



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

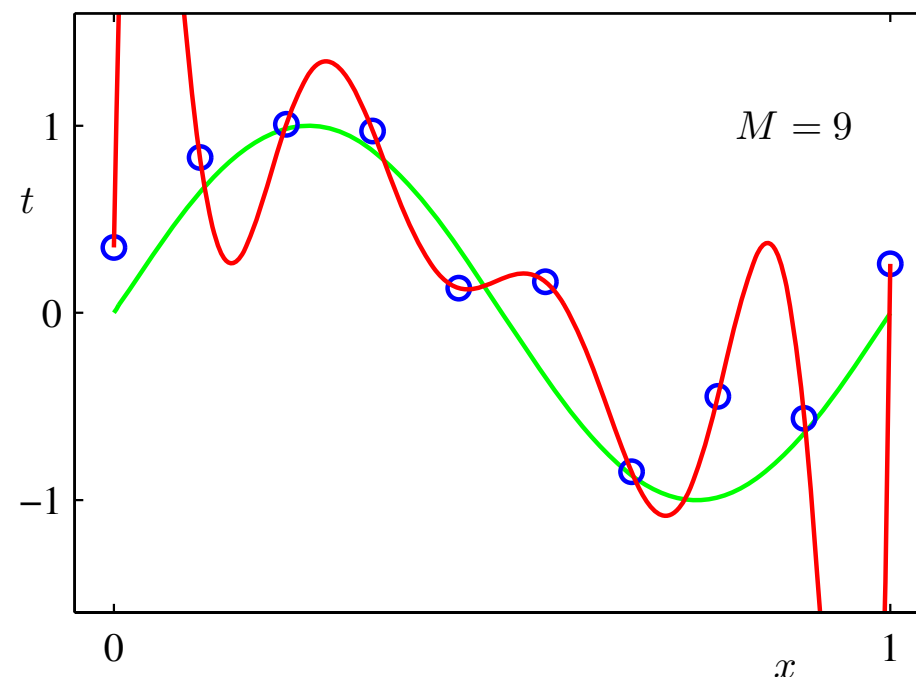
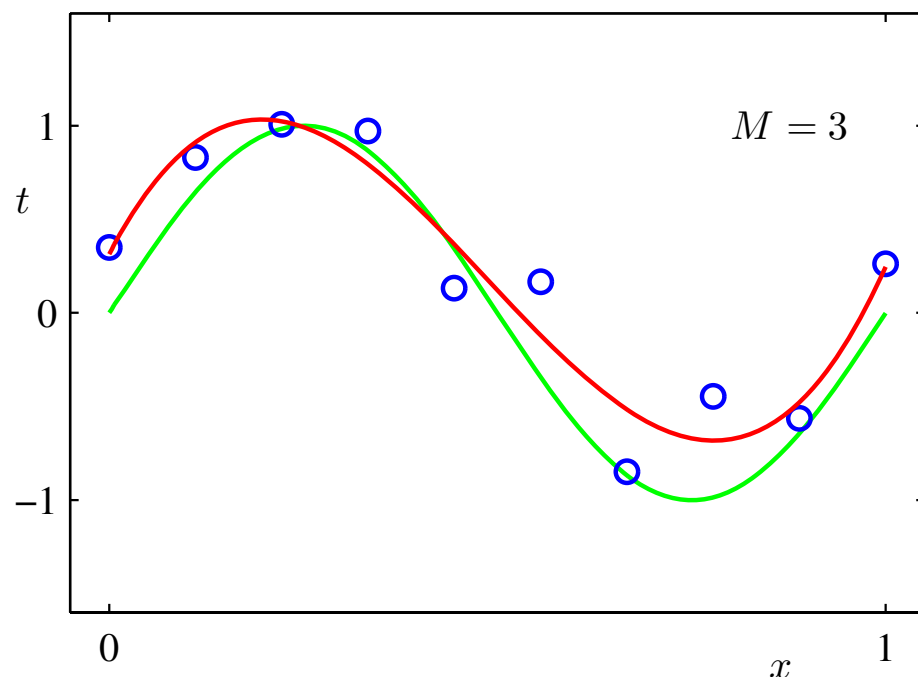
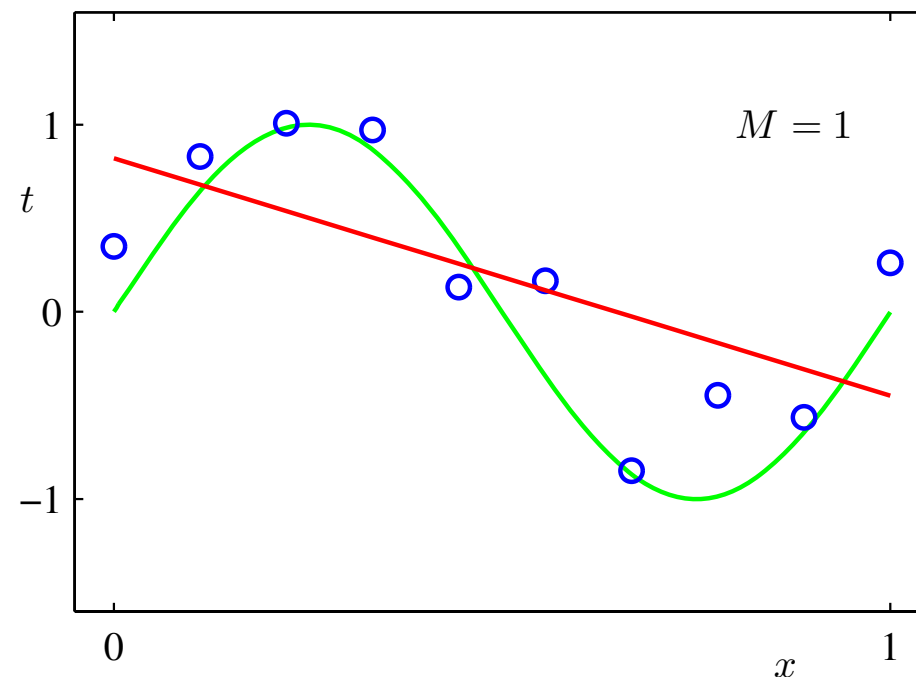
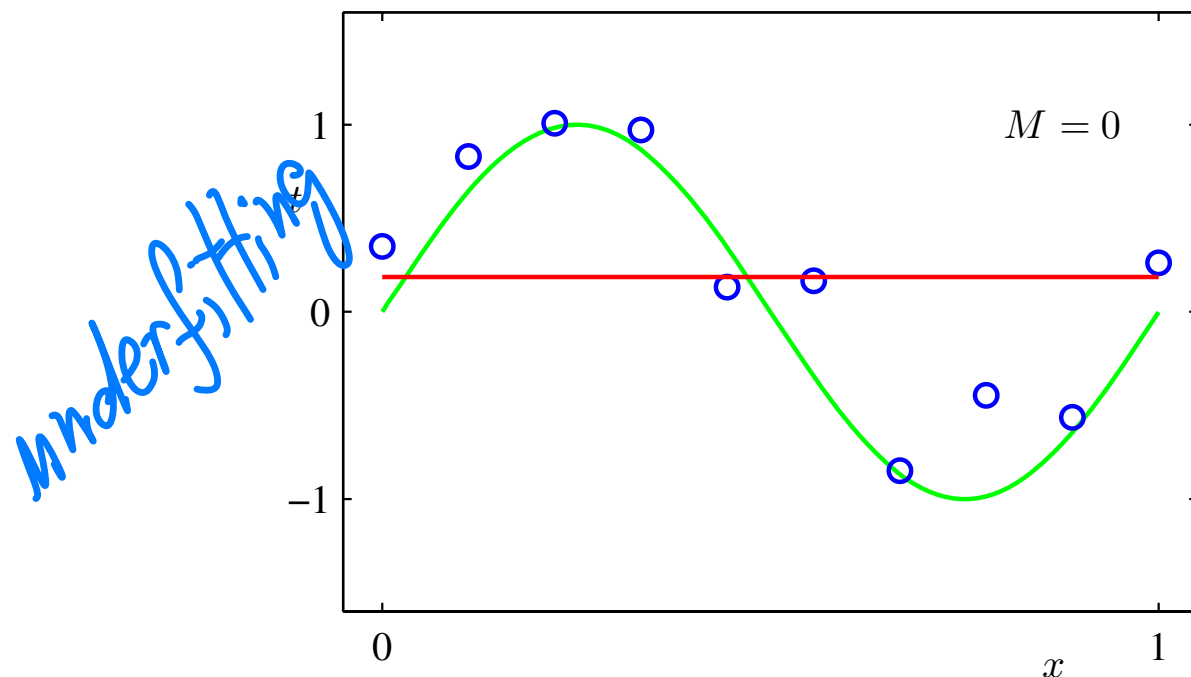
(Bishop 3.1.4)



Example: Overfitting and Underfitting

$$t = \sin(2\pi x) + \varepsilon \quad \varepsilon \sim \mathcal{N}(0, \beta^{-1})$$

$$y(x, \mathbf{w}) = w_0 + \sum_{i=1}^M w_i x^i$$



overfitting

Figure: Fits of different polynomials (Bishop 1.4)

Example: Overfitting (M=9)

| | $M = 0$ | $M = 1$ | $M = 3$ | $M = 9$ |
|---------|---------|---------|---------|-------------|
| w_0^* | 0.19 | 0.82 | 0.31 | 0.35 |
| w_1^* | | -1.27 | 7.99 | 232.37 |
| w_2^* | | | -25.43 | -5321.83 |
| w_3^* | | | 17.37 | 48568.31 |
| w_4^* | | | | -231639.30 |
| w_5^* | | | | 640042.26 |
| w_6^* | | | | -1061800.52 |
| w_7^* | | | | 1042400.18 |
| w_8^* | | | | -557682.99 |
| w_9^* | | | | 125201.43 |

Table: Polynomial coefficients (Bishop 1.1)

Regularized Least Squares

- ▶ Instead of manually constraining the number of parameters for small datasets, **add penalty term** for large parameter values:

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \{t_i - y(\mathbf{x}_i, \mathbf{w})\}^2 + \frac{\lambda}{2} \sum_{i=1}^M w_i^2$$

- Ridge regression
- ℓ_2 regularization
- weight decay

- ▶ The bias term w_0 is often not included in regularization

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**:

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

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

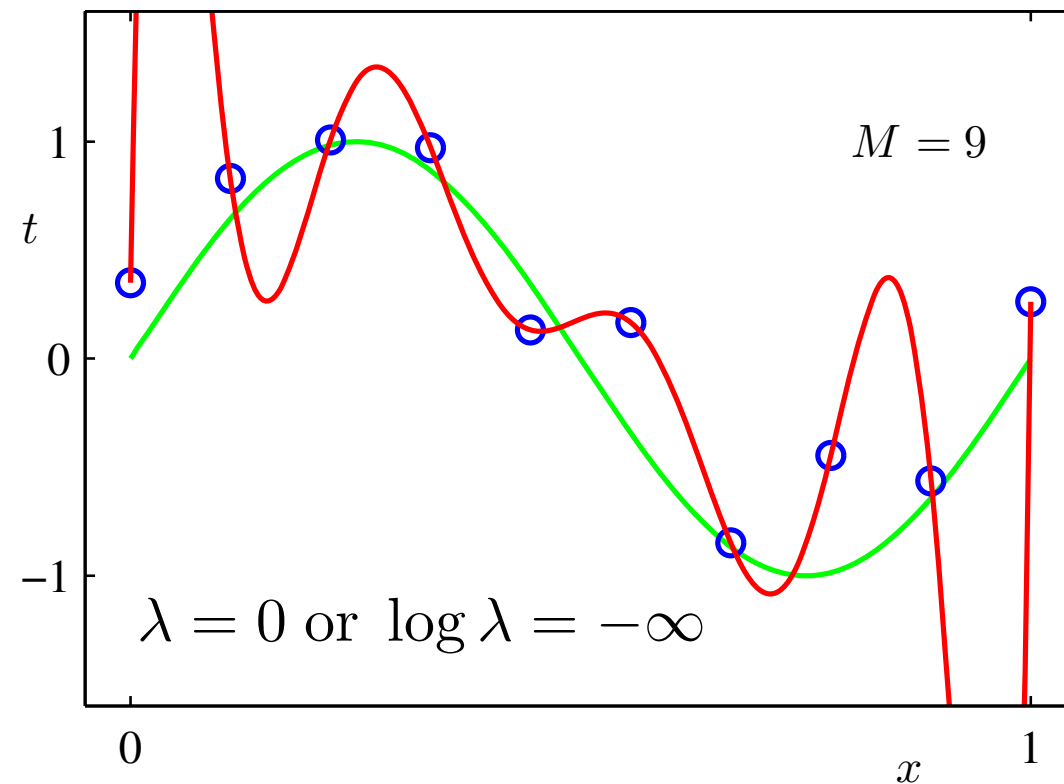
$$p(\mathbf{w} | \alpha) = \prod_{i=1}^M \mathcal{N}(w_i | 0, \alpha^{-1})$$

