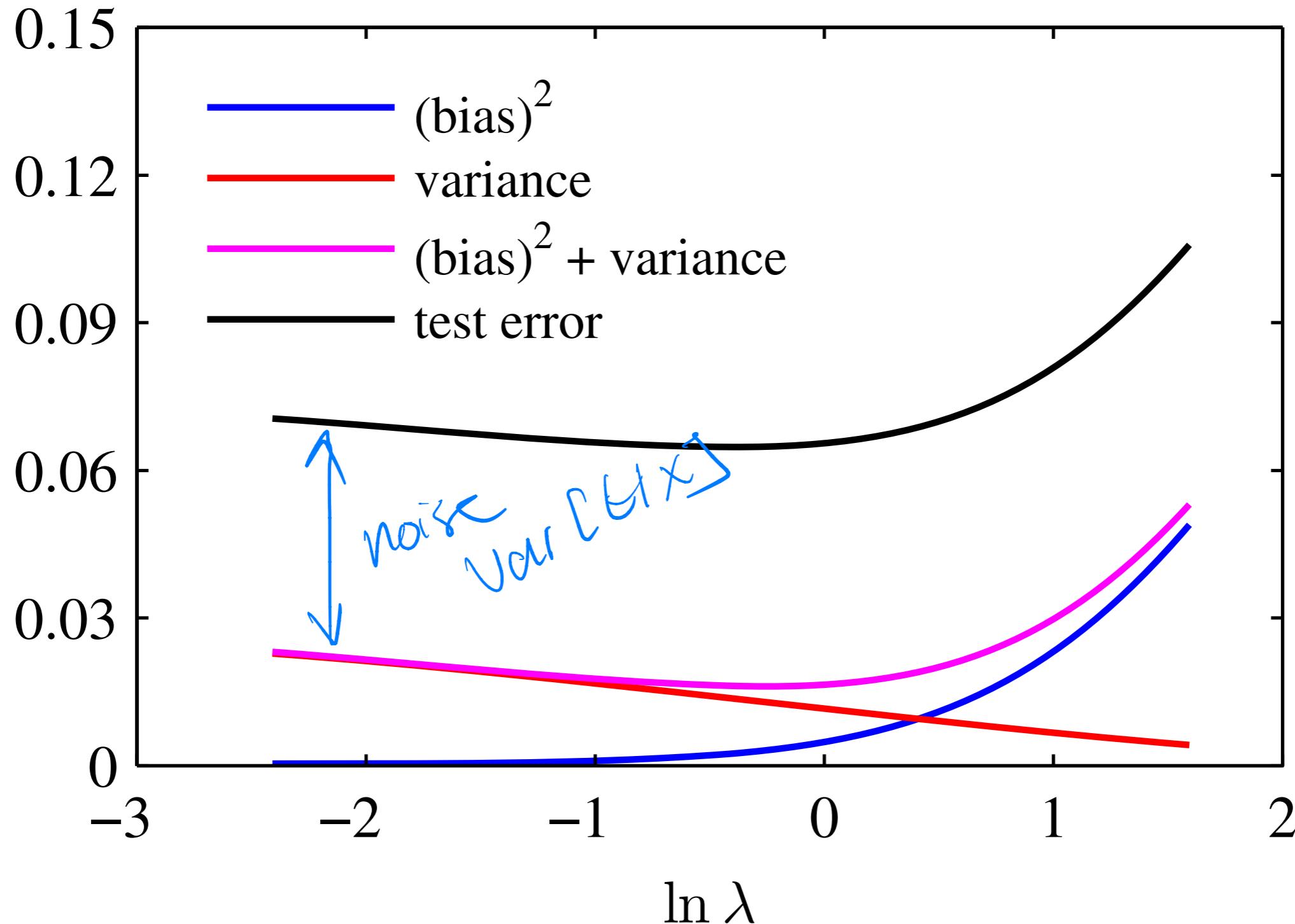


Figure: bias-variance decomposition (Bishop 3.6)

