

Machine Learning 1

Lecture 5 - Classification

Recap/Catch-up:

Bias-Variance

Gaussian Posteriors,

Sequential Bayesian Learning,

Bayesian Predictive Distributions

Equivalent Kernel

Continue with:

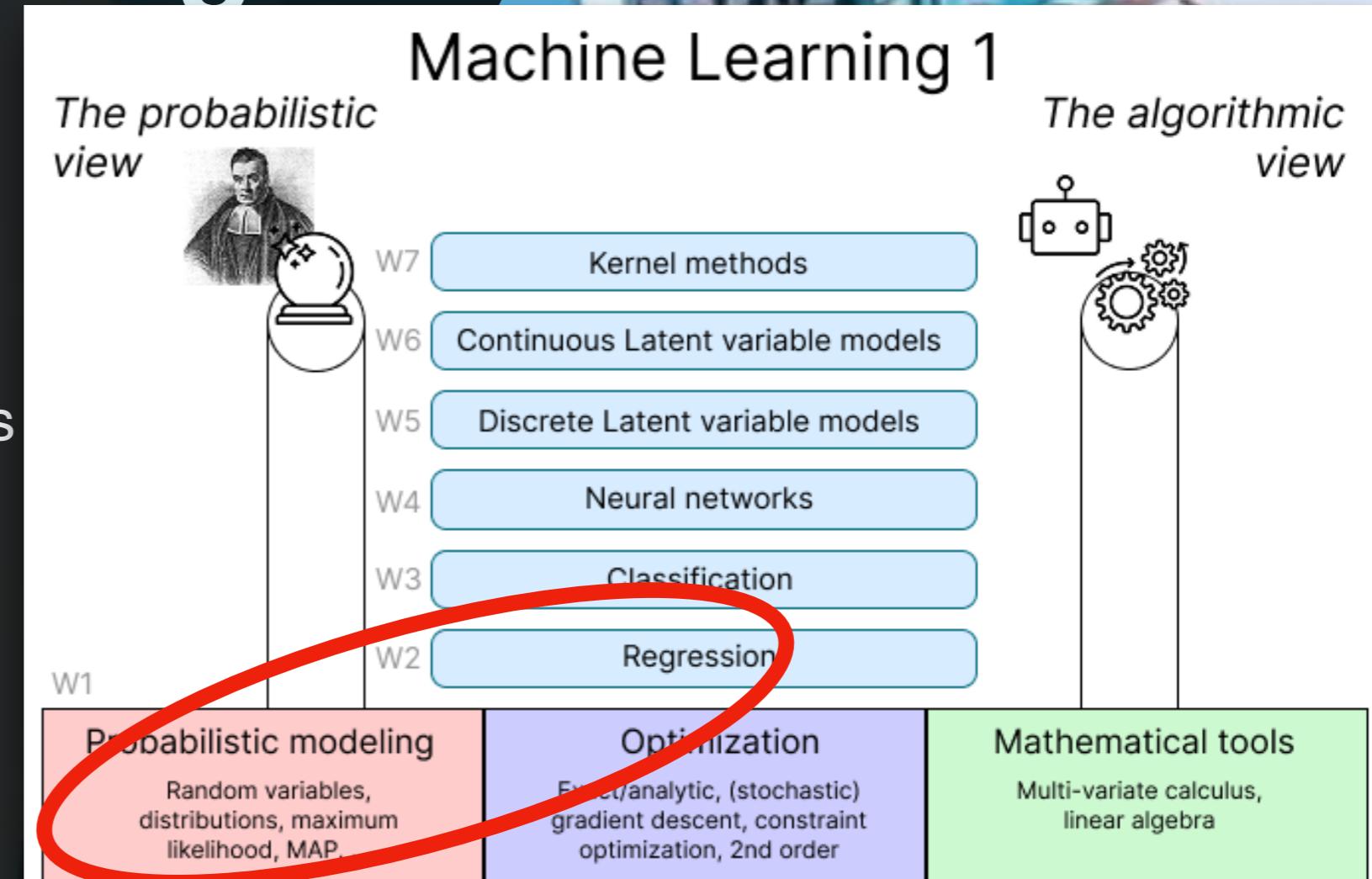
~~Bayesian Model Comparison~~

~~Empirical Bayes~~

Classification

Decision Theory

Generative Models



Erik Bekkers

Machine Learning 1

Lecture 4.2 - Supervised Learning
Bias Variance Decomposition

Erik Bekkers

(Bishop 1.5.5, 3.2)

Recap

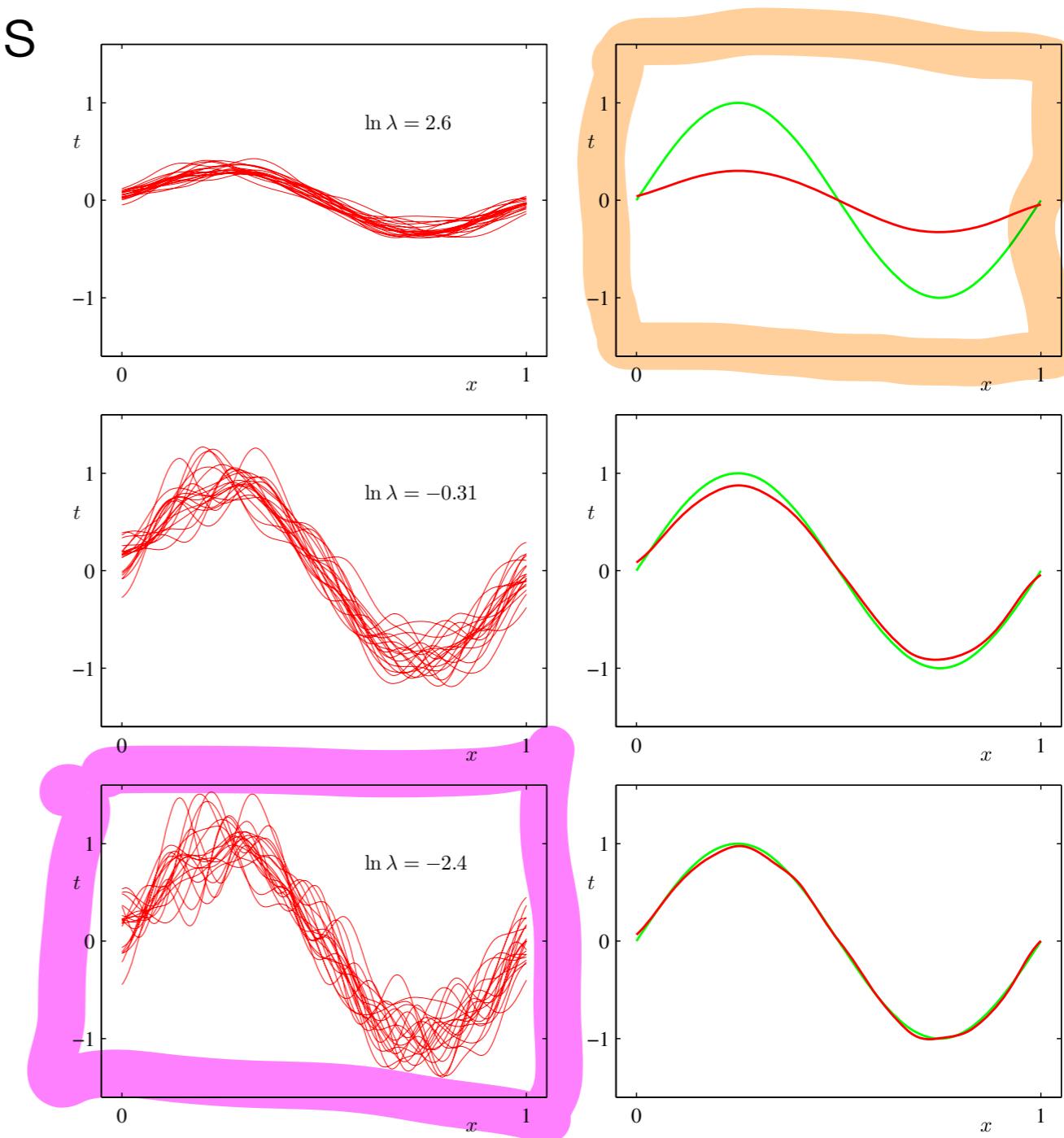


Bias-Variance Summary slide

- ▶ The bias-variance decomposition is a way to **analyze a model's prediction error**.
- ▶ It tells us **in expectation** (over all possible datasets) a model's error consists of three parts:

Expected Loss

$$= \text{Bias}^2 + \text{Variance} + \text{Noise}$$



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):

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

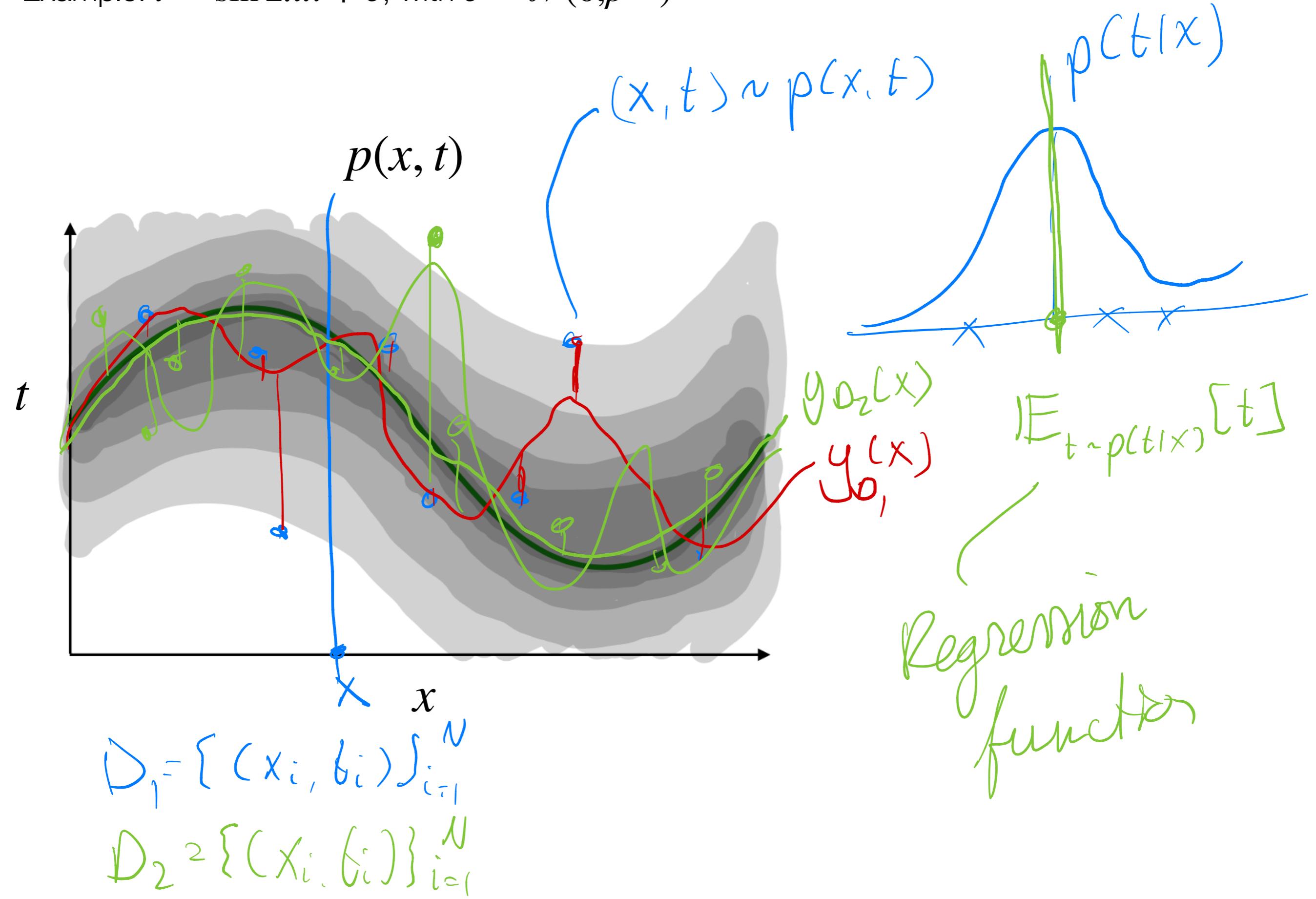
- Every time we make an observation of random variables (\mathbf{x}, t) we make a different error. The best model we can possibly have (Bishop 1.5.5):