When $\lambda = \frac{\alpha}{\beta}$
→ Ridge

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

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

Example: Regularized Polynomial Regression



no regularization

good λ

Figure: polynomial regression (Bishop 1.4)

large λ

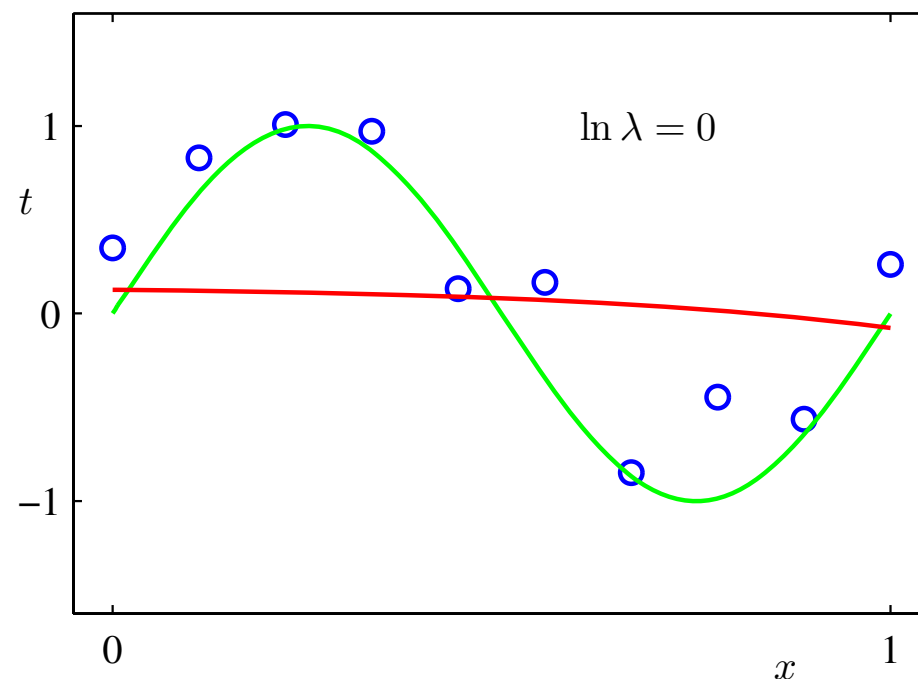
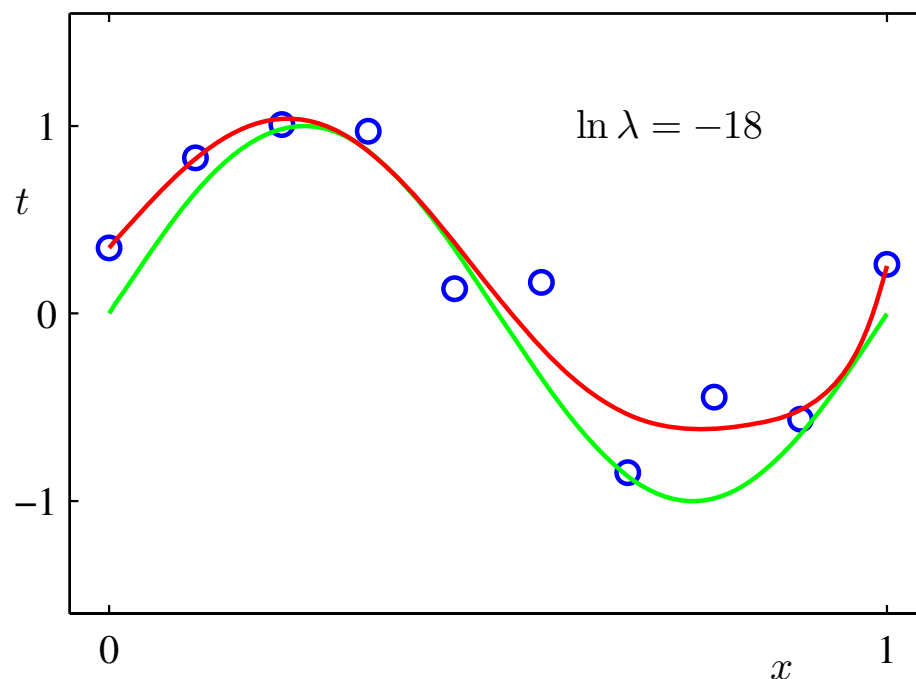


Figure: Regularized polynomial regression (Bishop 1.7)

Example: Regularized Polynomial Regression

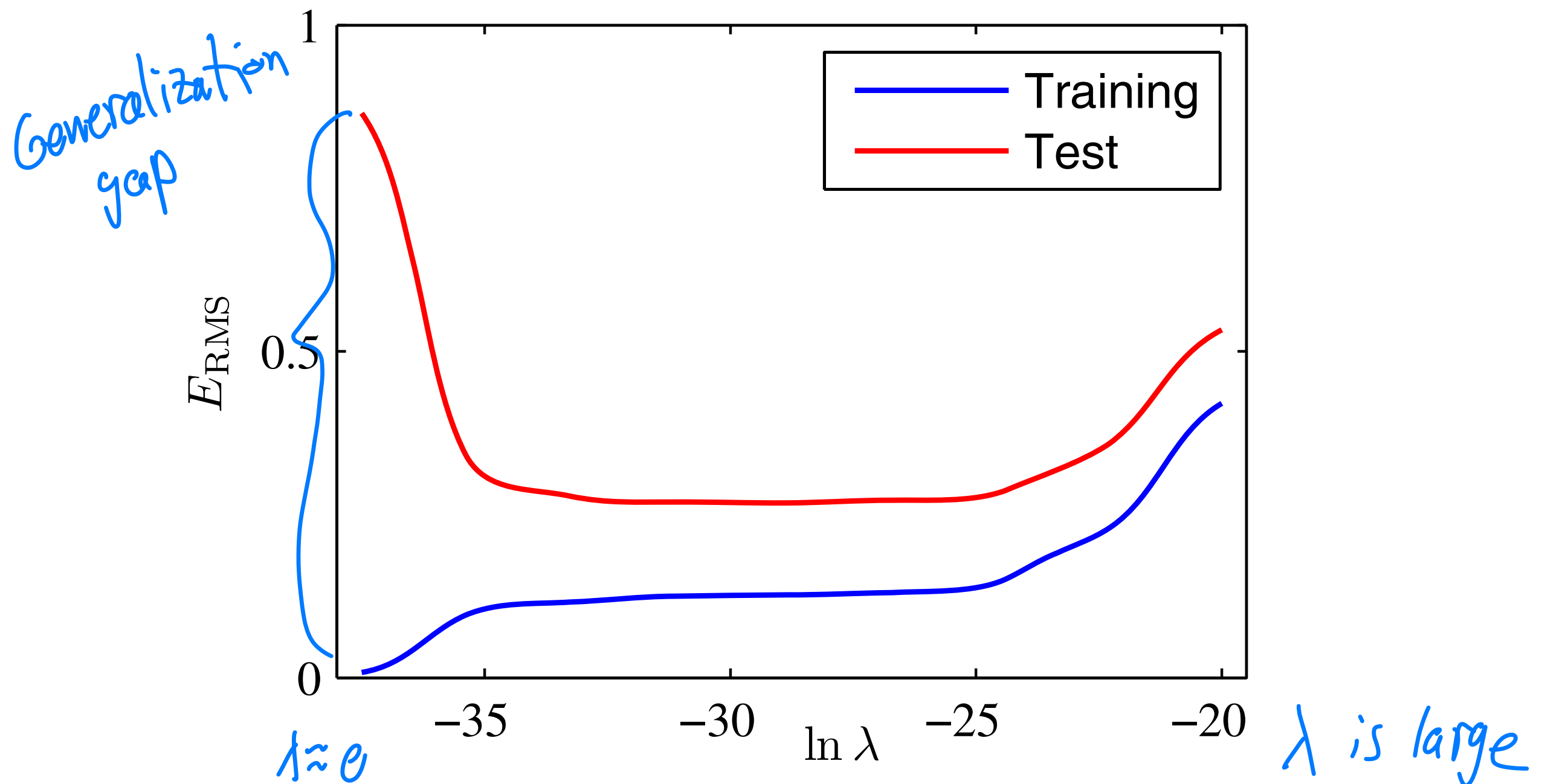


Figure: train and test errors for regularized $M=9$ polynomial regression (Bishop 1.8)

Regularized Least Squares (II)

Weight decay: $\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \{t_i - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i)\}^2 + \frac{\lambda}{2} \sum_{i=1}^{M-1} |w_i|^2$

General penalty:

$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{i=1}^N \{t_i - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_i)\}^2 + \frac{\lambda}{2} \sum_{i=1}^{M-1} |w_i|^q$$

Case $q = 1$: Lasso

$+\frac{\lambda}{2} \sum_{i=1} |w_i|$
- sparsification

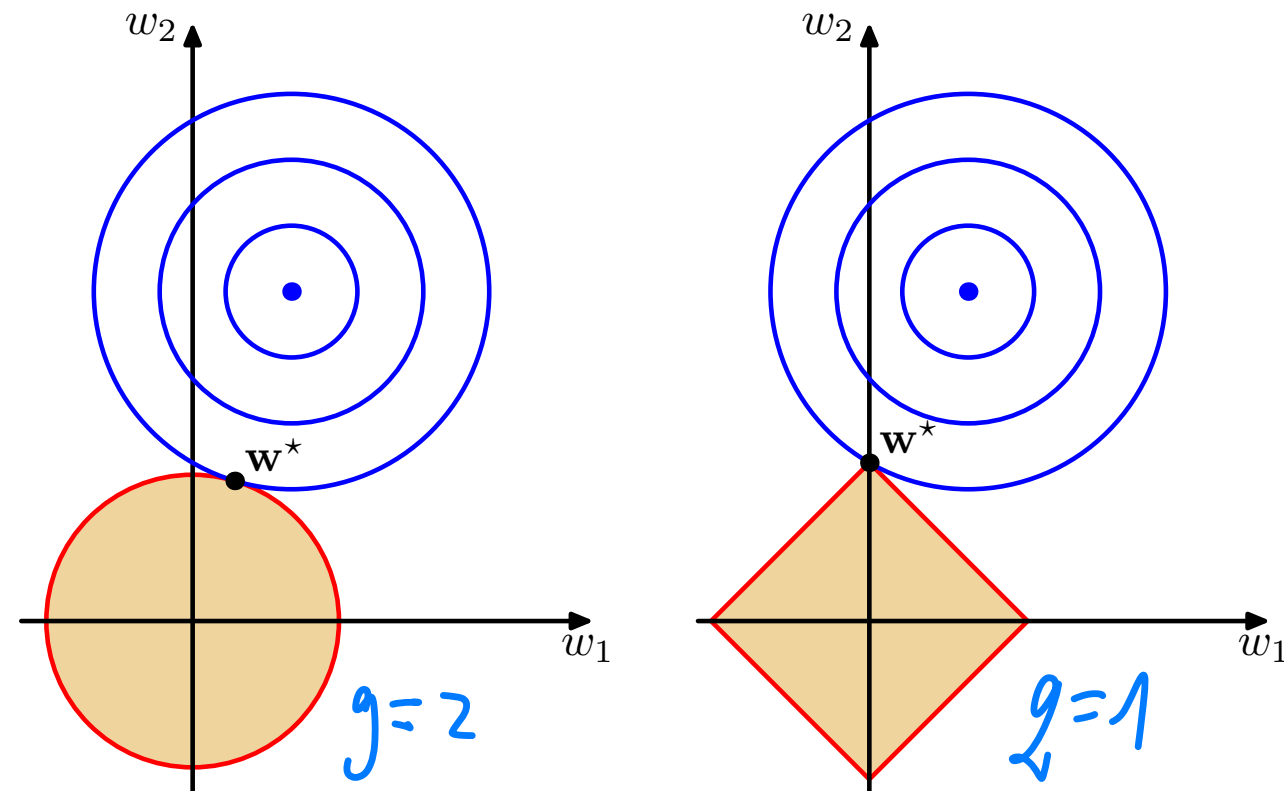


Figure: regularization as constrained optimization (Bishop 3.4)

Equivalent to some constraint optimization problem (Bishop App. E)

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

Regularized Least Squares: sparse weights

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

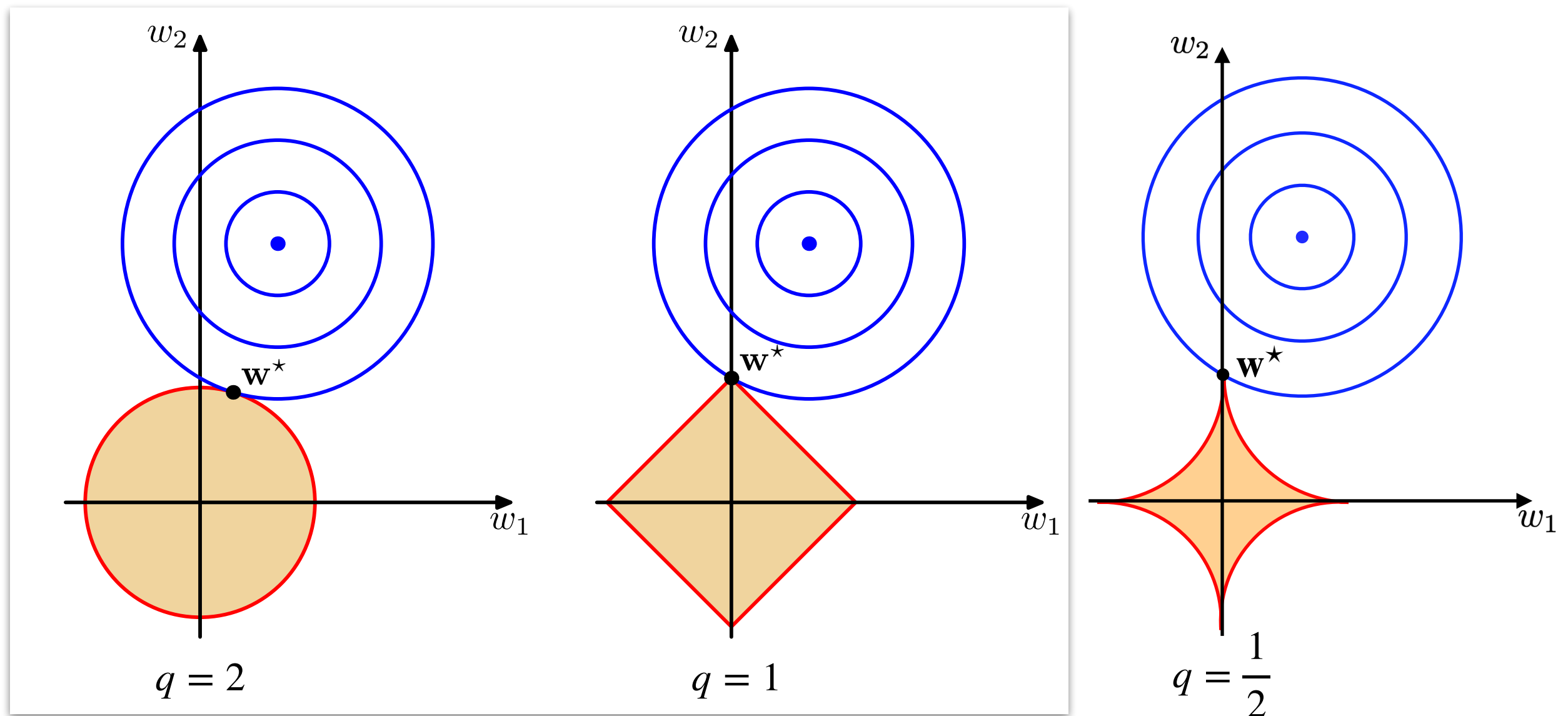


Figure: regularization as constrained optimization (Bishop 3.4)