Machine Learning 1

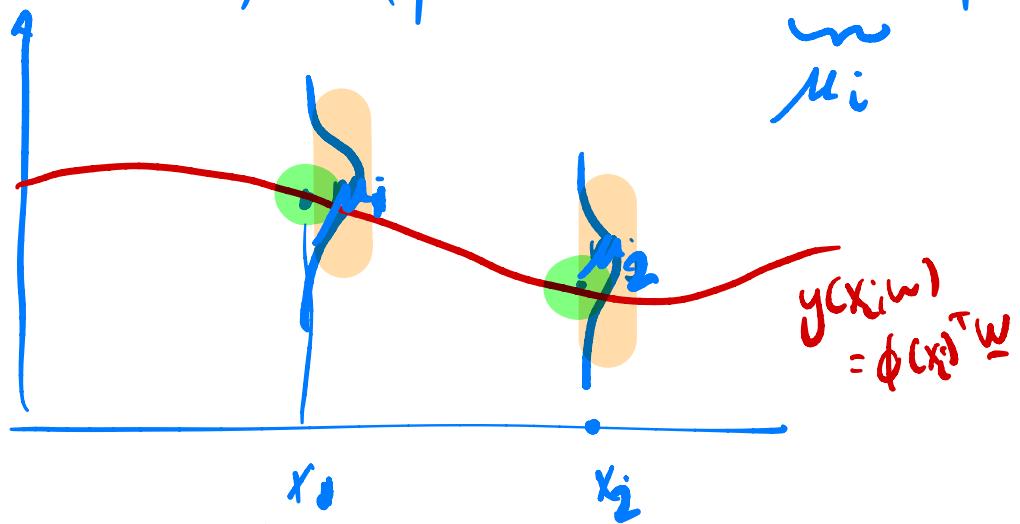
Lecture 4.3 - Supervised Learning
Bayesian Linear Regression - **Gaussian
Posteriors**

Erik Bekkers

(Bishop 3.3.1 (and 2.3.3))



$$(4) P(t_i | x_i, \underline{w}, \beta) = N(b_i | \underline{\phi}_i^T \underline{w}, \beta)$$



$$\underline{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_N \end{pmatrix} = \begin{pmatrix} \underline{\phi}(x_1)^T w \\ \underline{\phi}(x_2)^T w \\ \vdots \\ \underline{\phi}(x_N)^T w \end{pmatrix} = \begin{pmatrix} \phi_1(x_1) & \phi_2(x_1) & \dots \\ \vdots & \vdots & \vdots \\ \phi_1(x_N) & \phi_2(x_N) & \dots \end{pmatrix} w$$

$= \Phi \underline{w}$ predict all means at once

$$\Sigma = \beta^{-1} I_N = \begin{pmatrix} \beta^{-1} & 0 & \dots \\ 0 & \beta^{-1} & \dots \\ \vdots & \vdots & \ddots \end{pmatrix}$$

Since independent and shared β

$$(*) \prod_{i=1}^N e^{-\frac{\beta}{2} (\underline{t}_i - \underline{\phi}_i^\top \underline{w})^2}$$

$$= e^{-\frac{\beta}{2} \sum_{i=1}^N (\underline{t}_i - \underline{\phi}_i^\top \underline{w})^2}$$

$$= e^{-\frac{\beta}{2} (\underline{t} - \underline{\Phi} \underline{w})^\top (\underline{t} - \underline{\Phi} \underline{w})}$$

$$= e$$

$$= e^{-\frac{\beta}{2} (\underline{t} - \underline{\Phi} \underline{w})^\top (\underline{t} - \underline{\Phi} \underline{w})}$$

$$\text{so } \underline{\mu} = \underline{\Phi} \underline{w} \quad \in \mathbb{R}^N$$

$$\Sigma = \underline{\Phi}^{-1} \underline{I}_N \quad \in \mathbb{R}^{N \times N}$$

Machine Learning 1

Lecture 4.4 - Supervised Learning
Bayesian Linear Regression - **Sequential
Bayesian Learning**

Erik Bekkers

(Bishop 3.3.1)



Machine Learning 1

Lecture 4.4 - Supervised Learning
Bayesian Linear Regression - **Sequential
Bayesian Learning**

Erik Bekkers

(Bishop 3.3.1)



Machine Learning 1

Lecture 4.5 - Supervised Learning
Bayesian Linear Regression - **Predictive
Distribution**

Erik Bekkers

(Bishop 3.3.2)

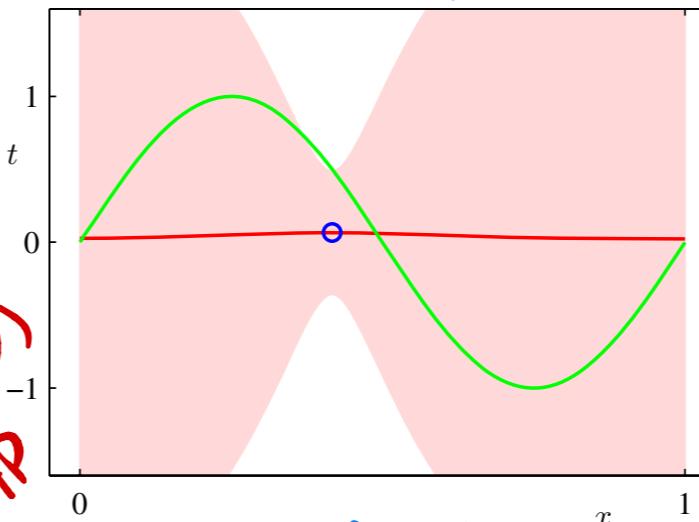


Predictive Distribution

active learning?