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

- We now consider the expected loss:

$$\mathbb{E}_{(\mathbf{x}, t) \sim p(\mathbf{x}, t)}[L(t, y(\mathbf{x}))] = \underbrace{\int \int (t - y(\mathbf{x}))^2 p(\mathbf{x}, t) d\mathbf{x} dt}_{\text{(Bias)}^2 + \text{Variance} + \text{Noise}}$$

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



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:

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

average model

Machine Learning 1

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

Erik Bekkers

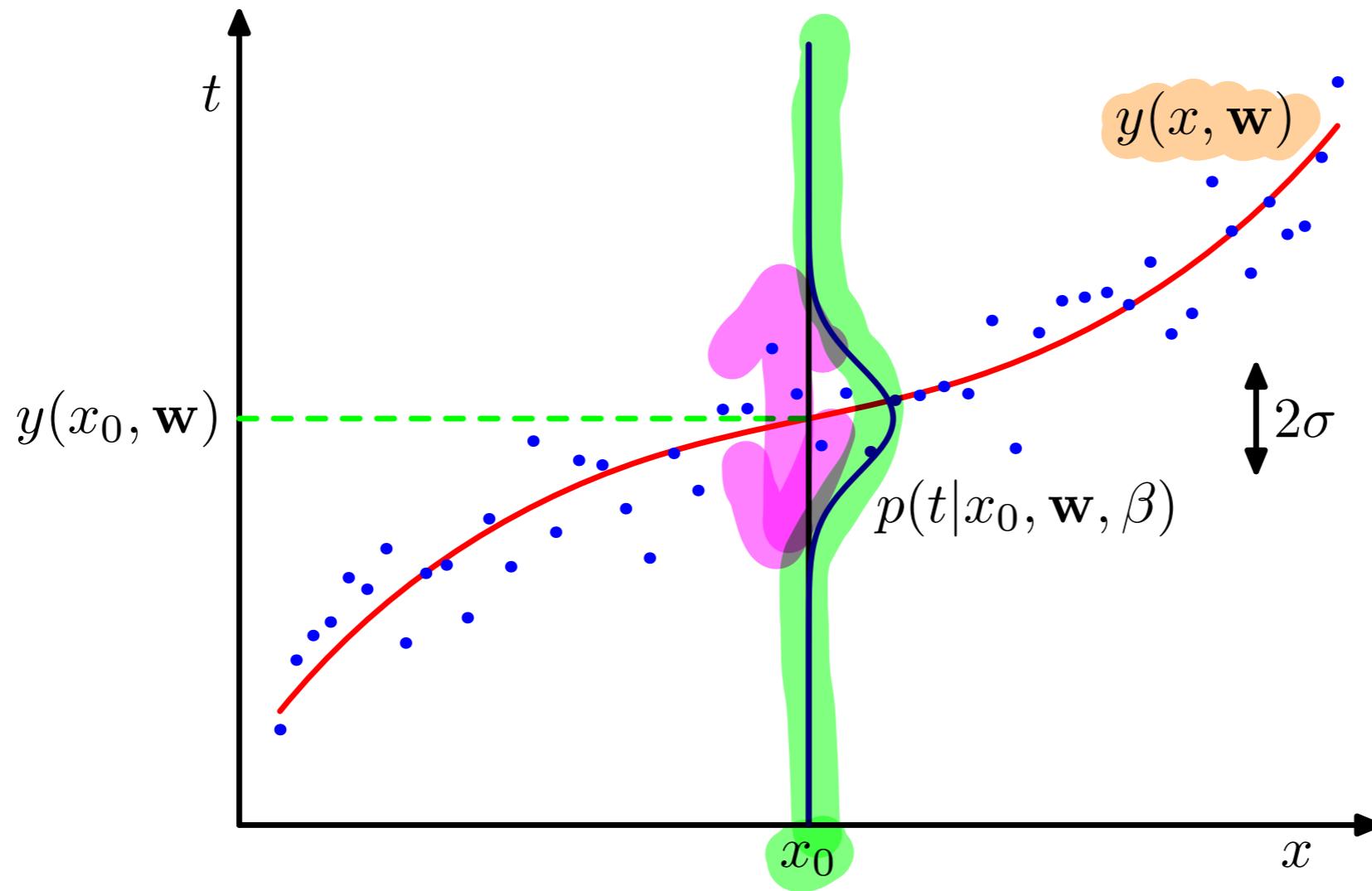
(Bishop 3.3.1 (and 2.3.3))



Bayesian Linear Regression

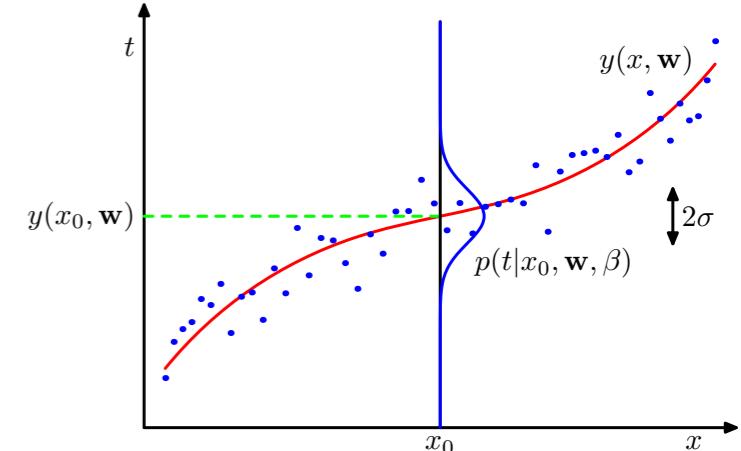
- Regression problem with:

- Data: $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$, $\mathbf{t} = (t_1, \dots, t_N)^T$
- Predictive distribution $p(t' | \mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t' | \mathbf{w}^T \phi(\mathbf{x}'), \beta^{-1})$



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

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

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

conjugate here means
that the posterior will
be of the same family of
distributions (Gaussian)

Bishop Ch 2.3, Eq. 2.116

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

Bayesian Linear Regression

- Regression problem with:

So, with a **Gaussian likelihood and **Gaussian** prior**

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

- Predictive distribution $p(t|x_0, \mathbf{w}, \beta)$ explicit expression for the posterior
(which is again **Gaussian**)

- Probabilistic model with Gaussians:

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

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

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



Machine Learning 1

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

Erik Bekkers

(Bishop 3.3.1)

Recap



Example: Sequential Bayesian Learning

- Data come in as 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: $p(t' | x', \mathbf{w}, \beta) = \mathcal{N}(t' | y(x', \mathbf{w}), \beta^{-1}), \quad \beta^{-1} =$

- 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!

$$p(\mathbf{w} | D_N) = \frac{p(\mathbf{x}_N | \mathbf{w}) p(\mathbf{w} | D_{N-1})}{p(x_N)}$$

$$D_{N-1} = \{(x_1, t_1), \dots, (x_{N-1}, t_{N-1})\}$$

size is $N-1$

$$D_N = D_{N-1} \cup (x_N, t_N)$$

Example: Sequential Bayesian Learning

- Data generated by

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

$\omega_0 + \omega_1 x + \epsilon$

- Prior

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

- Sample 1 datapoint

(x_1, t_1)

- Likelihood

$$p(t_1 | x_1, \mathbf{w}, \beta) = \mathcal{N}(t_1 | w_0 + w_1 x_1, \beta^{-1})$$

- Posterior

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

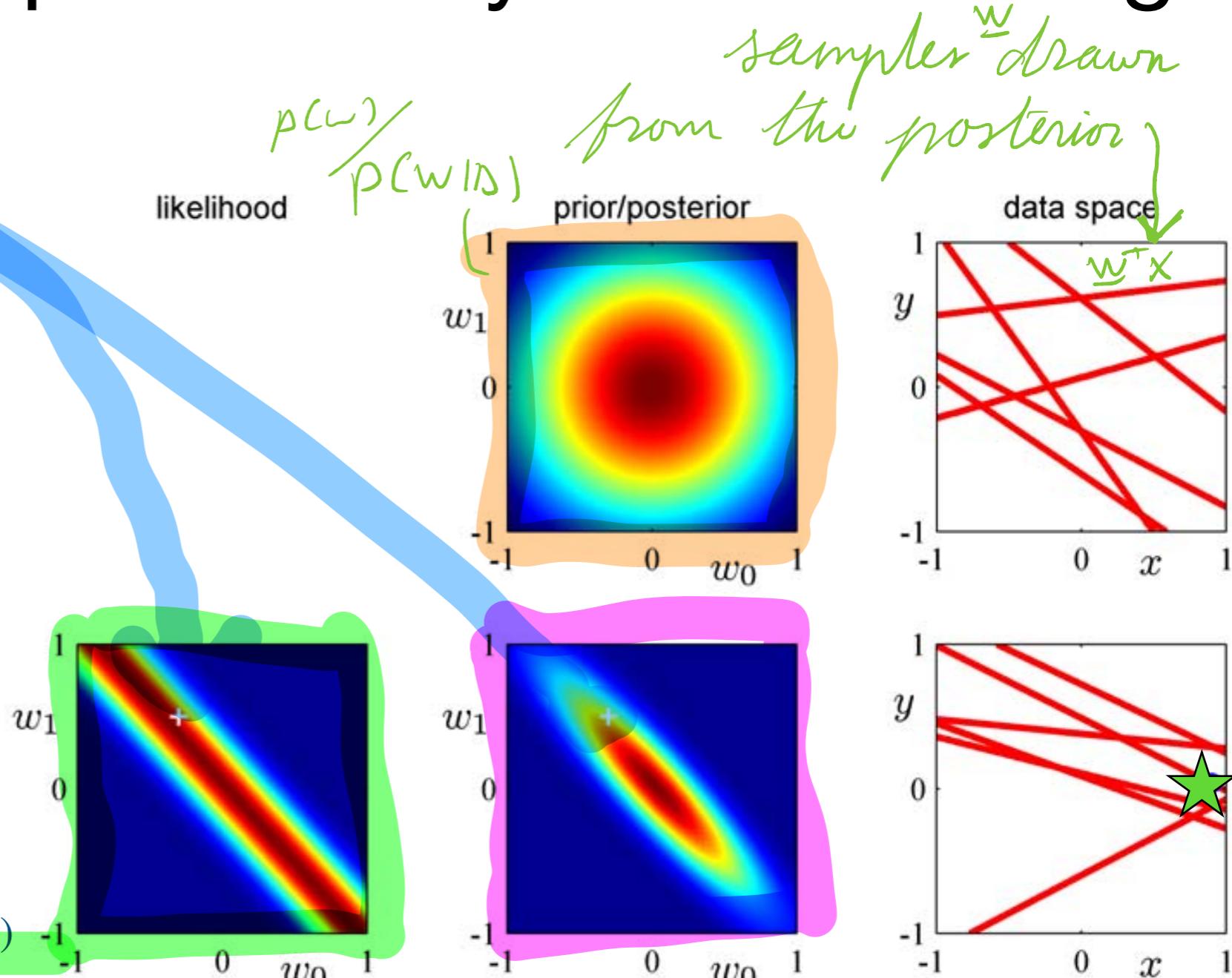


Figure: Sequential Bayesian learning (Bishop 3.7)