Example: Prostate specific antigen prediction

$q=2$ (Ridge regression)

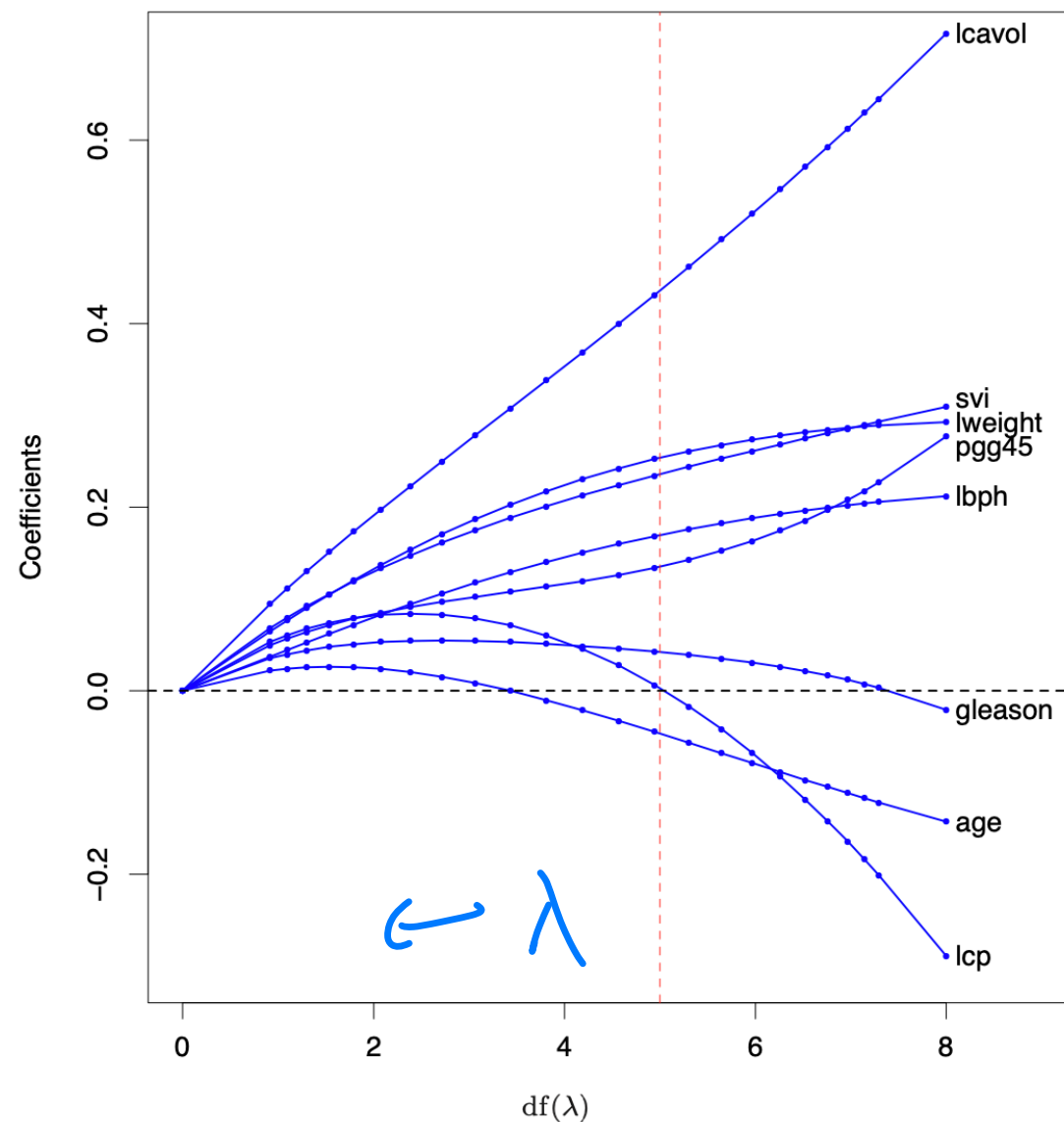


FIGURE 3.8. Profiles of ridge coefficients for the prostate cancer example, as the tuning parameter λ is varied. Coefficients are plotted versus $df(\lambda)$, the effective degrees of freedom. A vertical line is drawn at $df = 5.0$, the value chosen by cross-validation.

$q=1$ (Lasso regression)

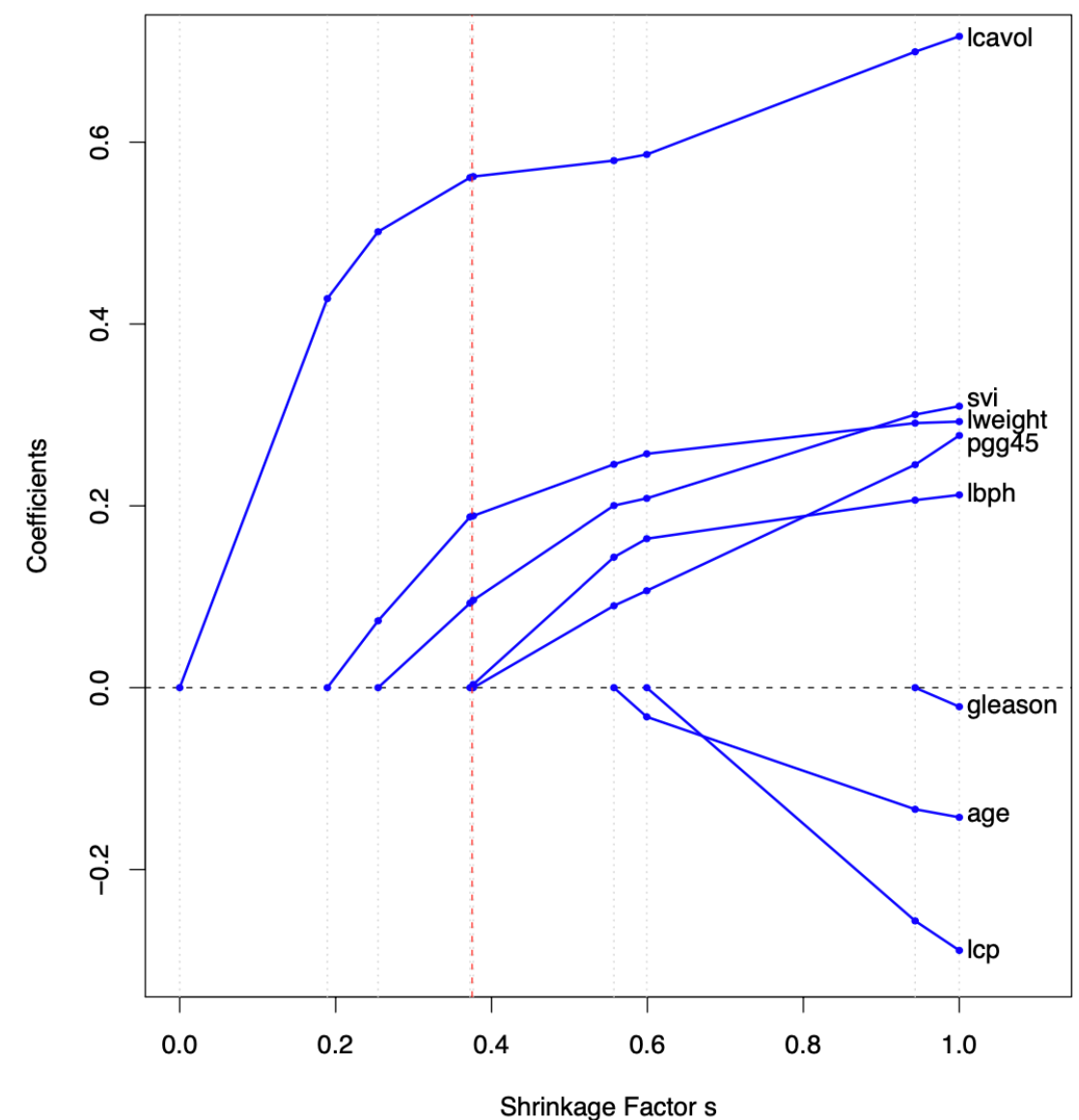


FIGURE 3.10. Profiles of lasso coefficients, as the tuning parameter t is varied. Coefficients are plotted versus $s = t / \sum_1^p |\hat{\beta}_j|$. A vertical line is drawn at $s = 0.36$, the value chosen by cross-validation. Compare Figure 3.8 on page 65; the lasso profiles hit zero, while those for ridge do not. The profiles are piece-wise linear, and so are computed only at the points displayed; see Section 3.4.4 for details.

Figures from the Elements of Statistical Learning (ESL - Hastie et al.)



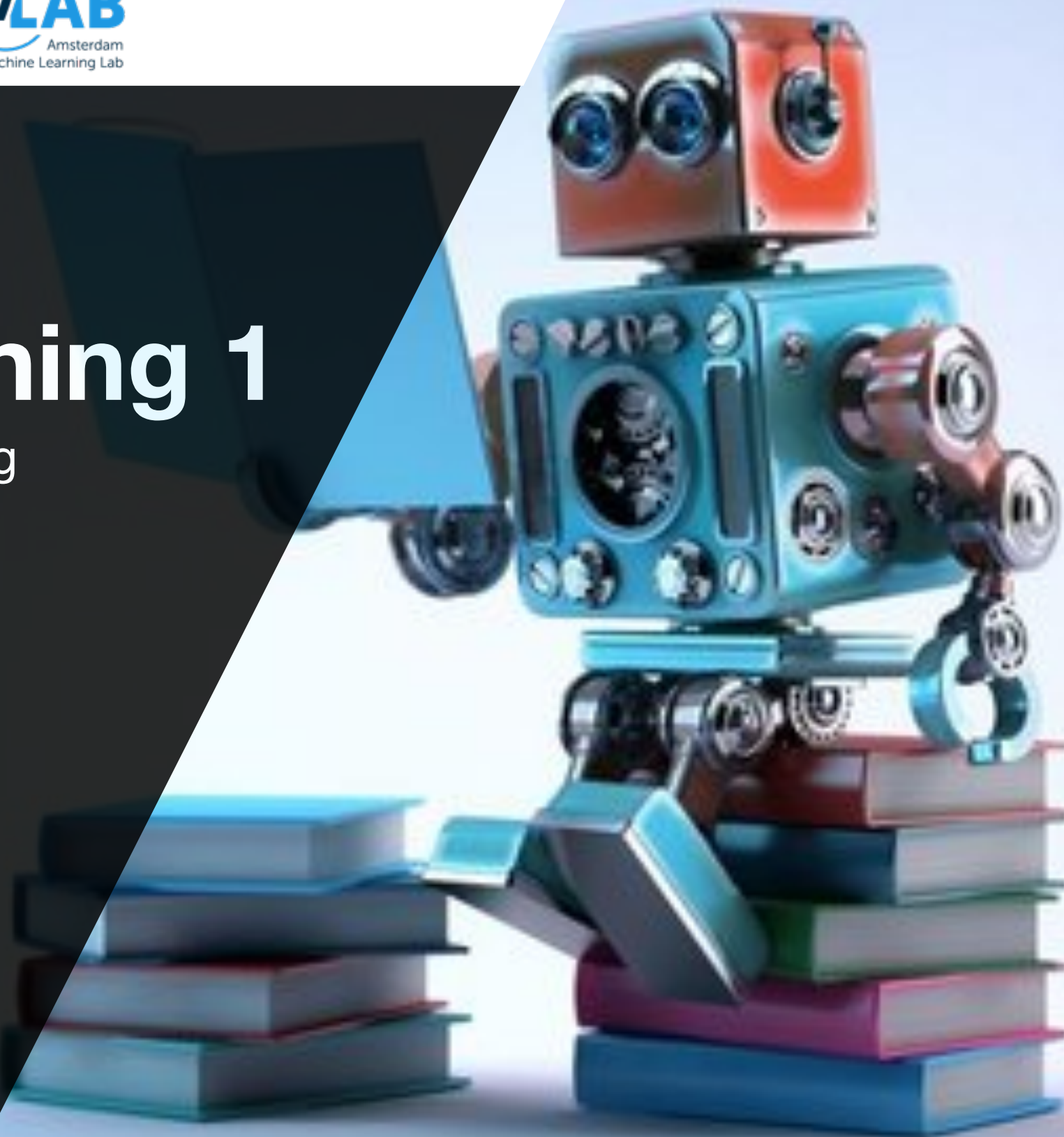
Machine Learning 1

Lecture 4.1 - Supervised Learning

Model Selection

Erik Bekkers

(Bishop 1.3)



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?