$N=1$

- ▶ Datasets:
 - ▶ $t = \sin(2\pi x) + \epsilon$
 - ▶ $\epsilon \sim \mathcal{N}(0, \beta^{-1})$



- ▶ Dataset sizes:

▶ $N = 1, 2, 4, 25$

- ▶ Model:
 - ▶ $y(x, \mathbf{w}) = \boldsymbol{\phi}(x)^T \mathbf{w}$
 - ▶ $\boldsymbol{\phi}_j(x)$: Gaussian basis functions

mean
given by $y(x, \mathbf{w}_{MAP})$

$N=4$

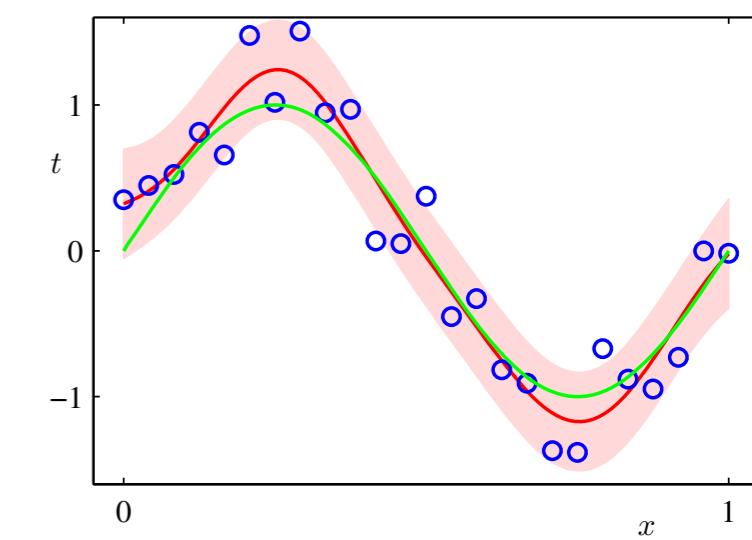
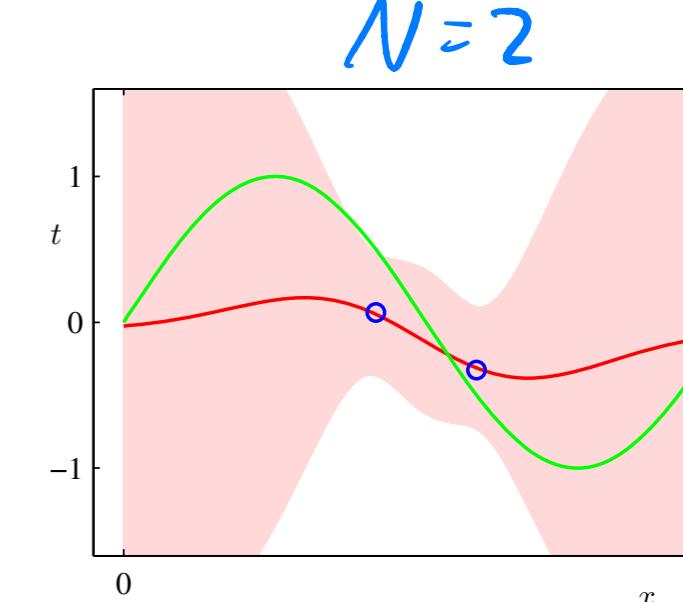
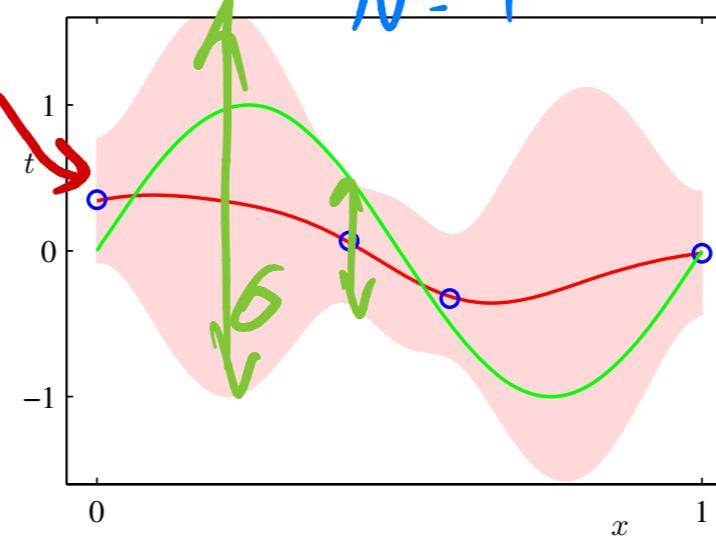


Figure: Predictive distribution (Bishop 3.8)