proportional to the same up to a constant

Example: Sequential Bayesian Learning

- Sample 2nd data point ★

- Posterior → Prior

- Likelihood

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

- Posterior

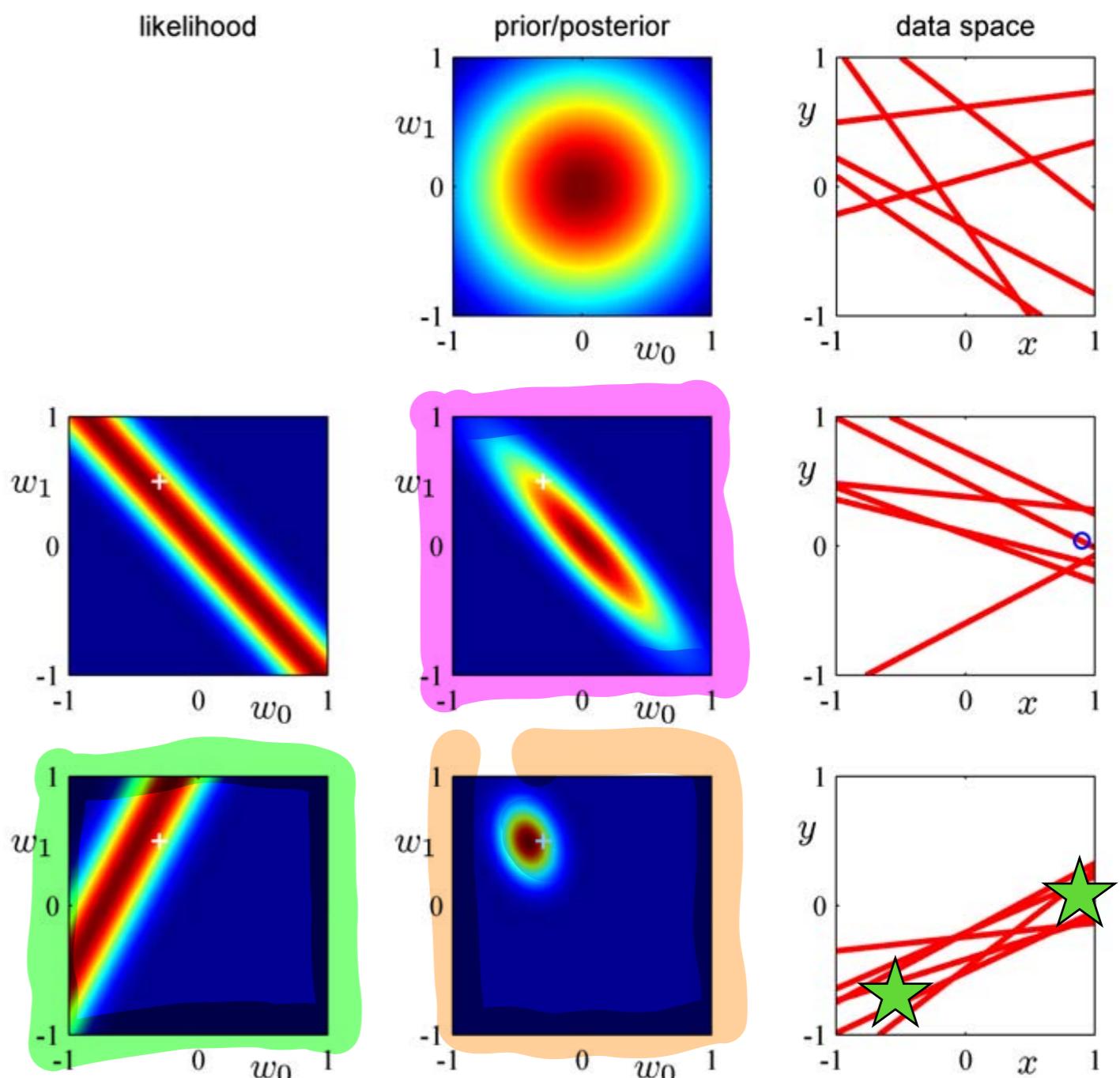
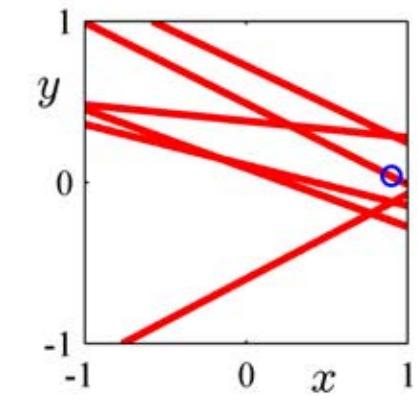
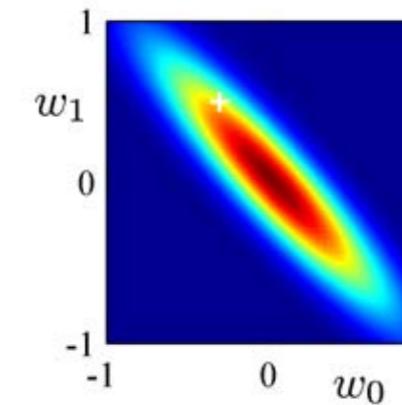
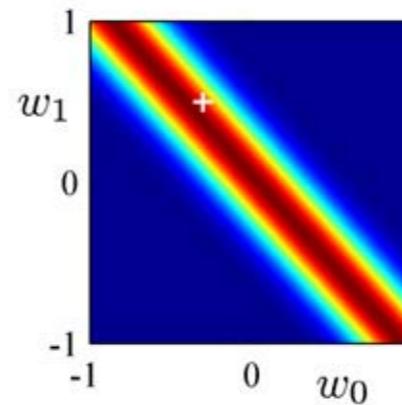


Figure: Sequential Bayesian learning (Bishop 3.7)

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

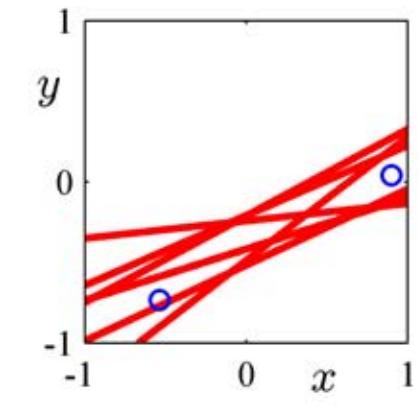
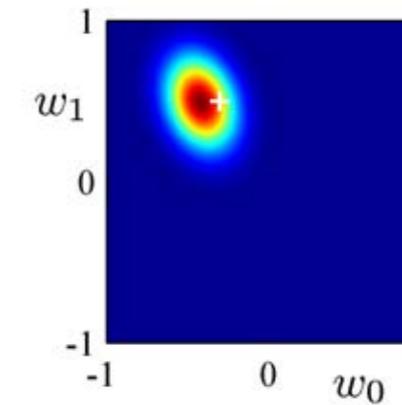
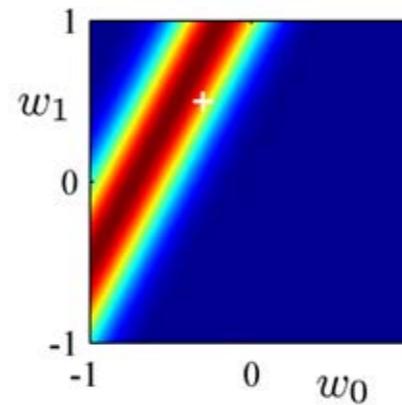
Example: Sequential Bayesian Learning

- After 19 data points



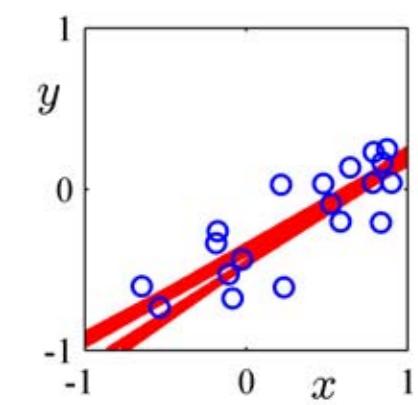
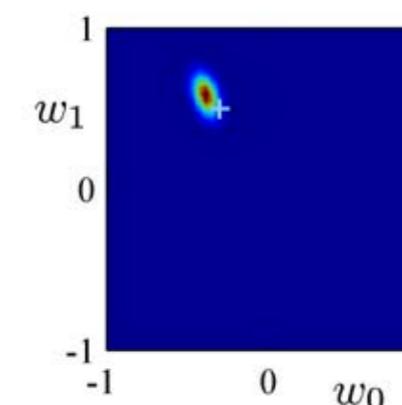
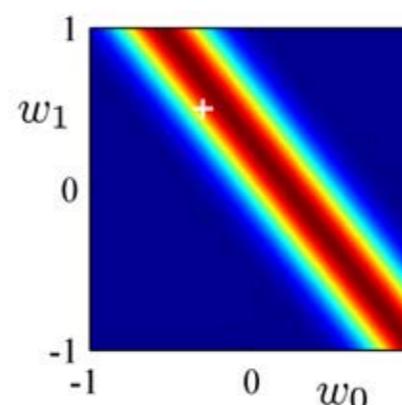
- Prior

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



- Likelihood

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



- Posterior

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

Figure: Sequential Bayesian learning (Bishop 3.7)