Model selection | Dataset splits

Divide data $D = \{(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N)\}$ in 3 groups:

- ▶ **Training** set D_{train} (± 80 % of D):

- ▶ Minimize the error $E(y(\mathbf{x}, \mathbf{w}), t)$
for every $(\mathbf{x}, t) \in D_{\text{train}}$

} obtain \mathbf{w}^*

- ▶ **Validation** set D_{val} (± 10 % of D):

- ▶ Used to *estimate* generalization error
for every $(\mathbf{x}_{\text{val}}, t_{\text{val}}) \in D_{\text{val}}$

$E(y(\mathbf{x}_{\text{val}}, \mathbf{w}^*), t_{\text{val}})$

} model selection
estimate performance

- ▶ **Test** set D_{test} (± 10 % of D):

- ▶ final test/generalization error estimate $E(y(\mathbf{x}_{\text{test}}, \mathbf{w}^*), t_{\text{test}})$
for every $(\mathbf{x}_{\text{test}}, t_{\text{test}}) \in D_{\text{test}}$

cannot be used for model selection

Supervised Learning: Small Datasets

- ▶ Small dataset → 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) → \hat{y}^{-k}

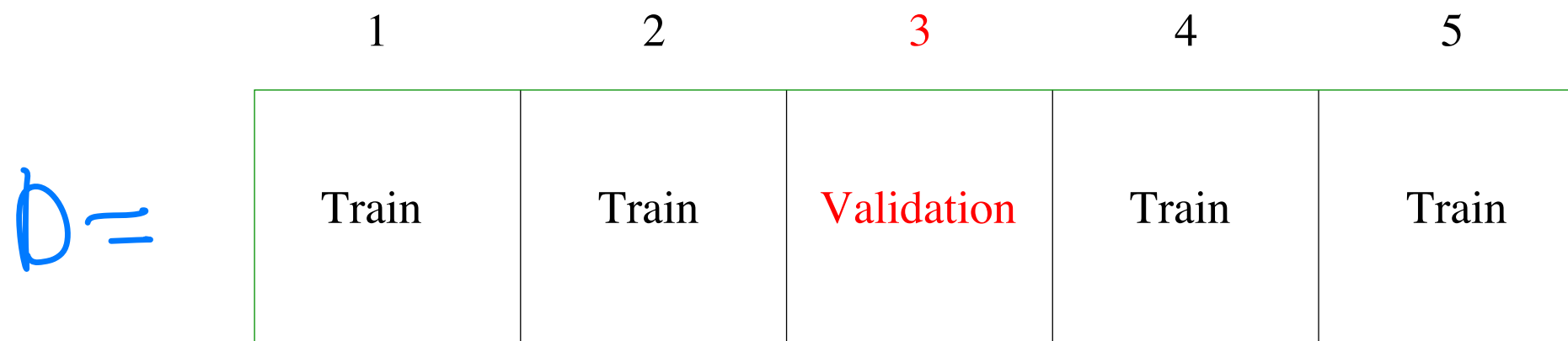


Figure: K-fold splitting of dataset (ESL 7.10)

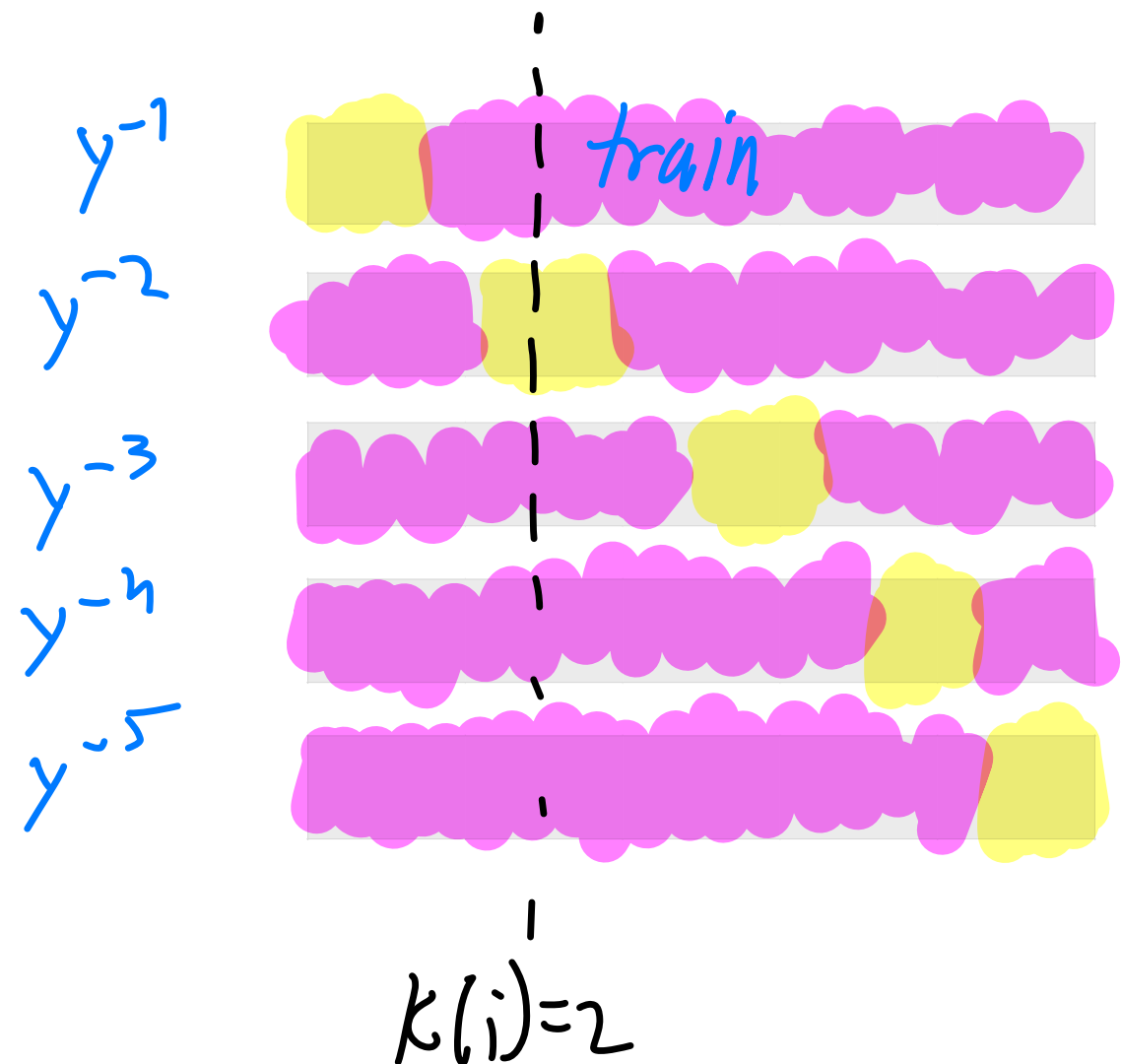
- ▶ Leave-one-out cross validation: $K = N$

Cross-Validation

- ▶ K trained functions \hat{y}^{-k}
- ▶ Use indexing function $\kappa : \{1, \dots, N\} \rightarrow \{1, \dots, K\}$
- ▶ Cross validation error:

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

model selection



Cross-Validation: Model Selection

- Hyperparameter selection:

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

- $$\alpha^* = \underset{\alpha}{\operatorname{argmax}} CV(\hat{y}_\alpha)$$

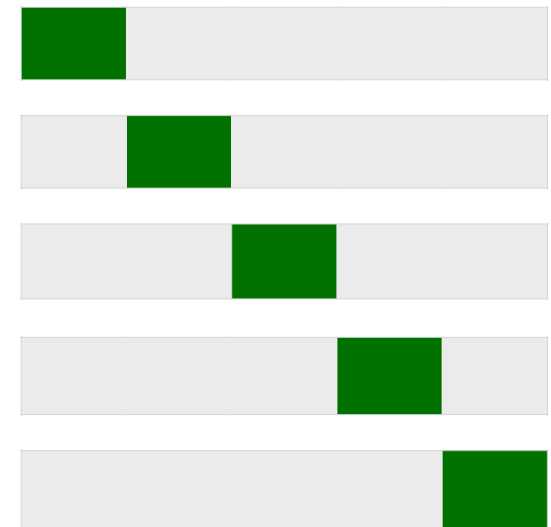
- Multiple hyperparameters: $\beta \in \{\beta_1, \beta_2\}, \gamma \in \{\gamma_1, \gamma_2, \gamma_3\}$

- How many times should CV be performed?

$$2 \times 3 = 6$$

- Total number of training runs?

$$6 \times K$$



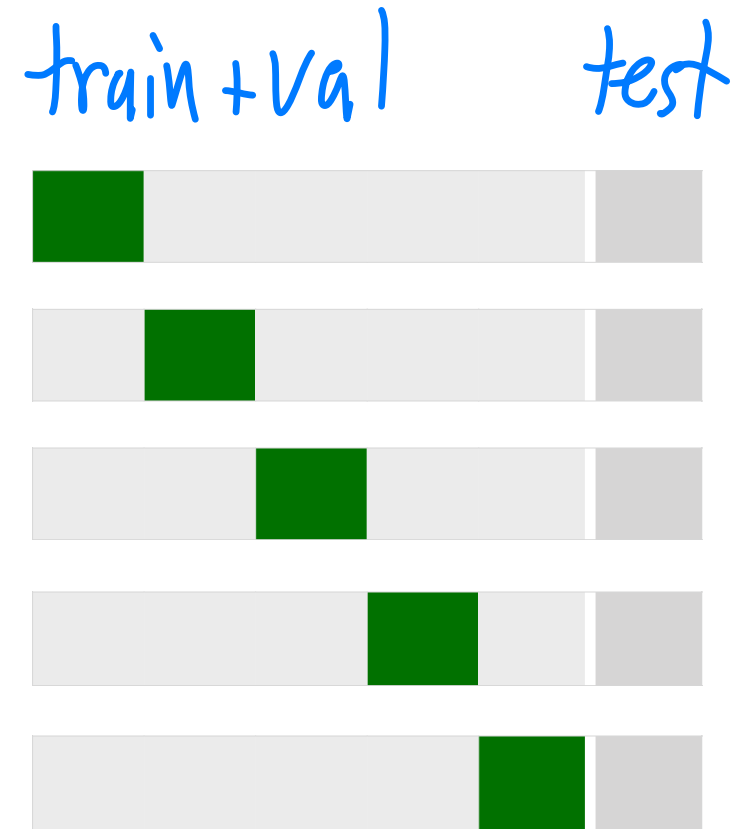
Cross-Validation: Test Error Estimation

- ▶ After Model selection we obtain some

$$\alpha^*, \beta^*$$

use ensemble method

- ▶ Retrain y on all train data with α^*, β^*
- ▶ Evaluate model on held-out test set
 - ▶ Still noisy estimate of generalization error if test set is small
- ▶ Nested cross validation!



Nested Cross-Validation

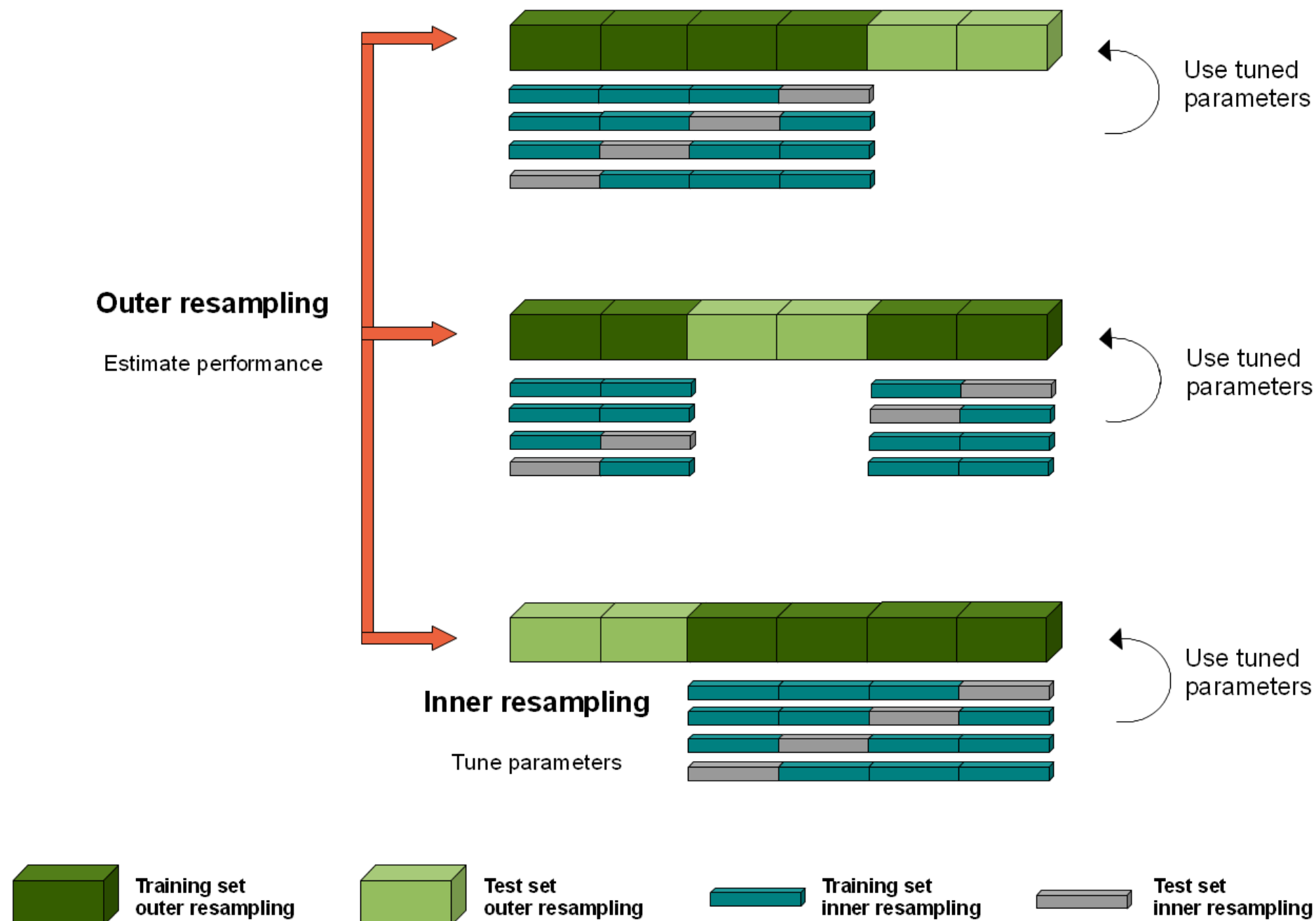


Figure: Nested cross-validation

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)



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

Why do models make errors?