- ▶ Predictive distribution:

$$p(t' | x', \mathbf{X}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(t' | \mathbf{x}', \boldsymbol{\phi}(\mathbf{x}')^T \mathbf{m}_N, \sigma_N^2(\mathbf{x}'))$$

$y(x, \mathbf{w}_{MAP})$

$$\sigma_N^2(\mathbf{x}') = \beta^{-1} + \boldsymbol{\phi}(\mathbf{x}')^T \mathbf{S}_N \boldsymbol{\phi}(\mathbf{x}'), \quad \mathbf{m}_N = \beta \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t}, \quad \mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}$$

Samples drawn from Bayesian Predictive Distribution

$$\underline{w} \sim p(\mathbf{w} | \mathbf{X}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N),$$

$$\text{with } \mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{t}$$
$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi$$

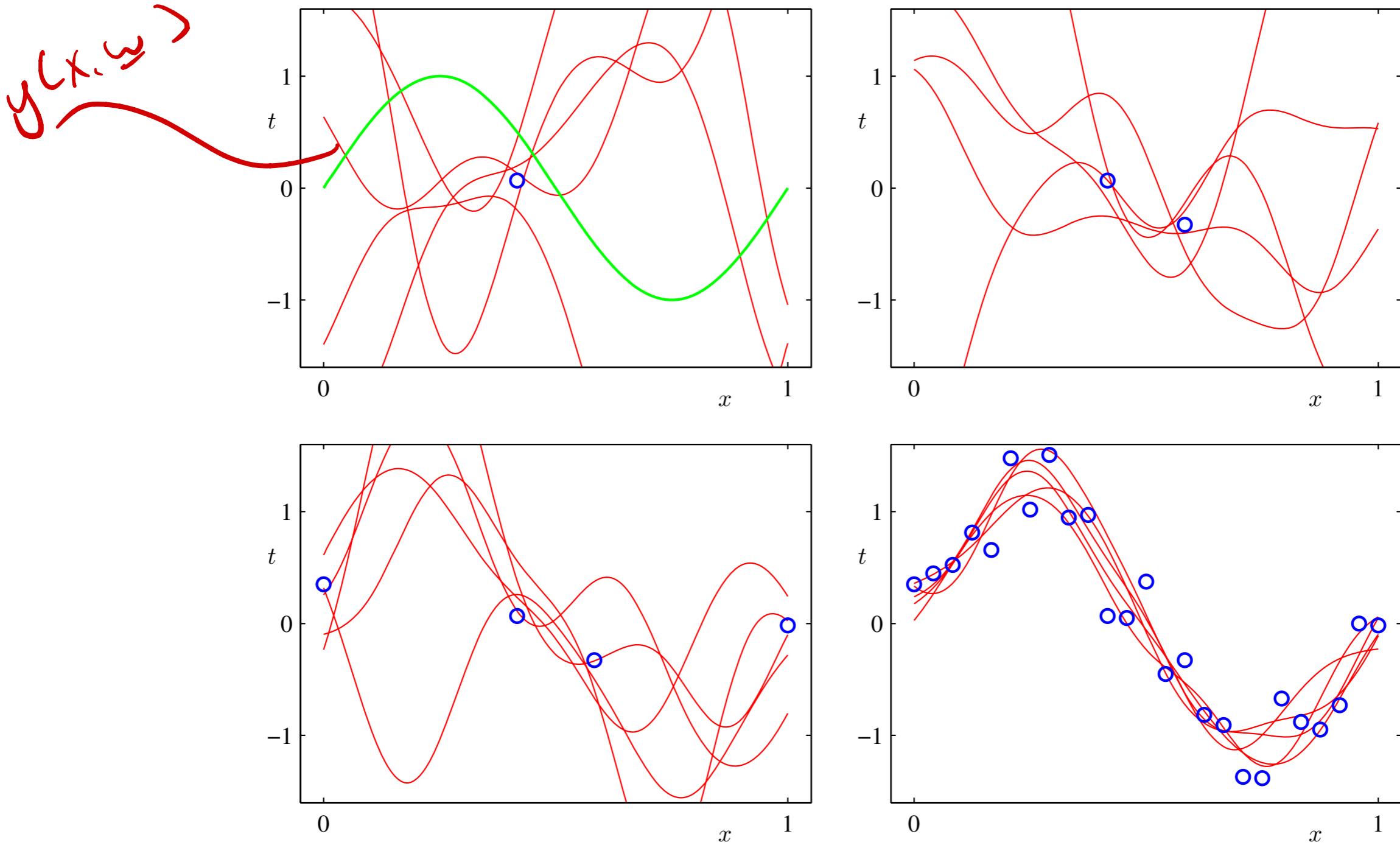


Figure: Sample functions $y(x, \mathbf{w})$ with \mathbf{w} sampled from posterior distribution (Bishop 3.9)