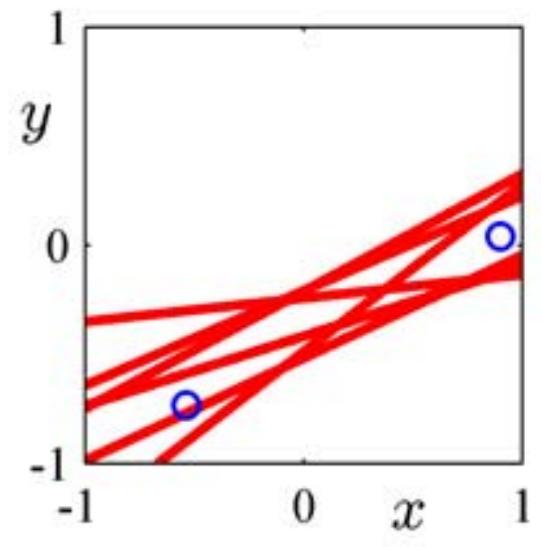
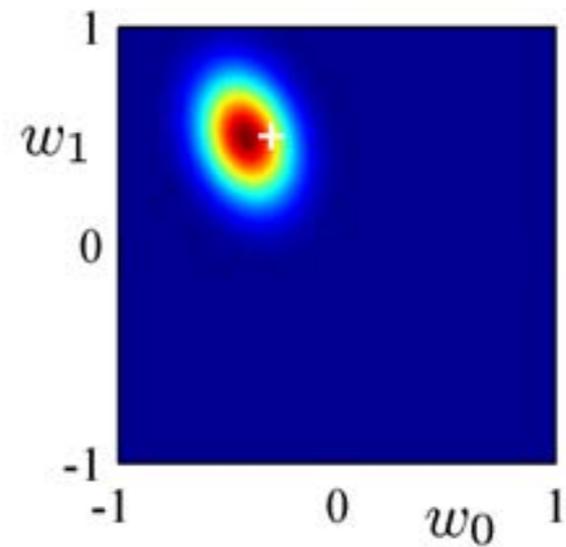
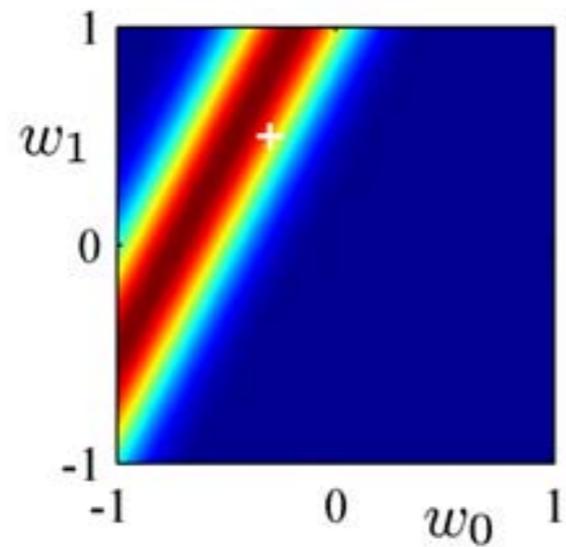
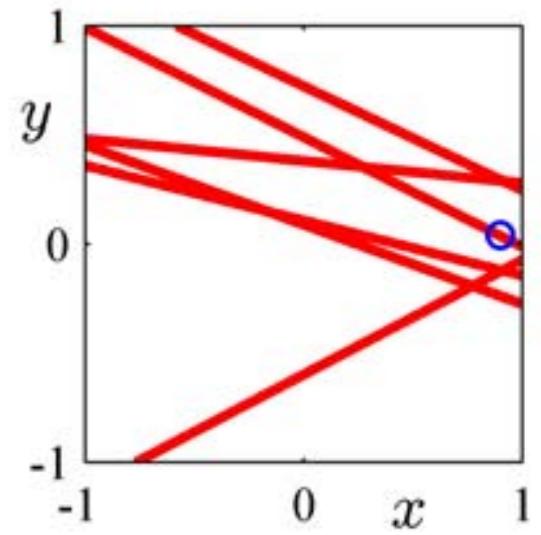
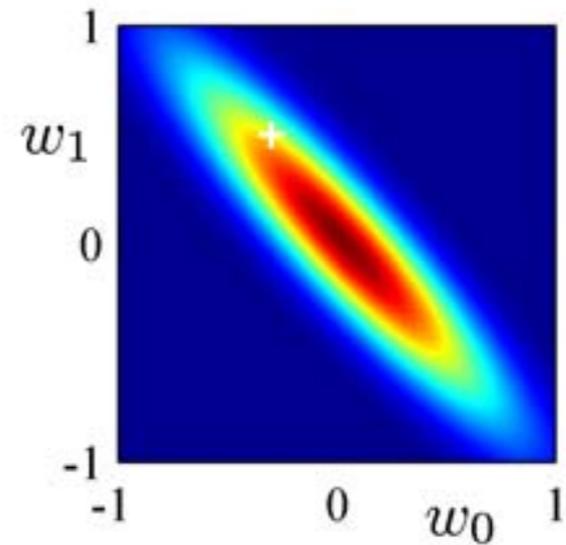
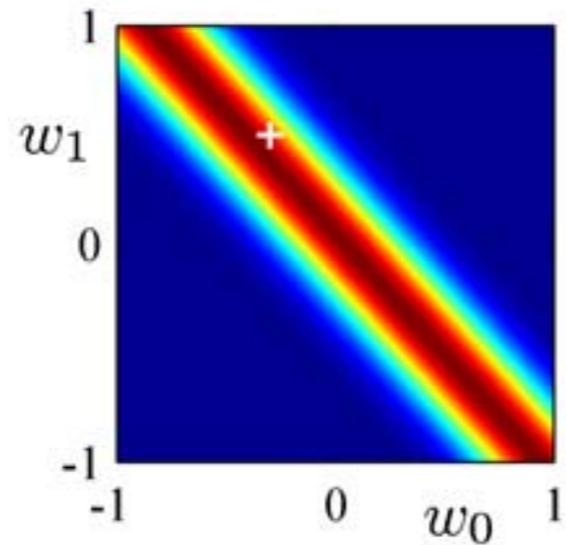
Bayesian Linear Regression

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

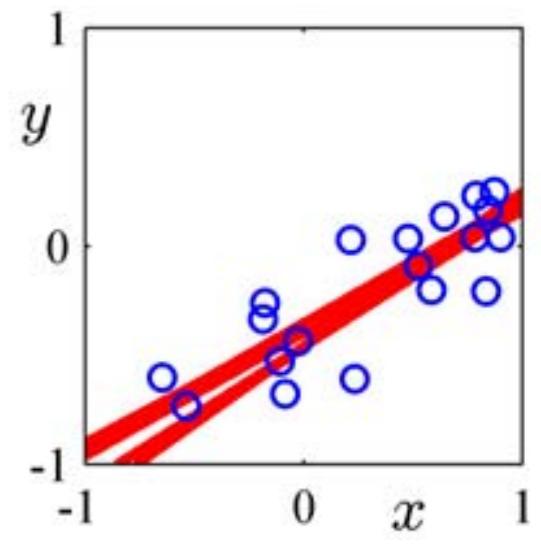
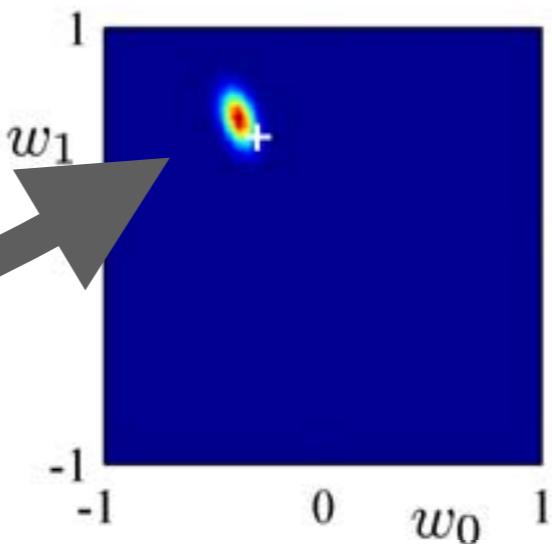
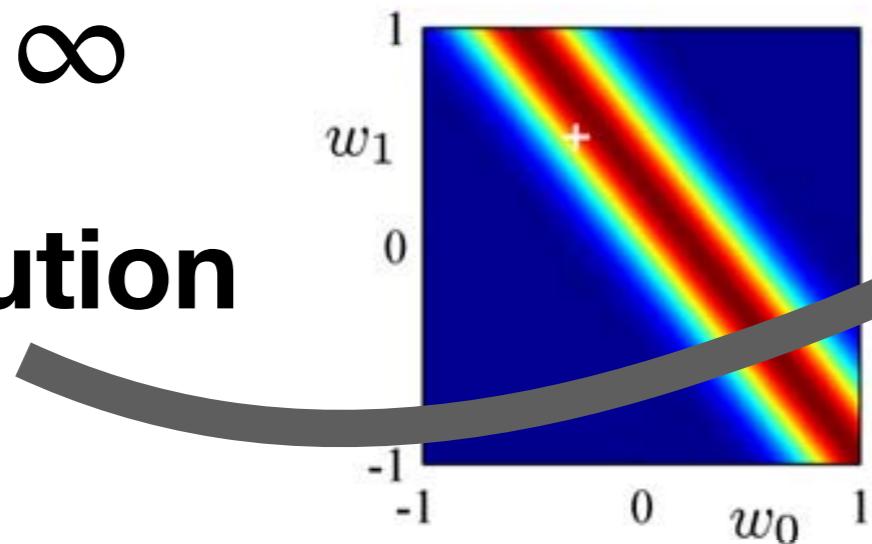
- ▶ After infinite amount of data ($N \rightarrow \infty$):

- ▶ $\lim_{N \rightarrow \infty} \mathbf{S}_N = \mathbf{0}$ (zero matrix)
- ▶ $\lim_{N \rightarrow \infty} \mathbf{m}_N = \underbrace{(\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}}_{\mathcal{W}_{ML}}$



$N \rightarrow \infty$

ML solution



Bayesian Linear Regression

- ▶ Important limits

- ▶ If $N \rightarrow \infty$

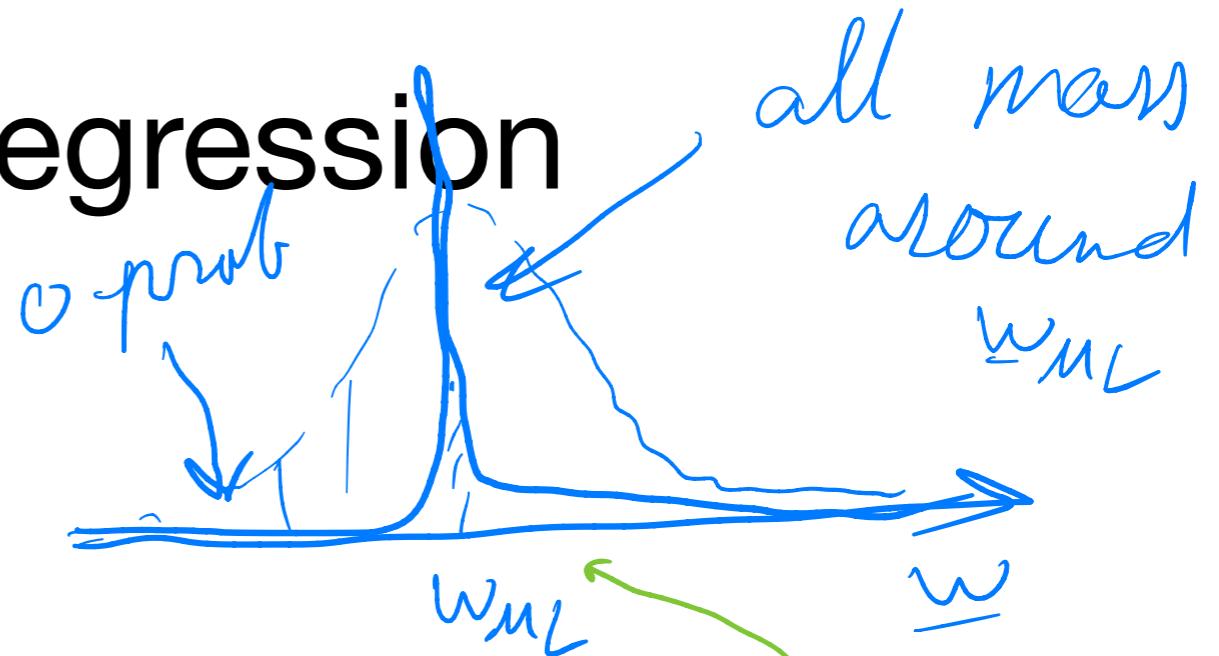
- ▶ MAP solution converges to ML solution!