What kind of errors can we expect?

- ▶ Consider dataset of observations $(\mathbf{x}, t) \sim p(\mathbf{x}, t)$ and a model $y(\mathbf{x})$
- ▶ The model makes errors (regression loss function):

$$L(t, y(\mathbf{x})) = (t - y(\mathbf{x}))^2$$

- ▶ Every time we make an observation of random variables (\mathbf{x}, t) we make a different error. We now consider the expected loss:

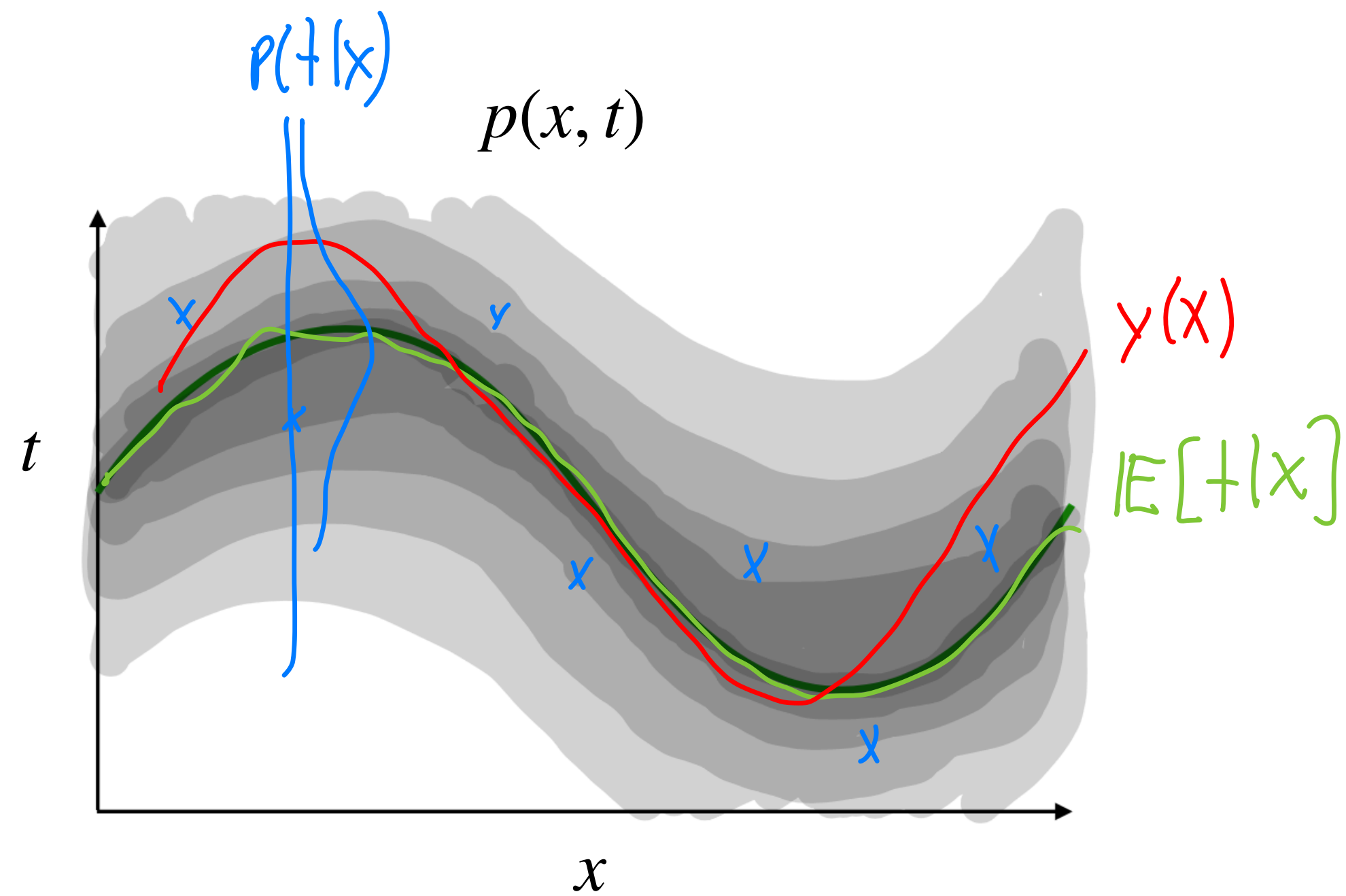
$$\mathbb{E}_{(\mathbf{x}, t) \sim p(\mathbf{x}, t)}[L(t, y(\mathbf{x}))] = \iint (t - y(\mathbf{x}))^2 p(\mathbf{x}, t) d\mathbf{x} dt$$

- ▶ The best model y we can possibly have (Bishop 1.5.5):

Regression
function

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

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 $\mathbb{E}[t | \mathbf{x}] := \mathbb{E}_{t \sim p(t|\mathbf{x})}[t | \mathbf{x}]$:

$$\mathbb{E}[L] = \int \int (y(\mathbf{x}) - \underbrace{\mathbb{E}[t | \mathbf{x}] + \mathbb{E}[t | \mathbf{x}] - t}_{=0})^2 p(\mathbf{x}, t) dt d\mathbf{x}$$

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 $\mathbb{E}[t | \mathbf{x}] := \mathbb{E}_{t \sim p(t|\mathbf{x})}[t | \mathbf{x}]$:

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

- Write out the square $((a + b)^2 = a^2 + 2ab + b^2)$ and use product rule ($p(\mathbf{x}, t) = p(t | \mathbf{x})p(\mathbf{x})$)

$$\begin{aligned} \mathbb{E}[L] &= \int \int (y(\mathbf{x}) - \mathbb{E}[t | \mathbf{x}])^2 p(t | \mathbf{x}) dt p(\mathbf{x}) d\mathbf{x} \\ &\quad + 2 \int \int (y(\mathbf{x}) - \mathbb{E}[t | \mathbf{x}]) (\mathbb{E}[t | \mathbf{x}] - t) p(t | \mathbf{x}) dt p(\mathbf{x}) d\mathbf{x} \stackrel{=0}{=} \\ &\quad \underbrace{\int t p(t | \mathbf{x}) dt}_{= \mathbb{E}[t | \mathbf{x}]} + \int \int (\mathbb{E}[t | \mathbf{x}] - t)^2 p(t | \mathbf{x}) dt p(\mathbf{x}) d\mathbf{x} \end{aligned}$$

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 $\mathbb{E}[t | \mathbf{x}] := \mathbb{E}_{t \sim p(t|\mathbf{x})}[t | \mathbf{x}]$:

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

- Write out the square $((a + b)^2 = a^2 + 2ab + b^2)$ and use product rule $(p(\mathbf{x}, t) = p(t | \mathbf{x})p(\mathbf{x}))$

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

- Thus the expected loss can be decomposed in two terms

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

Summary so far...

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

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

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

$$D_1 = \{ (x_1, t_1), \dots, (x_n, t_n) \} \rightarrow y_{D_1}$$

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

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

- What can we say about the expected loss of 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

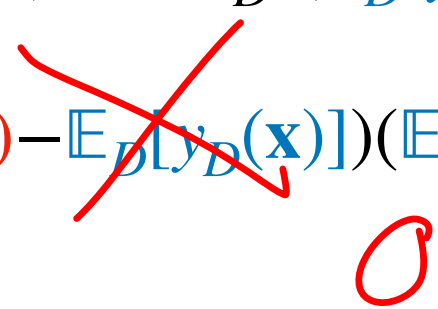
- Let's analyze performance of a learning algorithm by averaging the expected loss over learned y_D for different datasets D

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

- 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 $(a + b)^2 = a^2 + 2ab + b^2$)

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


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 $(a - b)^2 = a^2 - 2ab + b^2$)

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

(Compute expectation, rewrite)

$$\begin{aligned} &\text{variance} \quad \text{bias} \\ &= \mathbb{E}_D[(y_D(\mathbf{x}) - \mathbb{E}_D[y_D(\mathbf{x})])^2] + (\mathbb{E}_D[y_D(\mathbf{x})] - \mathbb{E}[t | \mathbf{x}])^2 \end{aligned}$$

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:

$$\begin{aligned}\mathbb{E}_D[\mathbb{E}[L]] = & \int \mathbb{E}_D[(y_D(\mathbf{x}) - \mathbb{E}_D[y_D(\mathbf{x})])^2] p(\mathbf{x}) d\mathbf{x} && \text{Variance} \\ & + \int (\mathbb{E}_D[y_D(\mathbf{x})] - \mathbb{E}[t | \mathbf{x}])^2 p(\mathbf{x}) d\mathbf{x} && \text{Bias}^2 \\ & + \int \text{var}[t | \mathbf{x}] p(\mathbf{x}) d\mathbf{x} && \text{Noise}\end{aligned}$$

Bias-Variance Decomposition: Example

- Generate $L = 100$ datasets of $N = 50$ points:

- $x \sim U(0,1)$
- $t = \sin(2\pi x) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \alpha^{-1})$
- $\mathbb{E}[t|x] = \sin(2\pi x)$

- Parametrize y with 24 Gaussian basis functions

The $L = 100$ datasets each gives a model

- $y^{(l)}(x) = (\mathbf{w}^{(l)})^T \boldsymbol{\phi}(x)$

- That minimizes:

- $$E_D = \frac{1}{2} \sum_{i=1}^N \{t_n - \mathbf{w}^T \boldsymbol{\phi}(x)\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

- Let $\bar{y}(x) = \mathbb{E}_D[y_D(x)]$ the average function

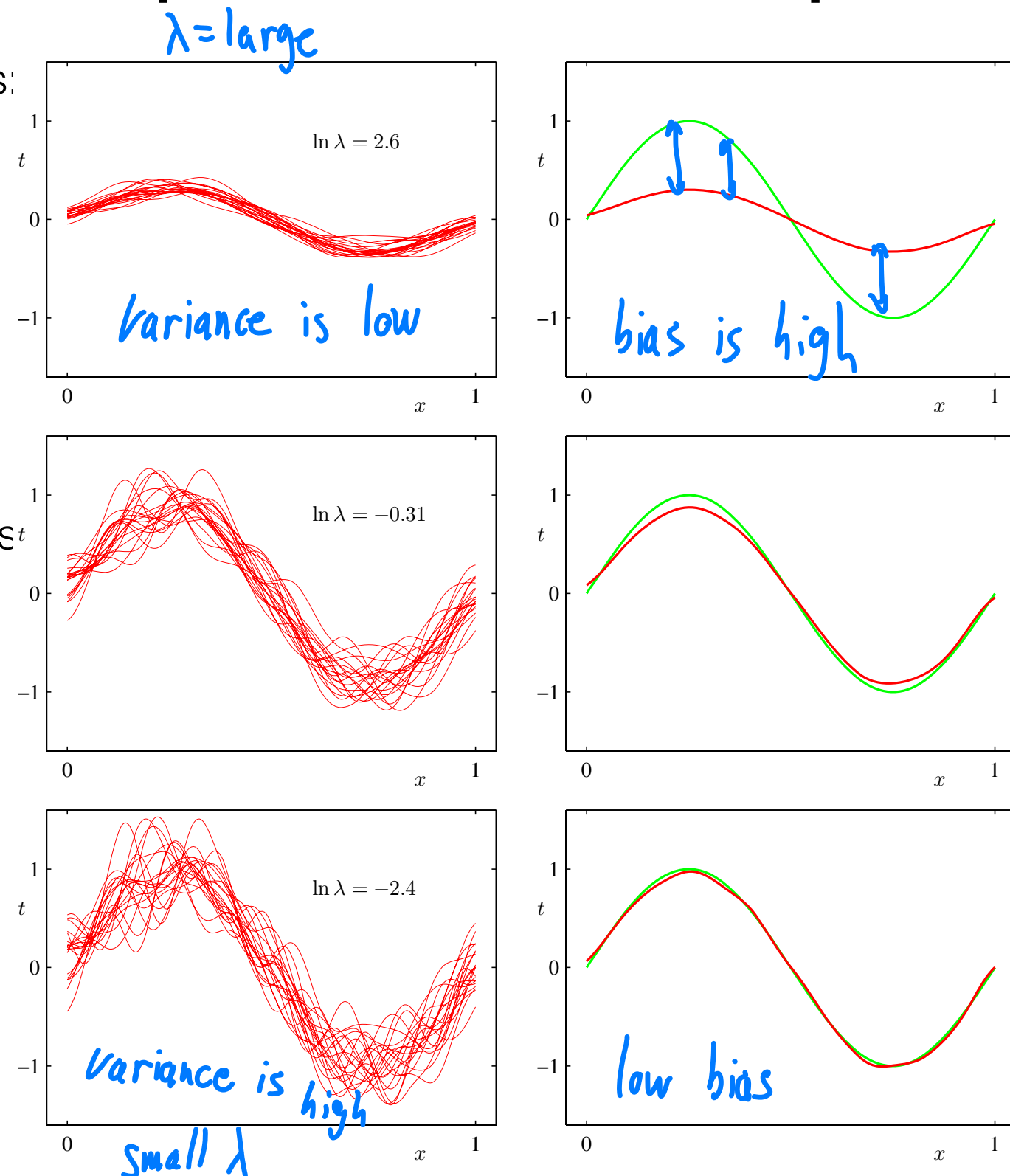


Figure: bias-variance decomposition (Bishop 3.5)