- ▶ Bayesian predictive distribution

$$p(t' | x', D) = \int p(t' | x', \mathbf{w}) p(\mathbf{w} | D) d\mathbf{w}$$

converges to?

$$= p(t' | x', w_{ML})$$



posterior
 $\sim \delta_{w_{ML}}^{(w)}$ when
 $N \rightarrow \infty$

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)

$$p(\mathbf{w} | \mathbf{X}, \mathbf{t}, \alpha, \beta) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N) \quad \text{with } \mathbf{m}_N = \beta \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t}$$

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

- Parametrized Gaussian predictive distribution:

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

- Gaussian Bayesian predictive distribution for new input

$$p(t' | \mathbf{x}', \mathbf{X}, \mathbf{t}, \alpha, \beta) = \int \mathcal{N}(t' | \boldsymbol{\phi}(\mathbf{x}')^T \mathbf{w}, \beta^{-1}) \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N) d\mathbf{w}$$

Bishop Eq.
2.115

this integral is analytic and leads to yet another Gaussian

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

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

analytic depends on \mathbf{x}' expr. for σ

Predictive Distribution

active learning
 ↳ collect more data
 in uncertain regions

- ▶ 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
- ▶ Predictive distribution:

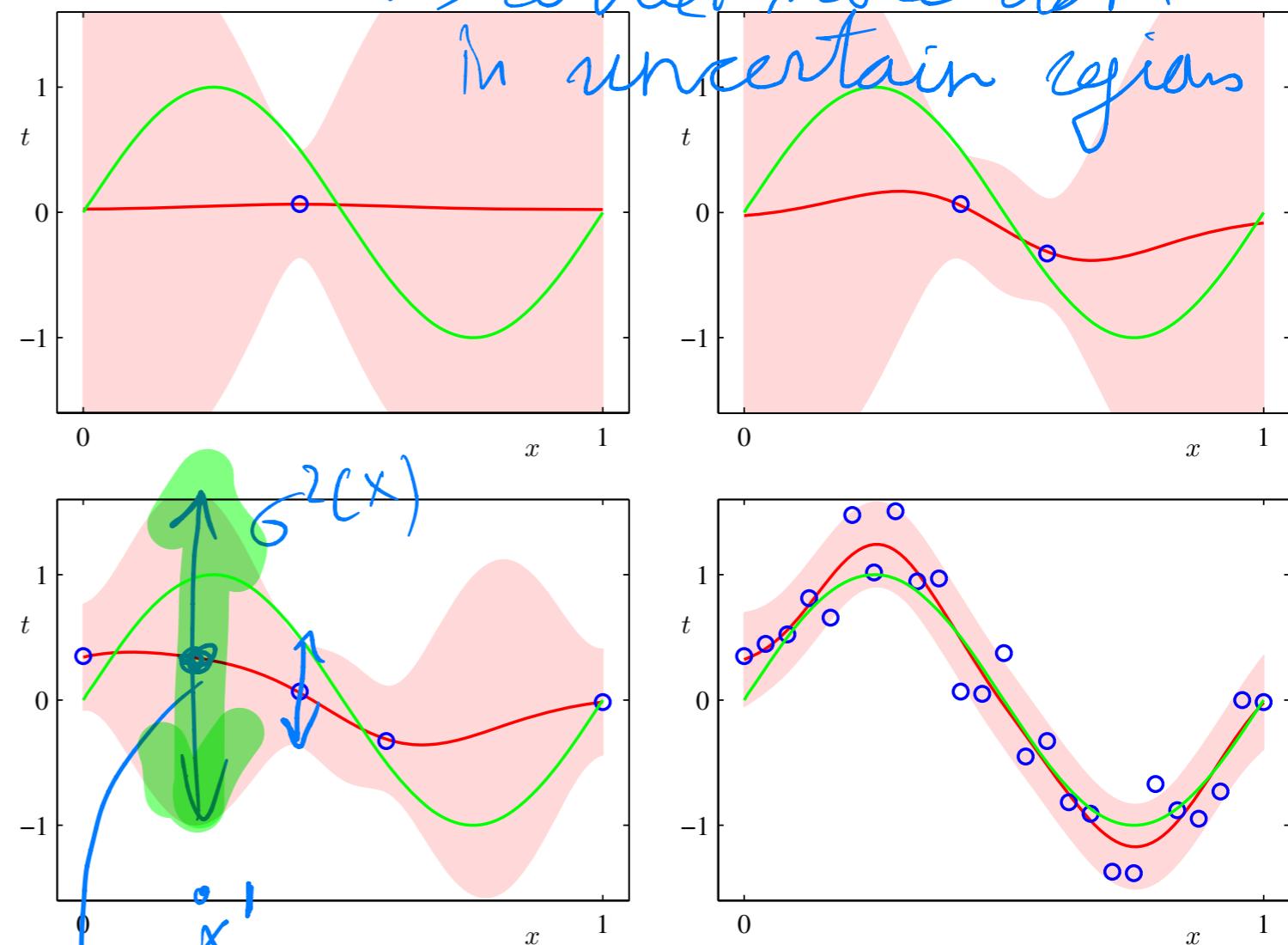


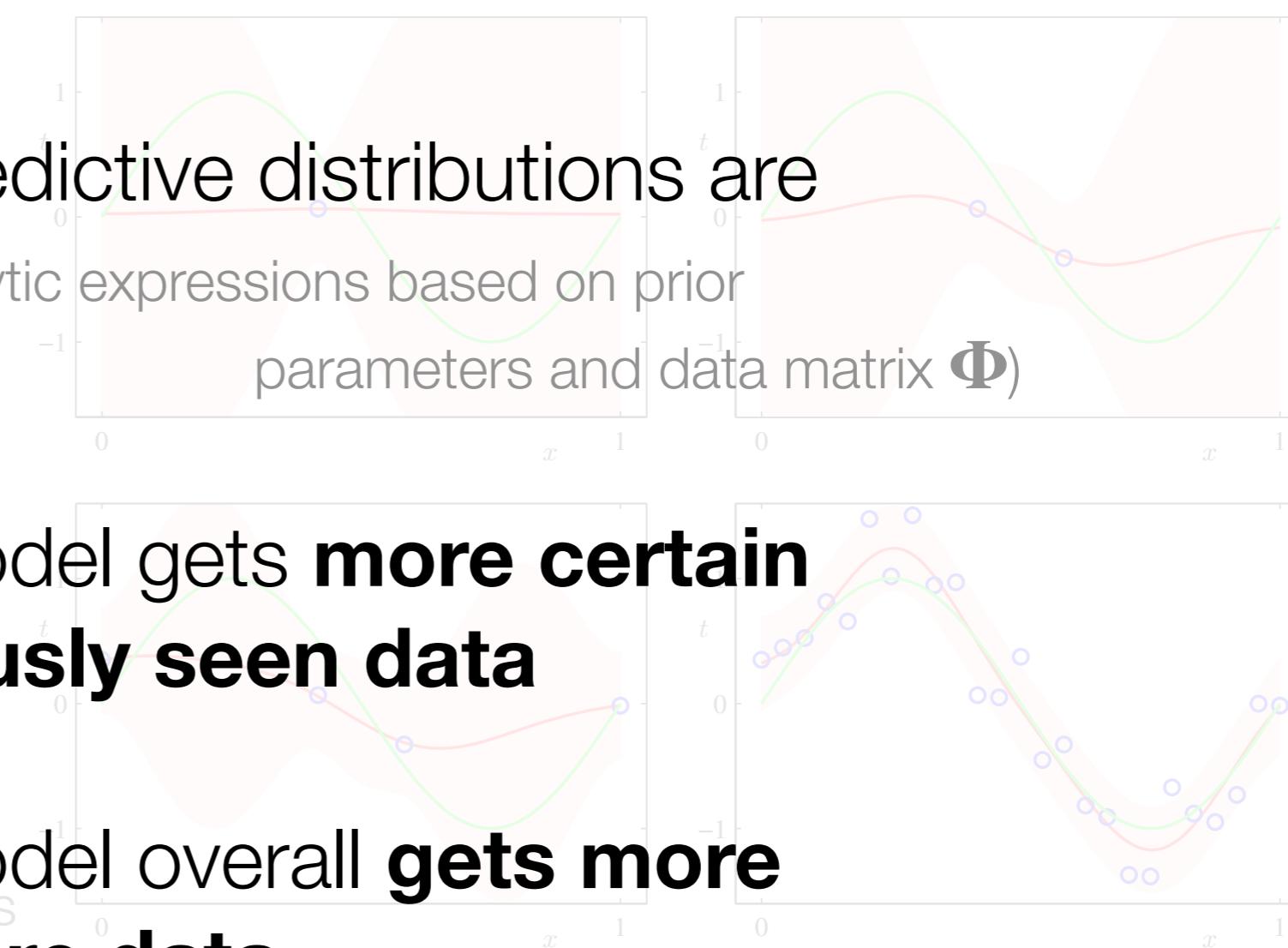
Figure: Predictive distribution (Bishop 3.8)

Bayesian

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

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

Predictive Distribution

- ▶ Datasets:
 - ▶ 1. The Bayesian predictive distributions are Gaussian (with analytic expressions based on prior parameters and data matrix Φ)
 - ▶ Dataset sizes:
 - ▶ 2. The Bayesian model gets **more certain close to previously seen data**
 - ▶ Model:
 - ▶ 3. The Bayesian model overall **gets more certain with more data**
 - ▶ Predictive distribution:
 - ▶ 4. The (mean) predictions themselves **get more accurate with more data**
- 
- Figure: Predictive distribution (Bishop 3.8)
- ▶ $y(x, w) = \phi(x)^T w$
 - ▶ $\phi_j(x)$: Gaussian basis functions
 - ▶ $\sigma_N^2(\mathbf{x}') = \beta^{-1} + \phi(\mathbf{x}')^T \mathbf{S}_N \phi(\mathbf{x}')$, $\mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{t}$, $\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi$

Samples drawn from Bayesian Predictive Distribution

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

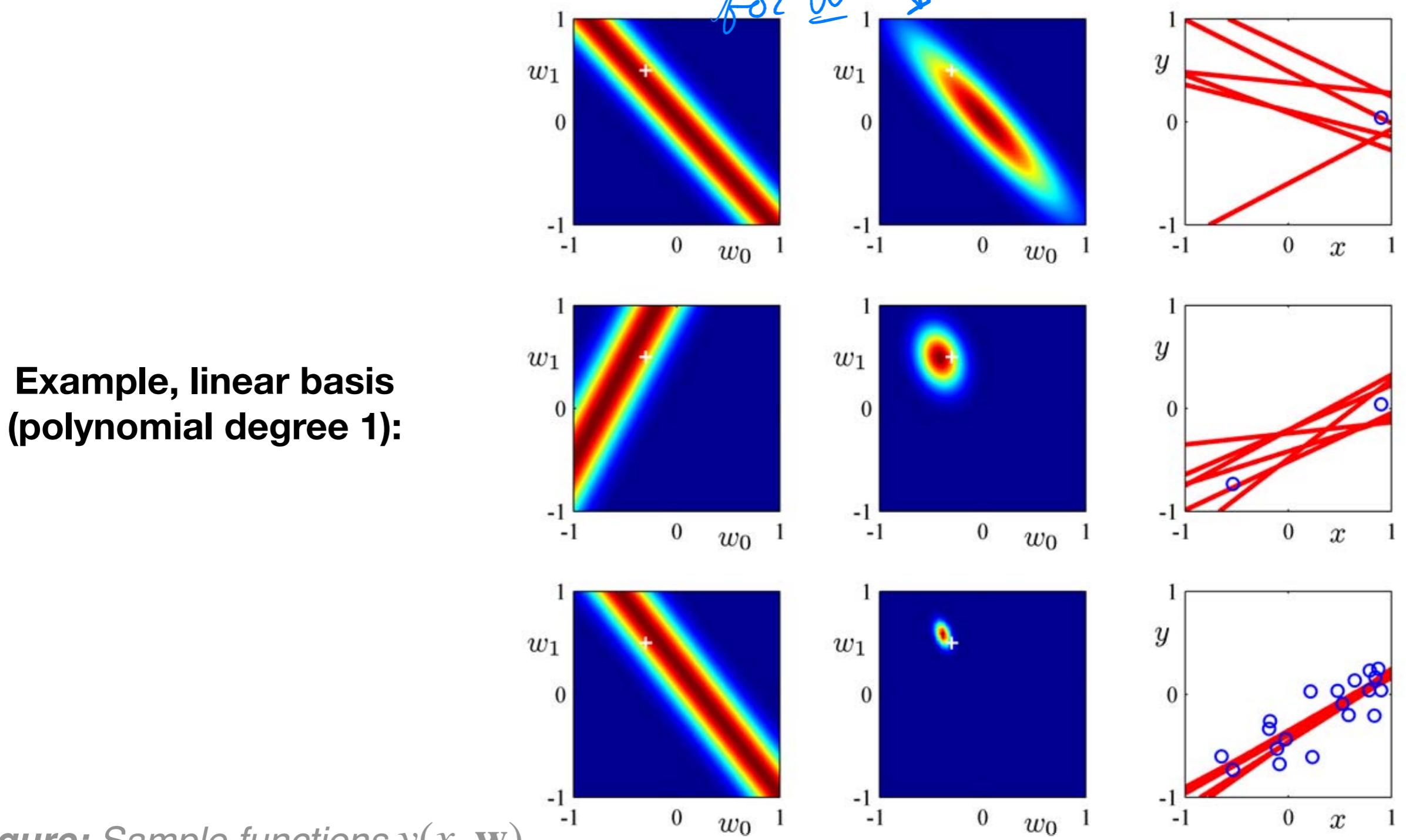


Figure: Sample functions $y(x, \mathbf{w})$ will be sampled from posterior distribution (Diagram 5.9)

Samples drawn from Bayesian Predictive Distribution

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

$$\text{with } \mathbf{m}_N = \beta \mathbf{S}_N \Phi^T \mathbf{t}$$

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

$y(x, w)$, $w \sim p(w | X, t, \alpha, \beta)$

**Example, linear model
w basis functions:**

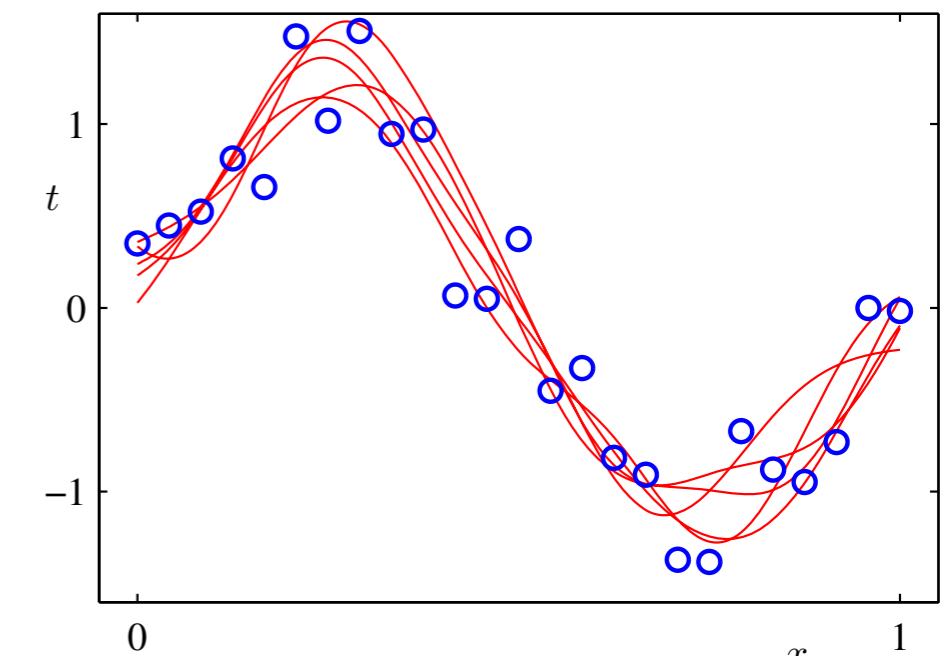
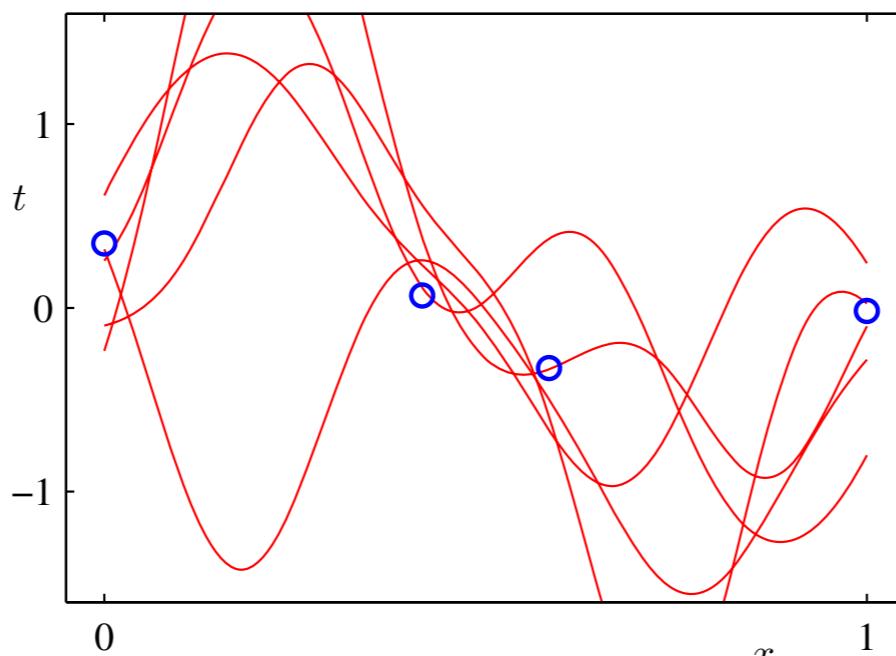
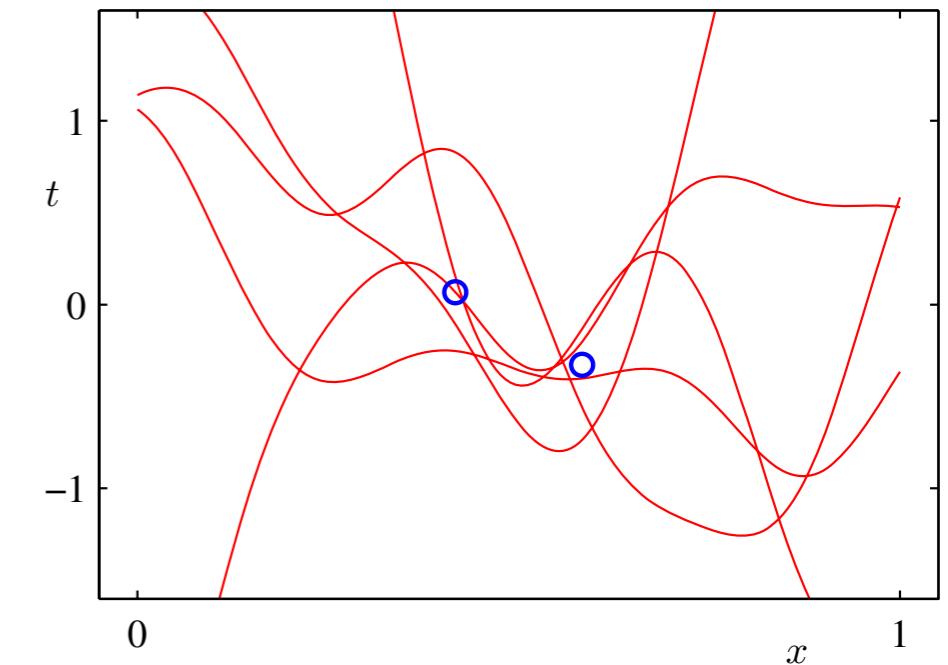
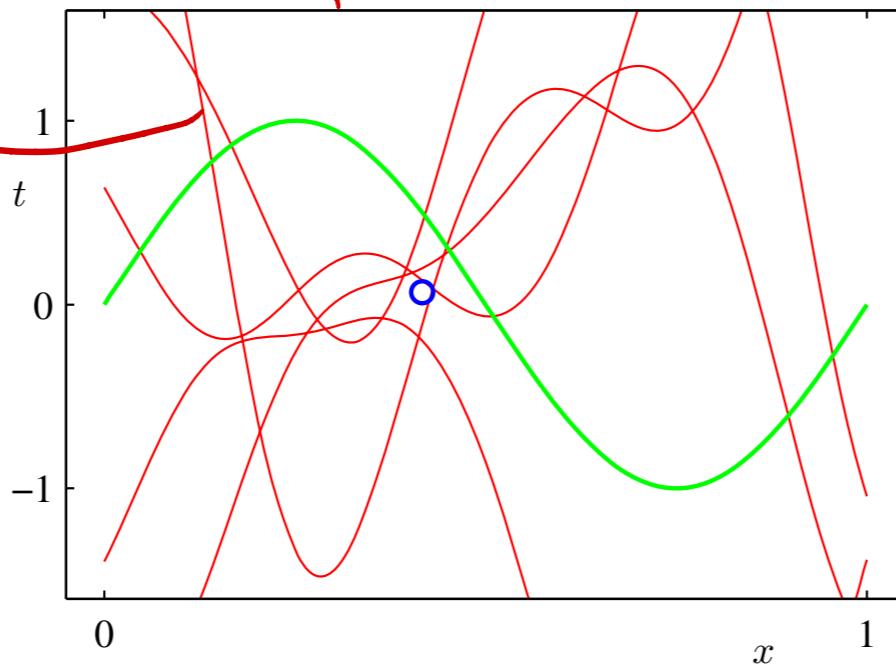


Figure: Sample functions $y(x, w)$ with w sampled from posterior distribution (Bishop 3.9)