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:

$$\begin{aligned}\mathbb{E}_D[\mathbb{E}[L]] = & \int \mathbb{E}_D[(y_D(x) - \mathbb{E}_D[y_D(x)])^2] p(x) dx && \text{Variance} \\ & + \int (\mathbb{E}_D[y_D(x)] - \mathbb{E}[t | x])^2 p(x) dx && \text{Bias}^2 \\ & + \int \text{var}[t | x] p(x) dx && \text{Noise}\end{aligned}$$

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:

$$\mathbb{E}_D[\mathbb{E}[L]] \approx \frac{1}{N} \sum_{n=1}^N \frac{1}{L} \sum_{l=1}^L (\mathbf{y}^{(l)}(x_n) - \bar{\mathbf{y}}(\mathbf{x}_n))^2 \quad \text{Variance}$$
$$+ \frac{1}{N} \sum_{i=1}^N \bar{\mathbf{y}}(x_n) - \sin 2\pi x_n)^2 \quad \text{Bias}^2$$

Bias-Variance Decomposition: Example

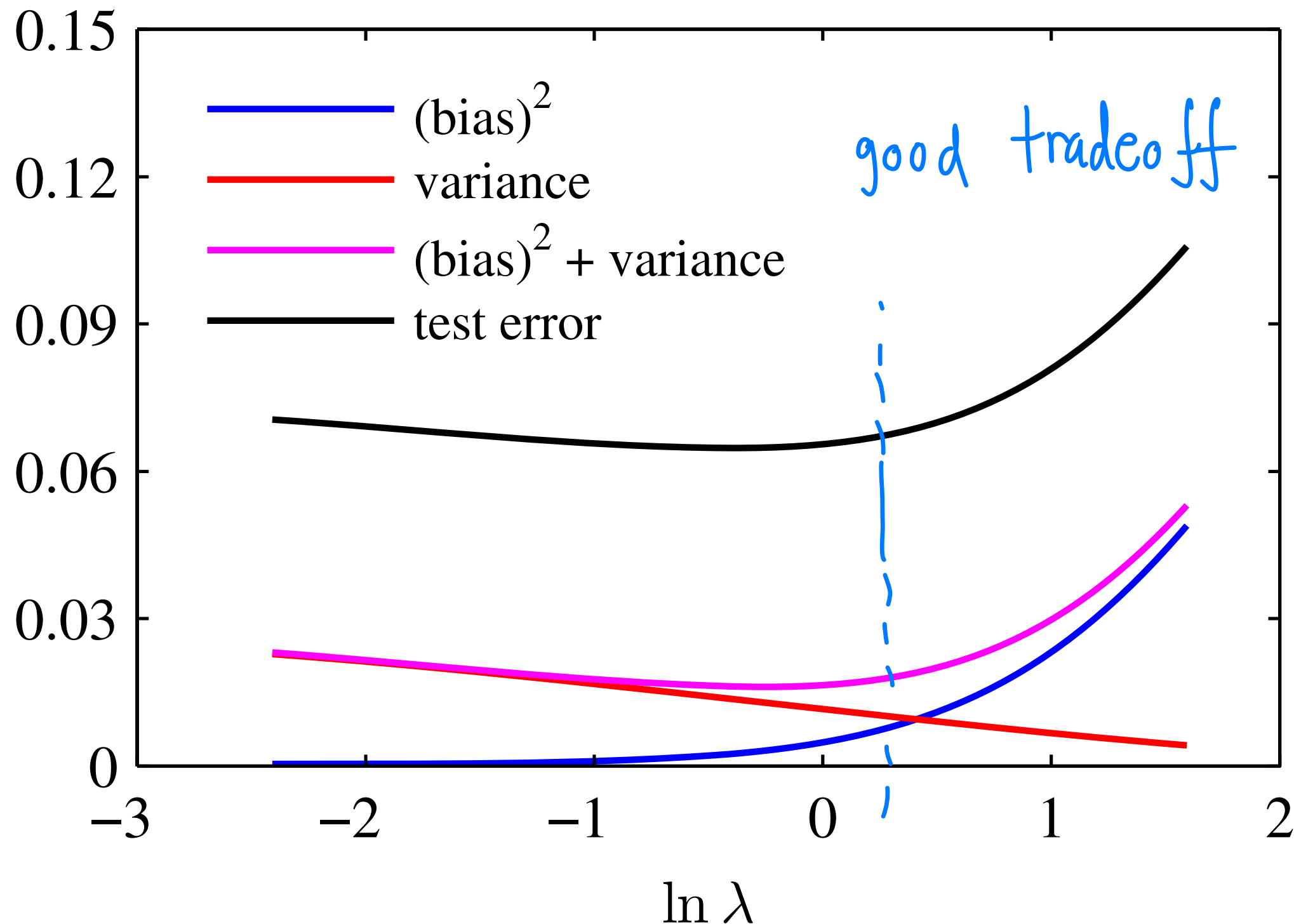


Figure: bias-variance decomposition (Bishop 3.6)

Bias-Variance Decomposition

- ▶ In practice we don't want to split our dataset into L datasets to determine the best model complexity (best value of λ)
- ▶ Better to keep large dataset
 - ▶ Less overfitting.
 - ▶ Different optimal model complexity!
- ▶ Model averaging? Bayesian regression!



Machine Learning 1

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

Erik Bekkers

(Bishop 3.3.1 (and 2.3.3))

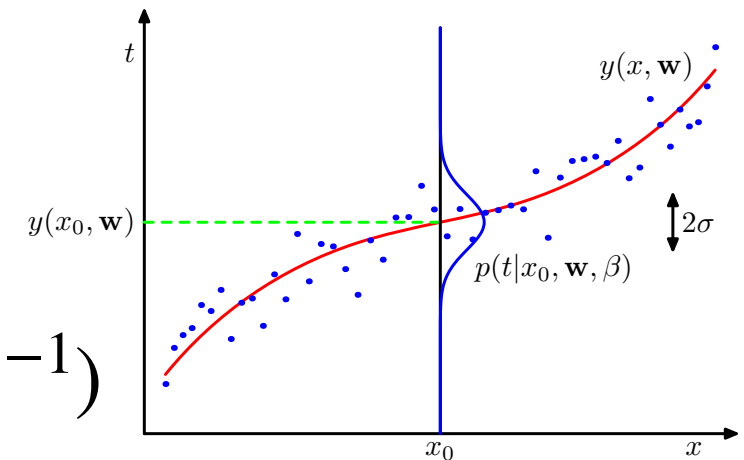


Bayesian Linear Regression

- Regression problem with:

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

- Predictive distribution $p(t' | \mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t' | \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}'), \beta^{-1})$



- Probabilistic model with Gaussians:

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

- Conjugate prior: $p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0)$

- Posterior:
$$p(\mathbf{w} | \mathbf{t}, \mathbf{X}) = \frac{p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) p(\mathbf{w})}{p(\mathbf{t} | \mathbf{X}, \beta)} = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N)$$

- Maximum A Posteriori estimate:

- $\mathbf{w}_{\text{MAP}} = \mathbf{m}_N$

Bishop Ch 2.3, Eq. 2.116

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}$$

$$\mathbf{m}_N = \mathbf{S}_N (\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \boldsymbol{\Phi}^T \mathbf{t})$$

Bayesian Linear Regression

- ▶ Special simple prior:

- ▶ $p(\mathbf{w} | \alpha) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I})$ ($\mathbf{m}_0 = \mathbf{0}$ and $\mathbf{S}_0 = \alpha^{-1} \mathbf{I}$)

- ▶ Posterior

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

- ▶ With $\mathbf{m}_N = \mathbf{S}_N (\cancel{\mathbf{S}_0^{-1}}^{\theta} \mathbf{m}_0 + \beta \Phi^T \mathbf{t}) = \beta \mathbf{S}_N \Phi^T \mathbf{t}$
 $\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \Phi^T \Phi = (\alpha \mathbf{I} + \beta \Phi^T \Phi)$

Bayesian Linear Regression

- ▶ Limiting cases of the posterior

- ▶ $p(\mathbf{w} | \mathbf{X}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N)$ with $\mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{t}$
 $\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi$

- ▶ Infinitely broad prior ($p(\mathbf{w} | \alpha) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I})$, $\alpha \rightarrow 0$):

(No restriction/assumption on \mathbf{w})

- ▶ $\lim_{\alpha \rightarrow 0} \mathbf{m}_N = \lim_{\alpha \rightarrow 0} \beta (\cancel{\alpha \mathbf{I}} + \beta \Phi^T \Phi)^{-1} \Phi^T \mathbf{t} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} = \mathbf{w}_{ML}$

- ▶ Infinitely narrow prior ($\alpha \rightarrow \infty$):

- ▶ $\lim_{\alpha \rightarrow \infty} \mathbf{m}_N = \lim_{\alpha \rightarrow \infty} \beta (\alpha \mathbf{I} + \beta \Phi^T \Phi)^{-1} \Phi^T \mathbf{t} = \dots = \underline{\mathbf{w}_0}$
- ▶ $\lim_{\alpha \rightarrow \infty} \mathbf{S}_N = \lim_{\alpha \rightarrow \infty} (\alpha \mathbf{I} + \beta \Phi^T \Phi)^{-1} = \mathbf{0} \quad \left(\begin{smallmatrix} \text{zero} \\ \text{matrix} \end{smallmatrix} \right)$



Machine Learning 1

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

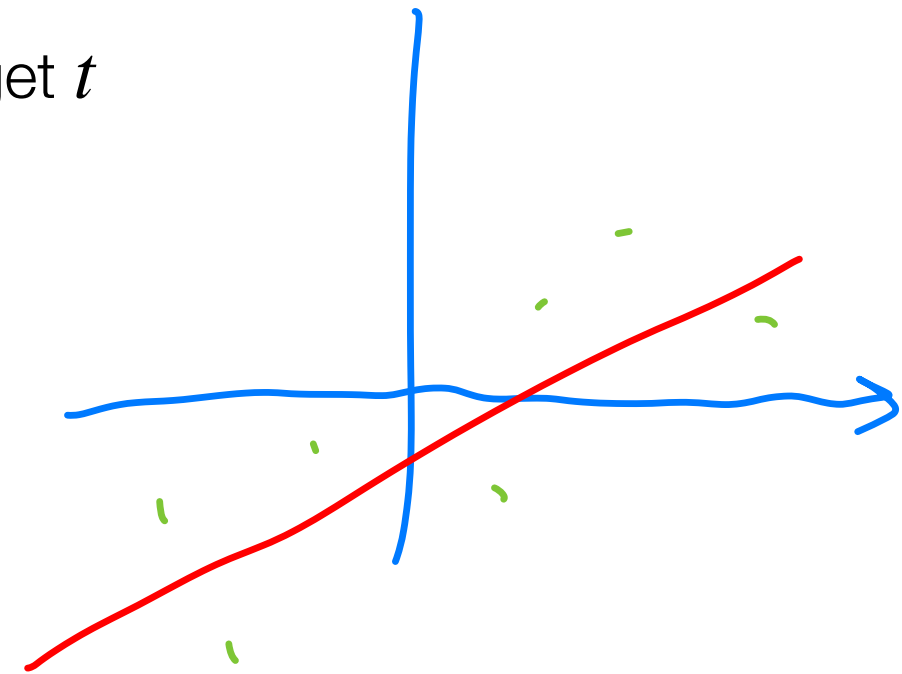
Erik Bekkers

(Bishop 3.3.1)



Example: Sequential Bayesian Learning

- ▶ Data come in as a sequences of observations of input x , target t
- ▶ **Synthetic data generated by**
 - ▶ $x \sim \mathcal{U}(x | -1, 1)$
 - ▶ $t = f(x, \mathbf{a}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, 0.2^2)$
 - ▶ $f(x, \mathbf{a}) = a_0 + a_1 x, \quad a_0 = -0.3, a_1 = 0.5$
- ▶ **Modeling choices**
 - ▶ Target distribution: $p(t' | x', \mathbf{w}, \beta) = \mathcal{N}(t' | y(x', \mathbf{w}), \beta^{-1}), \quad \beta^{-1} = 0.2^2$
 - ▶ Linear model: $y(x, \mathbf{w}) = w_0 + w_1 x$
 - ▶ Prior: $p(\mathbf{w} | \alpha) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I}), \quad \alpha = 2$
- ▶ **Sequential Bayesian Learning**: Posterior after $N - 1$ observations is prior for arrival of N^{th} datapoint!