Machine Learning 1

Lecture 5.1 - Supervised Learning

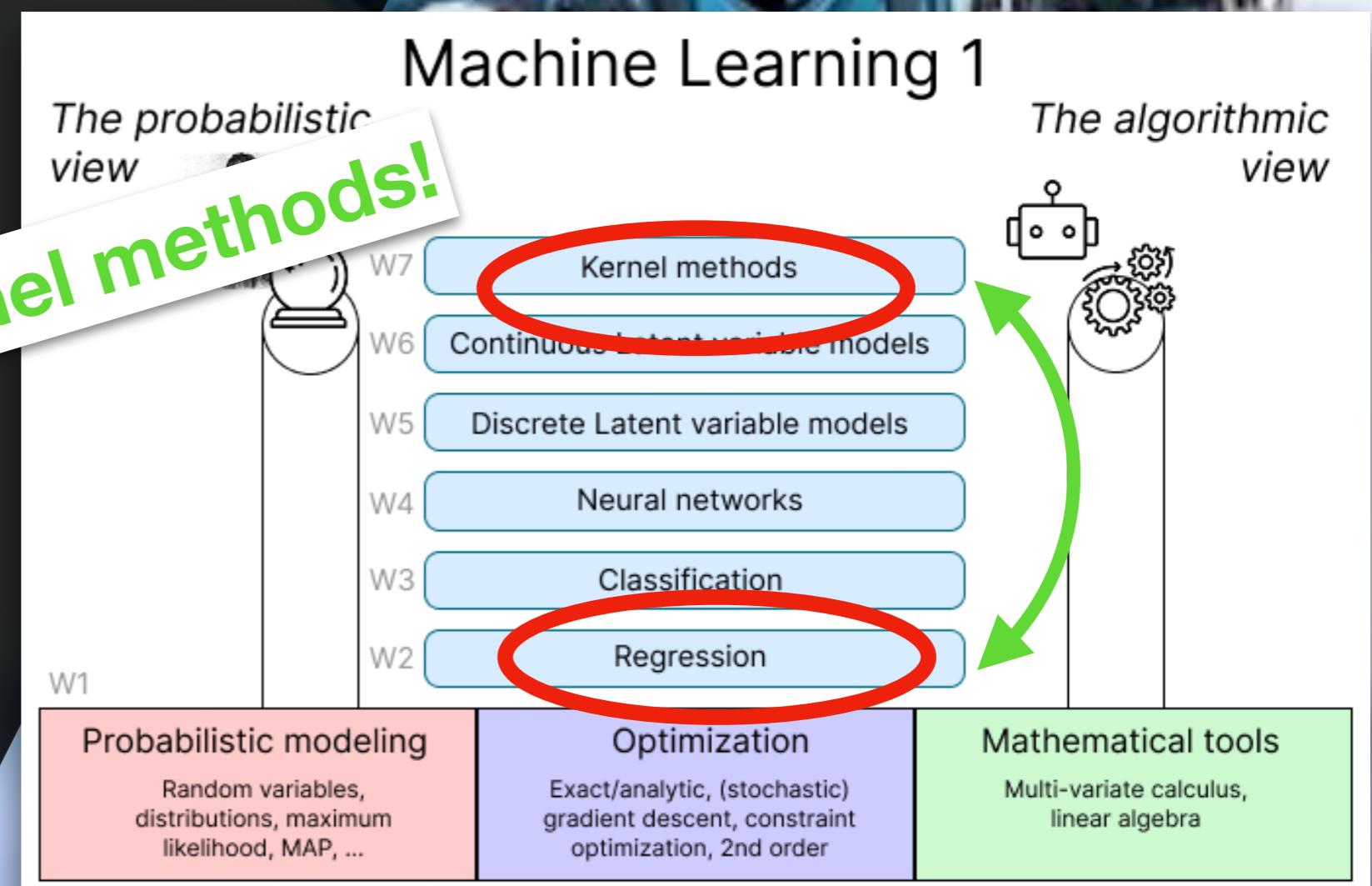
Bayesian Linear Regression - **The Equivalent Kernel**

Erik Bekkers

(Bishop 3.3.3)

Sneak preview to kernel methods!

for now understand that predictions can be made as a "weighted sum" of the data, without any explicit parameters



Equivalent Kernel Formulation

- predictive distribution

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

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

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

Equivalent Kernel Formulation

- predictive distribution

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

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

- predictive mean:

$$w = W_{MAP} = m_N$$

$$y(x', \mathbf{m}_N) = \boldsymbol{\phi}(x')^T \mathbf{m}_N$$

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

Equivalent Kernel Formulation

- predictive distribution

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

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

- predictive mean:

$$y(x', \mathbf{m}_N) = \boldsymbol{\phi}(x')^T \mathbf{m}_N$$

$$= \dots = \sum_{n=1}^N k(x', x_n) t_n$$

kernel provides weight

A weighted sum of training samples, no parameters!

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

Equivalent Kernel Formulation

- predictive distribution

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

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

- predictive mean:

$$y(x', \mathbf{m}_N) = \boldsymbol{\phi}(x')^T \mathbf{m}_N = \beta \boldsymbol{\phi}(x')^T \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t}$$

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

Equivalent Kernel Formulation

- predictive distribution

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

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

- predictive mean:

$$y(x', \mathbf{m}_N) = \boldsymbol{\phi}(x')^T \mathbf{m}_N = \beta \boldsymbol{\phi}(x') \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t} = \beta \boldsymbol{\phi}(x')^T \mathbf{S}_n \sum_{n=1}^N \boldsymbol{\Phi}_{n,:}^T t_n$$

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

Equivalent Kernel Formulation

- predictive distribution

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

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

- predictive mean:

$$y(x', \mathbf{m}_N) = \boldsymbol{\phi}(x')^T \mathbf{m}_N = \beta \boldsymbol{\phi}(x') \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t} = \beta \boldsymbol{\phi}(x')^T \mathbf{S}_N \sum_{n=1}^N \boldsymbol{\Phi}_{n,:}^T t_n$$

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

$$= \sum_{n=1}^N \beta \underbrace{\boldsymbol{\phi}(x')^T \mathbf{S}_N \boldsymbol{\phi}(x_n)}_{\text{kernel}} t_n$$

- Equivalent kernel

$$k(x', x) = \beta \boldsymbol{\phi}(x')^T \mathbf{S}_N \boldsymbol{\phi}(x_n)$$

Equivalent kernel for Gaussian Basis Functions

Figure: Equivalent kernel $k(x', x)$ (Bishop 3.10)

- Localized kernel

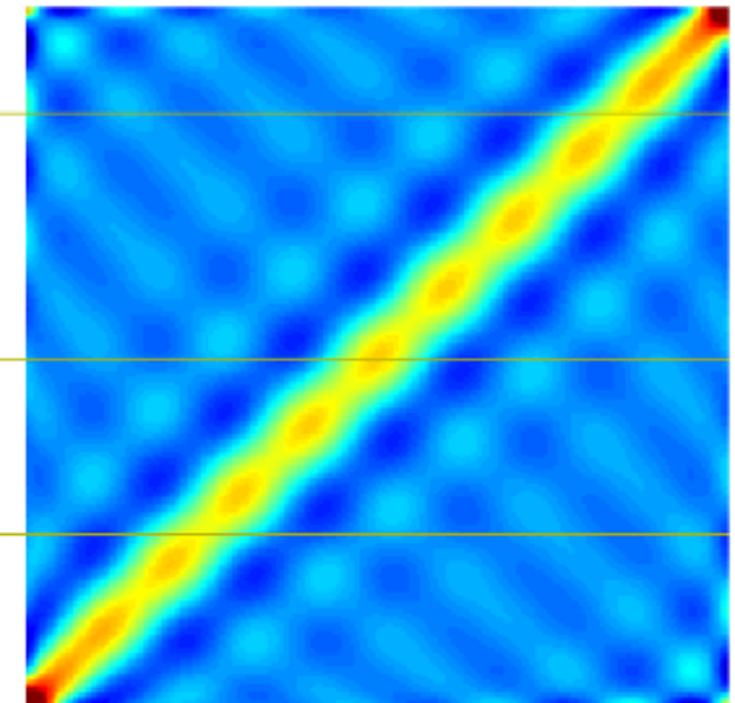
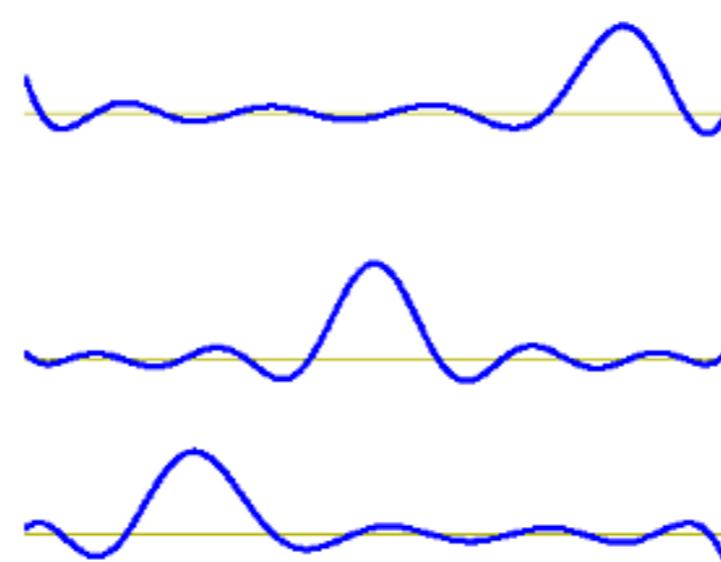
$$k(x', x) = \beta \phi(x')^T \mathbf{S}_N \phi(x)$$

- predictive mean

$$y(x', \mathbf{m}_N) = \sum_{n=1}^N k(x', x_n) t_n$$

- Training points x_n close to x' contribute more!

The kernel describes how two predictions co-vary



$$\text{cov}[t_1, t_2 | x_1, x_2] = \text{cov}_{\mathbf{w}}[y(x_1, \mathbf{w}), y(x_2, \mathbf{w})] = \text{cov}_{\mathbf{w}}[\phi(x_1)^T \mathbf{w}, \mathbf{w}^T \phi(x_2)]$$

Equivalent kernel for Gaussian Basis Functions

Figure: Equivalent kernel $k(x', x)$ (Bishop 3.10)

- ▶ Localized kernel

$$k(x', x) = \beta \phi(x')^T \mathbf{S}_N \phi(x)$$

- ▶ predictive mean

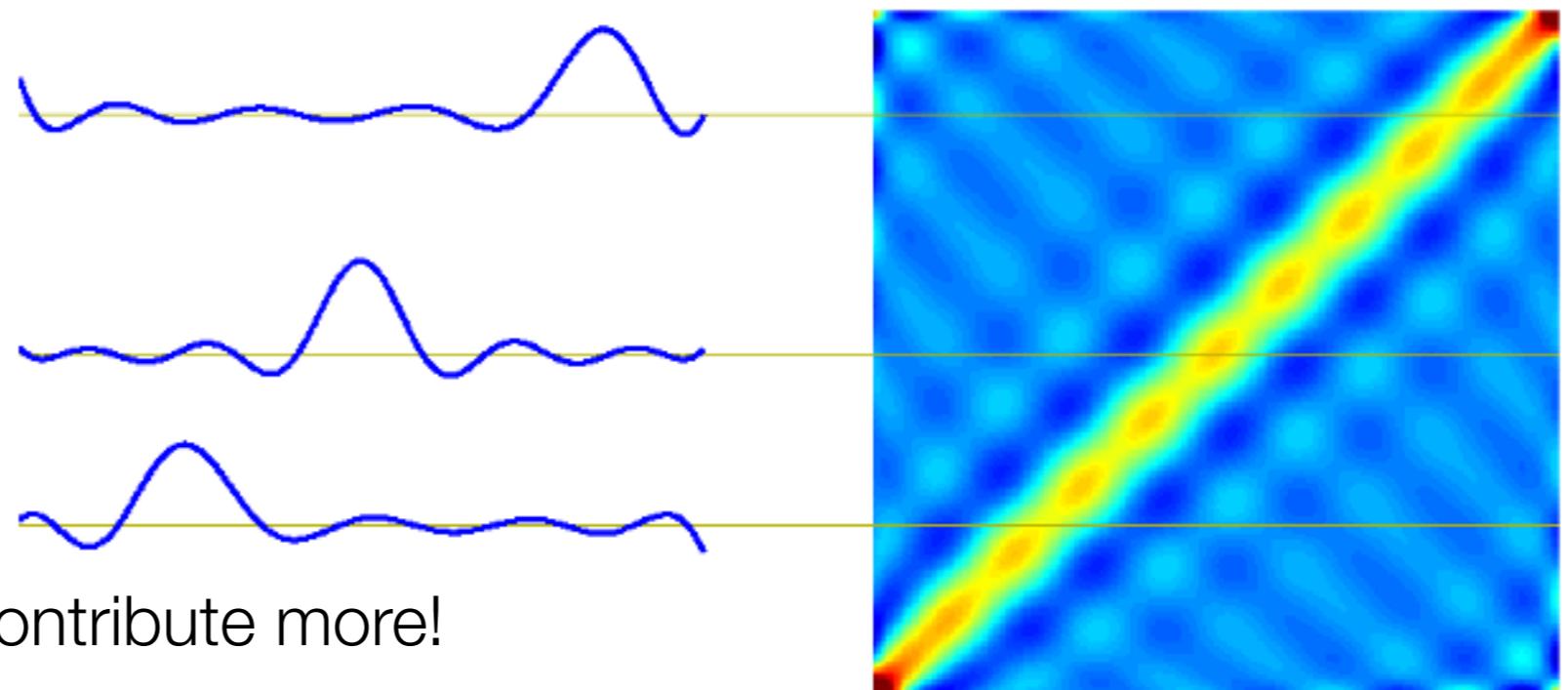
$$y(x', \mathbf{m}_N) = \sum_{n=1}^N k(x', x_n) t_n$$

- ▶ Training points x_n close to x' contribute more!