Example: Sequential Bayesian Learning

- Data generated by

$$t = -0.3 + 0.5x + \epsilon$$

- Prior

$$p(\mathbf{w} | \alpha) = \mathcal{N}(\mathbf{w} | \mathbf{0}, \alpha^{-1} \mathbf{I})$$

- Sample 1 datapoint ★

$$(x_1, t_1)$$

- Likelihood

$$p(t_1 | x_1, \mathbf{w}, \beta) = \mathcal{N}(t_1 | w_0 + w_1 x, \beta)$$

- Posterior

$$p(\mathbf{w} | x_1, t_1, \alpha, \beta) \propto p(t_1 | x_1, \mathbf{w}, \beta) \cdot p(\mathbf{w} | \alpha)$$

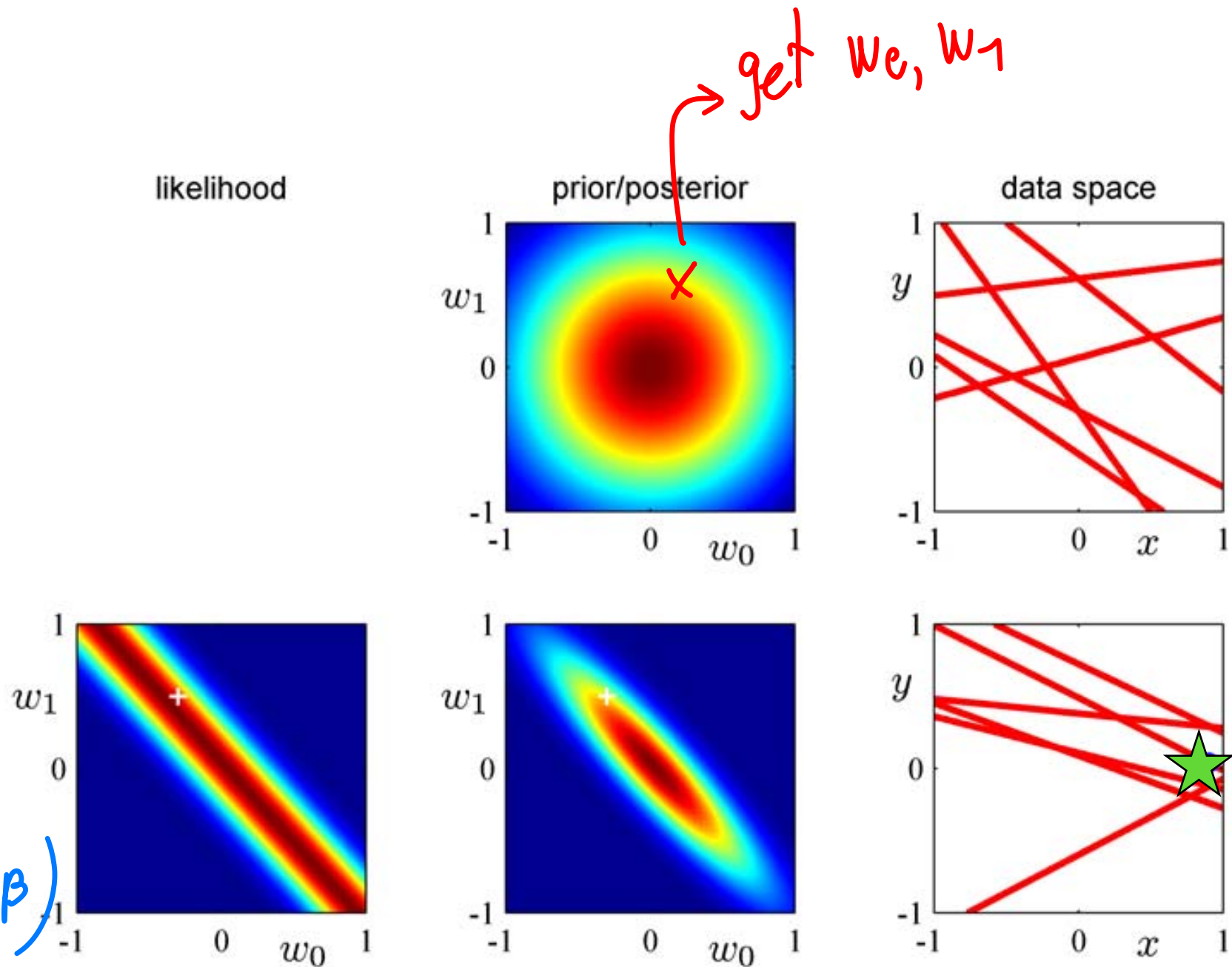


Figure: Sequential Bayesian learning (Bishop 3.7)

Example: Sequential Bayesian Learning

- Sample 2nd data point ★

- Posterior \rightarrow Prior

- Likelihood

$$p(t_2 | x_2, \mathbf{w}, \beta)$$

$$p(\mathbf{w} | (x_1, t_1), (x_2, t_2), \beta) =$$

$$= \frac{p(t_1 | \dots) \cdot p(t_2 | \dots) \cdot p(\mathbf{w} | \dots)}{p(t_1 | x) \cdot p(t_2 | x)}$$

- Posterior

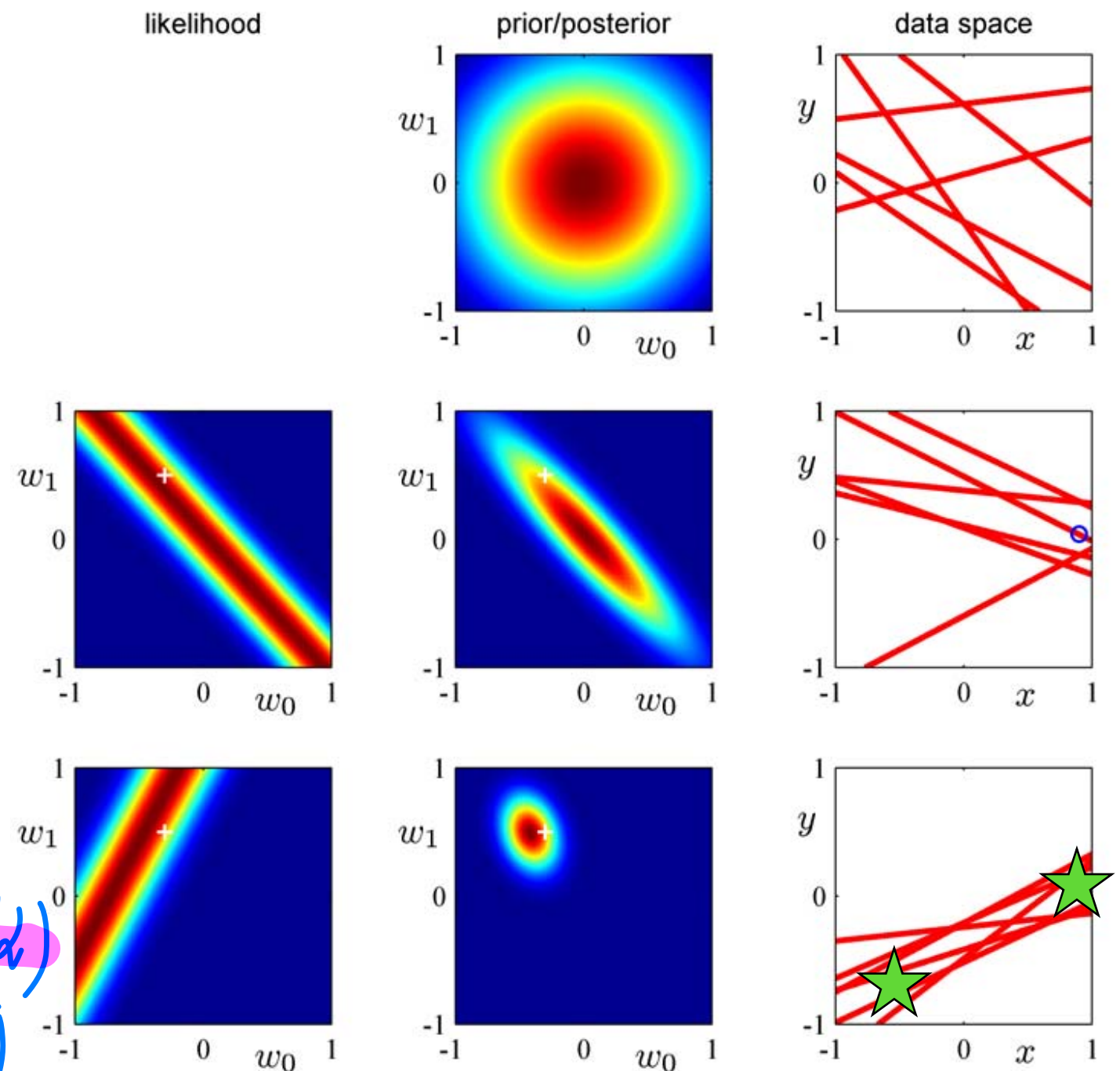


Figure: Sequential Bayesian learning (Bishop 3.7)

$$p(\mathbf{w} | (x_1, t_1), (x_2, t_2), \alpha, \beta) \propto p(t_2 | x_2, \mathbf{w}, \beta) p(\mathbf{w} | (x_1, t_1), \alpha, \beta)$$

Example: Sequential Bayesian Learning

‣ After 19 data points

‣ Prior

$$p(\mathbf{w} \mid \{(x_n, t_n)\}_{n=1}^{19}, \alpha, \beta)$$

‣ Likelihood

$$p(t_{20} \mid x_{20}, \mathbf{w}, \beta)$$

‣ Posterior

$$p(\mathbf{w} \mid \{(x_n, t_n)\}_{n=1}^{20}, \alpha, \beta) \propto p(t_{20} \mid x_{20}, \mathbf{w}, \beta) p(\mathbf{w} \mid \{(x_n, t_n)\}_{n=1}^{19}, \alpha, \beta)$$

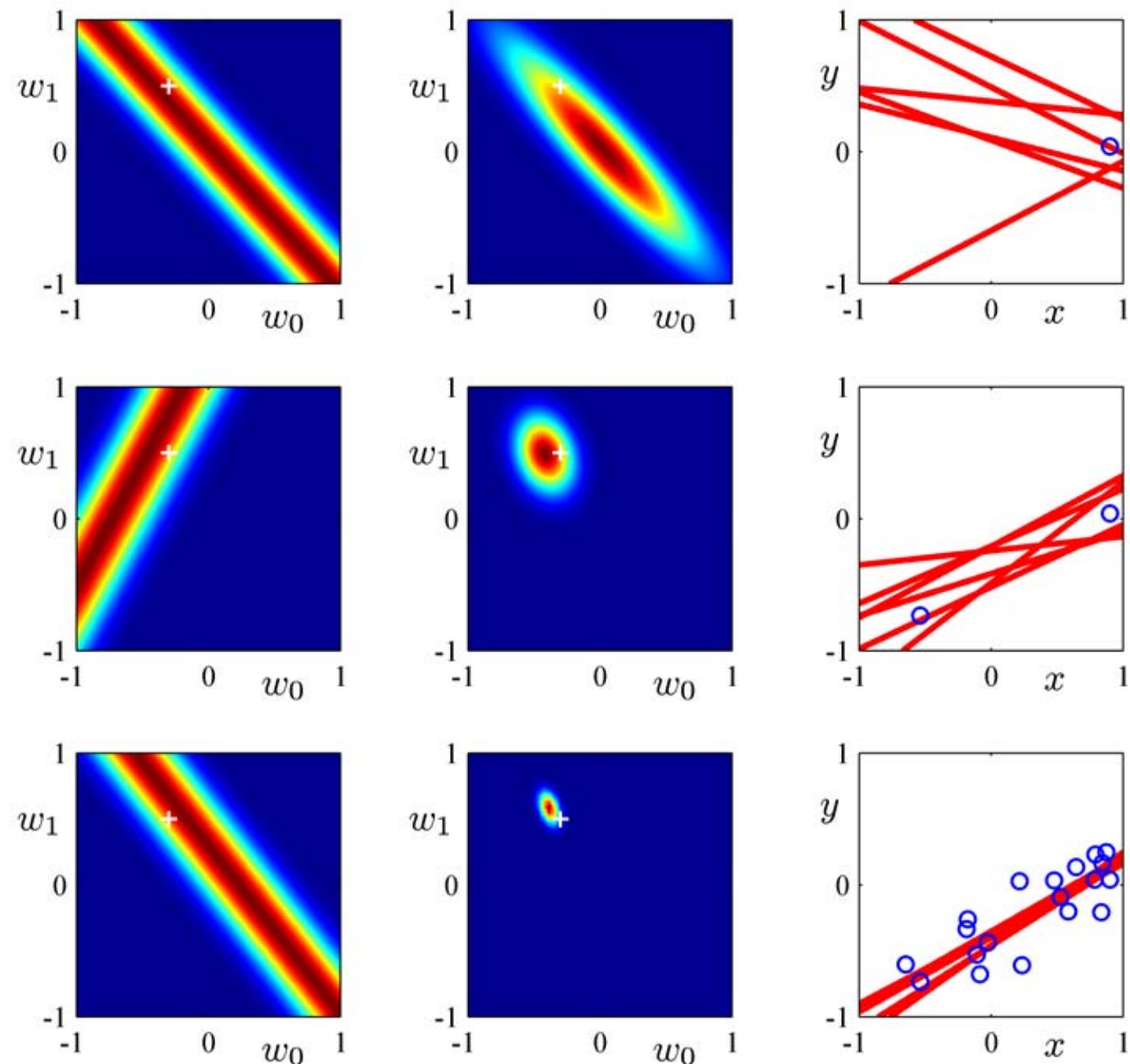


Figure: Sequential Bayesian learning (Bishop 3.7)

Bayesian Linear Regression

- ▶ Limiting cases of the posterior

- ▶ $p(\mathbf{w} | \mathbf{X}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N)$ with $\mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{t}$
 $\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi$

- ▶ After infinite amount of data ($N \rightarrow \infty$):

- ▶ $\lim_{N \rightarrow \infty} [\Phi^T \Phi]_{ij} = \lim_{N \rightarrow \infty} \sum_{n=1}^N \phi_i(\mathbf{x}_n) \phi_j(\mathbf{x}_n)$
 - ▶ $\lim_{N \rightarrow \infty} \mathbf{S}_N =$
 - ▶ $\lim_{N \rightarrow \infty} \mathbf{m}_N = \lim_{N \rightarrow \infty} \beta (\alpha \mathbf{I} + \beta \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$

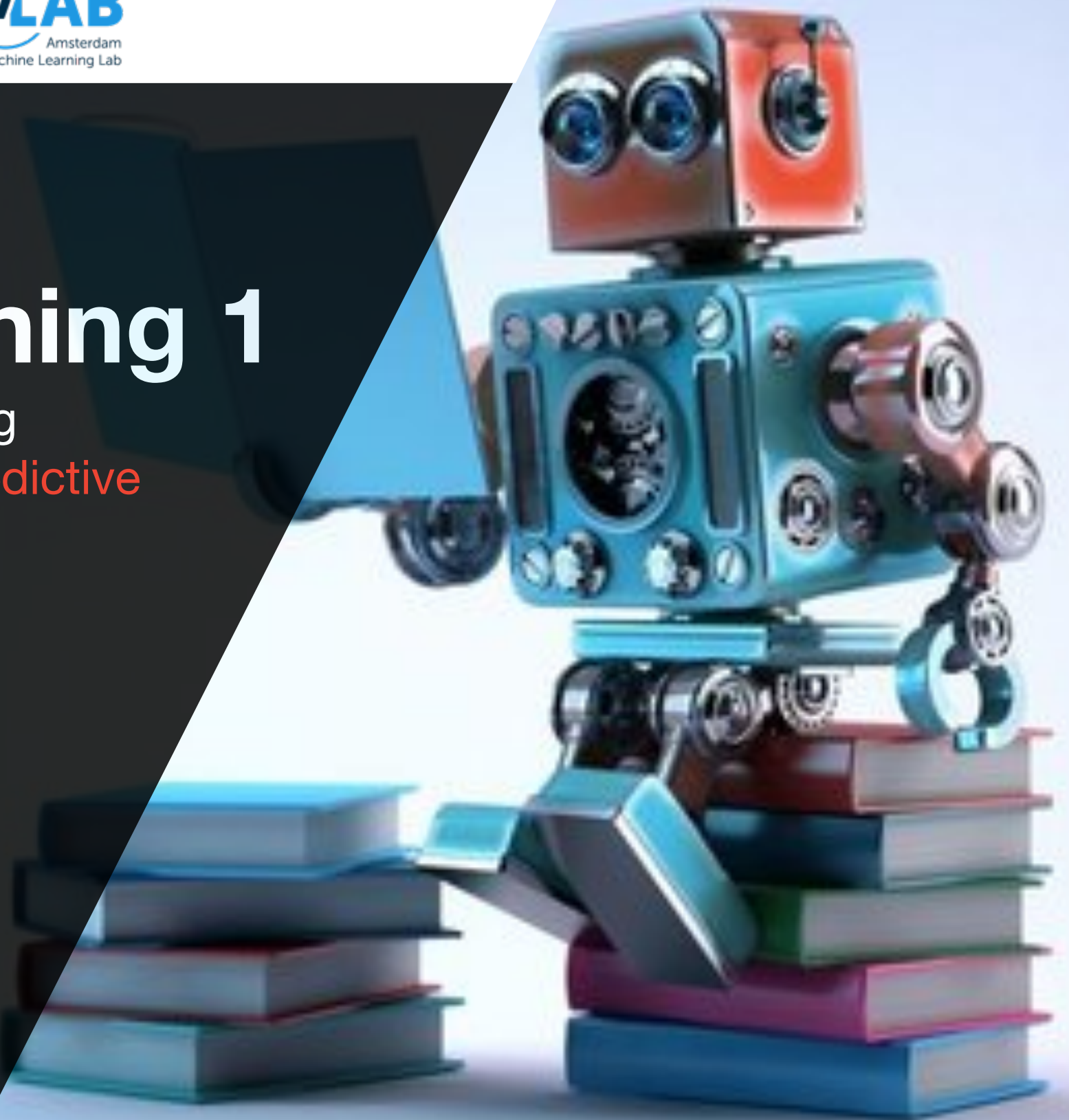


Machine Learning 1

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

Erik Bekkers

(Bishop 3.3.2)



Predictive Distribution

- ▶ Observed dataset with inputs $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$ and targets $\mathbf{t} = (t_1, \dots, t_N)^T$
- ▶ Gaussian Posterior distribution (from Gaussian prior and Gaussian likelihood)

$$\begin{aligned} \text{▶ } p(\mathbf{w} \mid \mathbf{X}, \mathbf{t}, \alpha, \beta) &= \mathcal{N}(\mathbf{w} \mid \mathbf{m}_N, \mathbf{S}_N) && \text{with } \mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{t} \\ & && \mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi \end{aligned}$$

- ▶ Parametrized Gaussian predictive distribution:

$$\text{▶ } p(t' \mid \mathbf{x}', \mathbf{w}, \beta) = \mathcal{N}(t' \mid \phi(\mathbf{x}')^T \mathbf{w}, \beta^{-1})$$

- ▶ Gaussian Bayesian predictive distribution for new input

$$\begin{aligned} \text{▶ } p(t' \mid x', \mathbf{X}, \mathbf{t}, \alpha, \beta) &= \int \mathcal{N}(t' \mid \phi(\mathbf{x}')^T \mathbf{w}, \beta^{-1}) \mathcal{N}(\mathbf{w} \mid \mathbf{m}_N, \mathbf{S}_N) d\mathbf{w} \\ &= \mathcal{N}(t' \mid \mathbf{x}', \phi(\mathbf{x}')^T \mathbf{m}_N, \sigma_N^2(\mathbf{x}')) \end{aligned}$$

Bishop Eq.
2.115

- ▶ With $\sigma_N^2(\mathbf{x}') = \beta^{-1} + \phi(\mathbf{x}')^T \mathbf{S}_N \phi(\mathbf{x}')$

Predictive Distribution

- Datasets:
 - $t = \sin(2\pi x) + \epsilon$
 - $\epsilon \sim \mathcal{N}(0, \beta^{-1})$
- Dataset sizes:
 - $N = 1, 2, 4, 25$
- Model:
 - $y(x, \mathbf{w}) = \boldsymbol{\phi}(x)^T \mathbf{w}$
 - $\phi_j(x)$: Gaussian basis functions

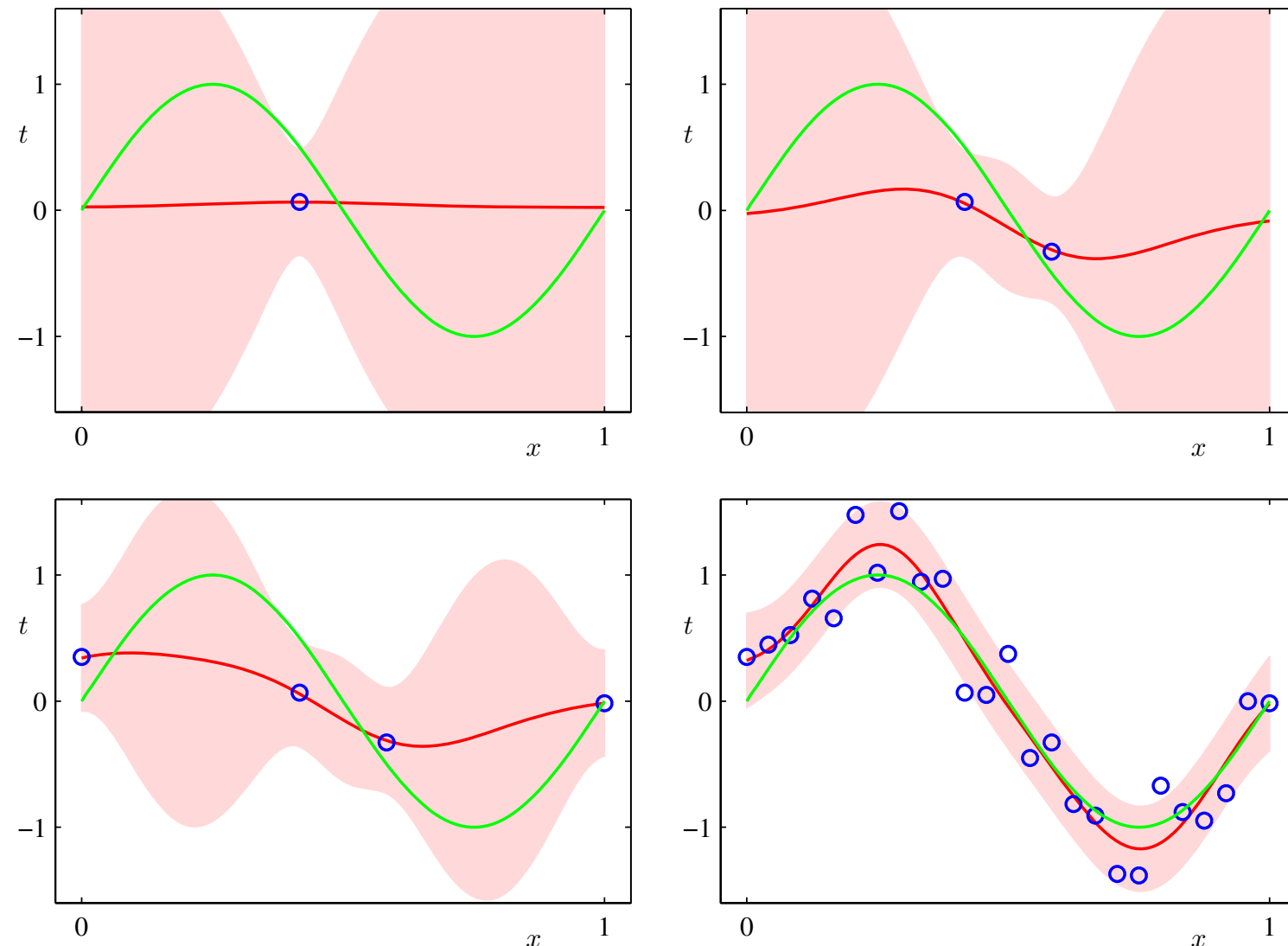


Figure: Predictive distribution (Bishop 3.8)

- Predictive distribution:
 - $p(t' | x', \mathbf{X}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(t' | \mathbf{x}', \boldsymbol{\phi}(\mathbf{x}')^T \mathbf{m}_N, \sigma_N^2(\mathbf{x}'))$
 - $\sigma_N^2(\mathbf{x}') = \beta^{-1} + \boldsymbol{\phi}(\mathbf{x}')^T \mathbf{S}_N \boldsymbol{\phi}(\mathbf{x}'), \quad \mathbf{m}_N = \beta \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t}, \quad \mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}$

Samples drawn from Bayesian Predictive Distribution

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

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

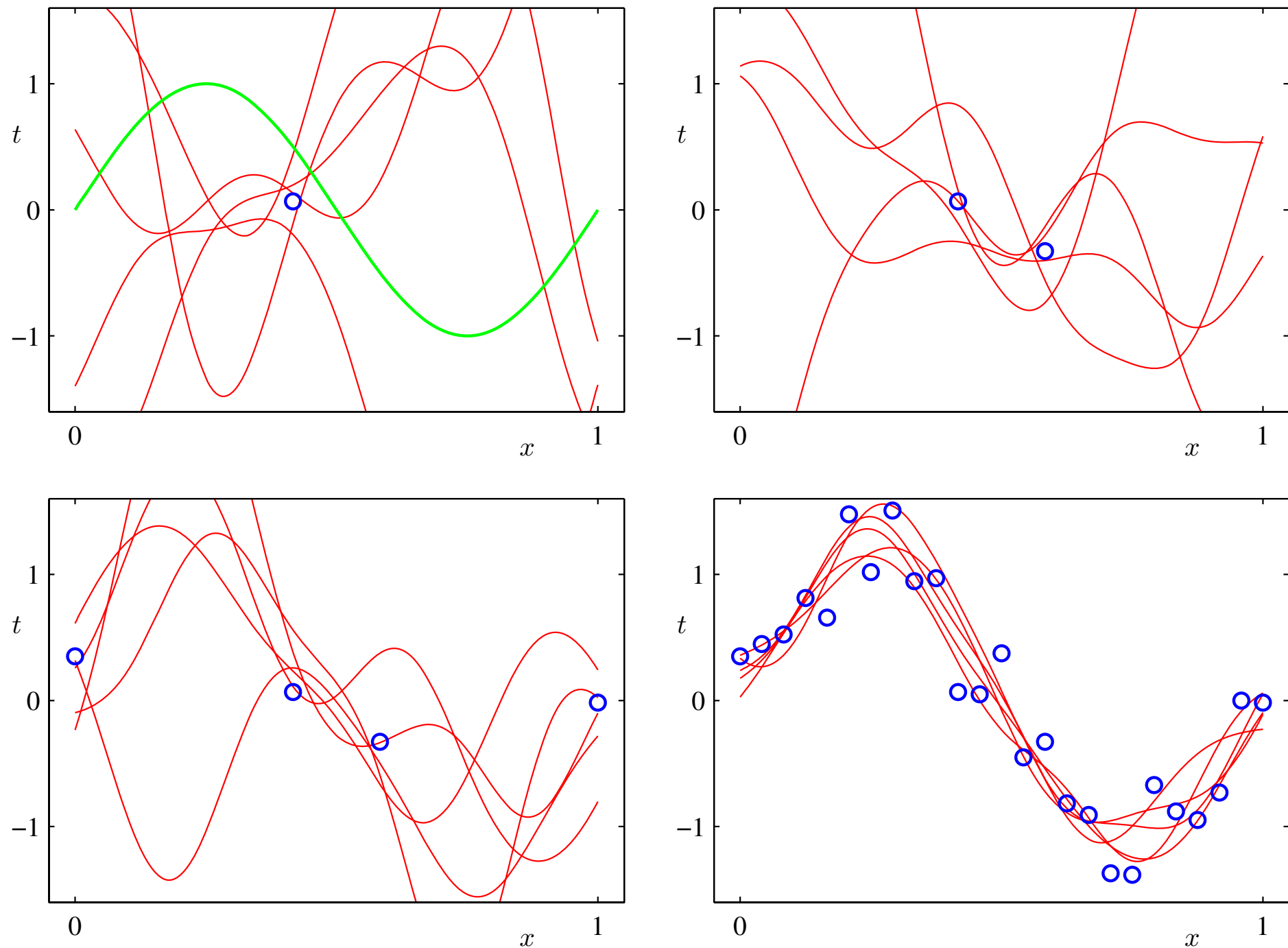


Figure: Sample functions $y(x, \mathbf{w})$ with \mathbf{w} sampled from posterior distribution (Bishop 3.9)