- ▶ Covariance of between predictions:

$$\text{cov}[t_1, t_2 | x_1, x_2] = \text{cov}_{\mathbf{w}}[y(x_1, \mathbf{w}), y(x_2, \mathbf{w})] = \text{cov}_{\mathbf{w}}[\phi(x_1)^T \mathbf{w}, \mathbf{w}^T \phi(x_2)]$$

$$= \mathbb{E}_{\mathbf{w}}[\phi(x_1)^T \mathbf{w} \mathbf{w}^T \phi(x_2)] - \mathbb{E}_{\mathbf{w}}[\phi(x_1)^T \mathbf{w}] \mathbb{E}_{\mathbf{w}}[\mathbf{w}^T \phi(x_2)]$$



Equivalent kernel for Gaussian Basis Functions

Figure: Equivalent kernel $k(x', x)$ (Bishop 3.10)

- Localized kernel

$$k(x', x) = \beta \phi(x')^T \mathbf{S}_N \phi(x)$$

- predictive mean

$$y(x', \mathbf{m}_N) = \sum_{n=1}^N k(x', x_n) t_n$$

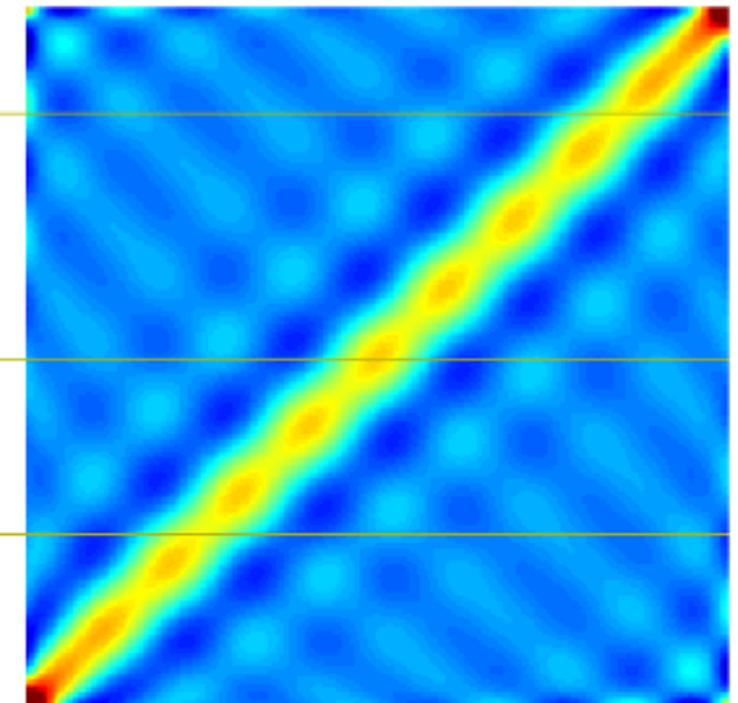
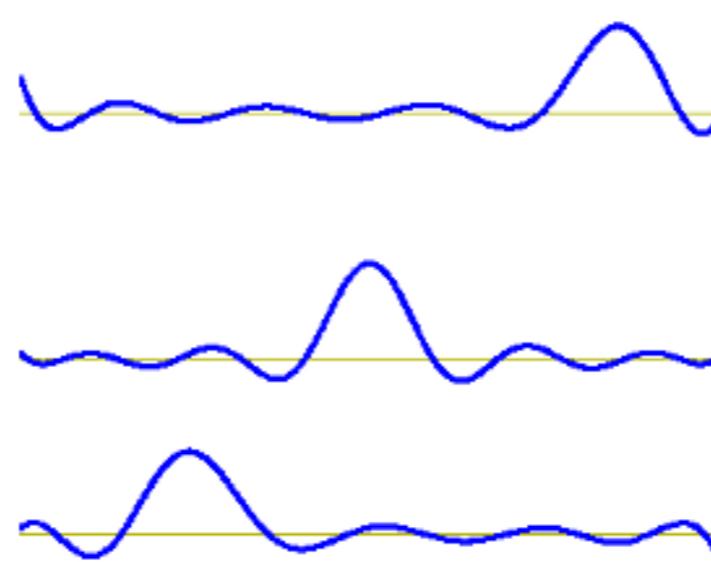
- Training points x_n close to x' contribute more!

- Covariance of between predictions:

$$\text{cov}[t_1, t_2 | x_1, x_2] = \text{cov}_{\mathbf{w}}[y(x_1, \mathbf{w}), y(x_2, \mathbf{w})] = \text{cov}_{\mathbf{w}}[\phi(x_1)^T \mathbf{w}, \mathbf{w}^T \phi(x_2)]$$

$$= \mathbb{E}_{\mathbf{w}}[\phi(x_1)^T \mathbf{w} \mathbf{w}^T \phi(x_2)] - \mathbb{E}_{\mathbf{w}}[\phi(x_1)^T \mathbf{w}] \mathbb{E}_{\mathbf{w}}[\mathbf{w}^T \phi(x_2)]$$

$$= \phi(x_1)^T (\mathbb{E}_{\mathbf{w}}[\mathbf{w} \mathbf{w}^T] - \mathbb{E}[\mathbf{w}] \mathbb{E}[\mathbf{w}]^T) \phi(x_2)$$



Equivalent kernel for Gaussian Basis Functions

Figure: Equivalent kernel $k(x', x)$ (Bishop 3.10)

- Localized kernel

$$k(x', x) = \beta \phi(x')^T \mathbf{S}_N \phi(x)$$

- predictive mean

$$y(x', \mathbf{m}_N) = \sum_{n=1}^N k(x', x_n) t_n$$

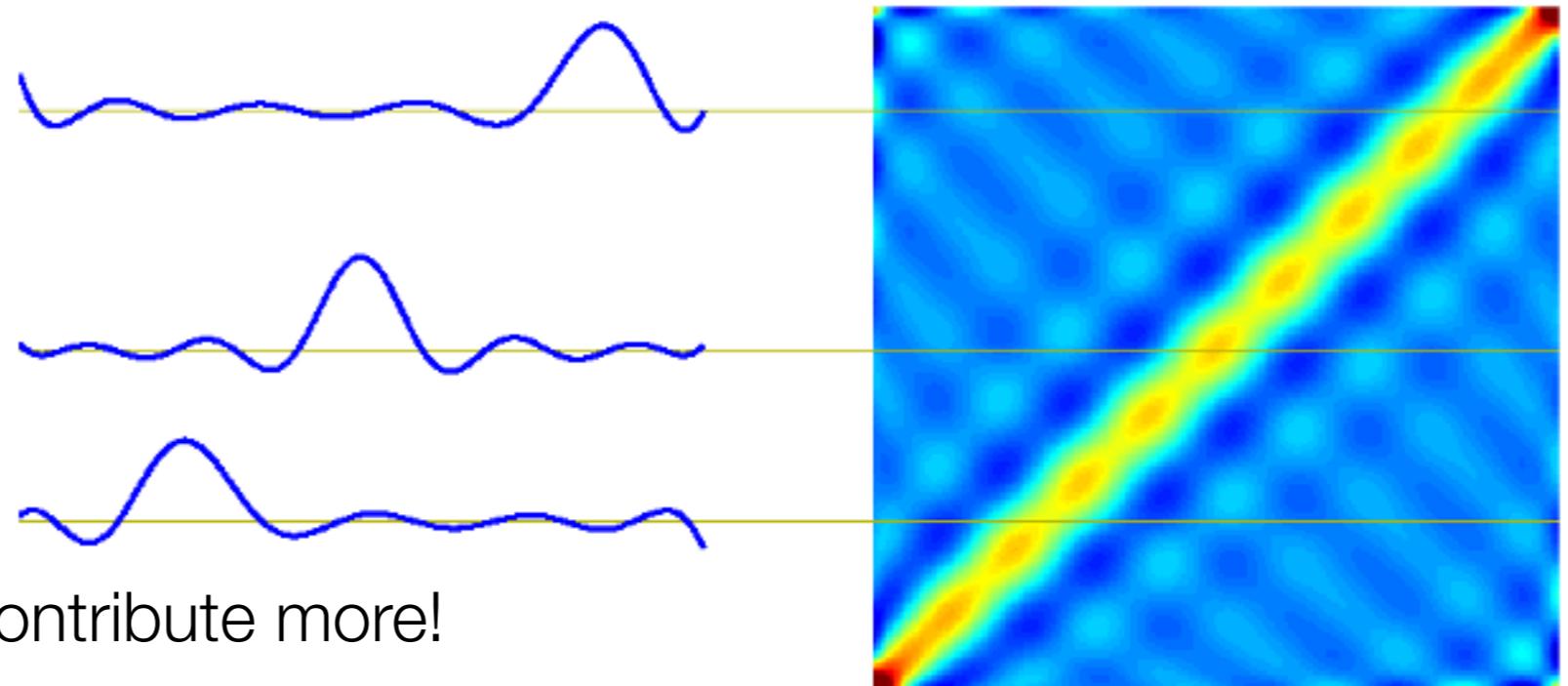
- Training points x_n close to x' contribute more!

- Covariance of between predictions:

$$\text{cov}[t_1, t_2 | x_1, x_2] = \text{cov}_{\mathbf{w}}[y(x_1, \mathbf{w}), y(x_2, \mathbf{w})] = \text{cov}_{\mathbf{w}}[\phi(x_1)^T \mathbf{w}, \mathbf{w}^T \phi(x_2)]$$

$$= \mathbb{E}_{\mathbf{w}}[\phi(x_1)^T \mathbf{w} \mathbf{w}^T \phi(x_2)] - \mathbb{E}_{\mathbf{w}}[\phi(x_1)^T \mathbf{w}] \mathbb{E}_{\mathbf{w}}[\mathbf{w}^T \phi(x_2)]$$

$$= \phi(x_1)^T (\mathbb{E}_{\mathbf{w}}[\mathbf{w} \mathbf{w}^T] - \mathbb{E}[\mathbf{w}] \mathbb{E}[\mathbf{w}]^T) \phi(x_2) = \phi(x_1)^T \text{cov}[\mathbf{w}, \mathbf{w}] \phi(x_2)$$



Equivalent kernel for Gaussian Basis Functions

Figure: Equivalent kernel $k(x', x)$ (Bishop 3.10)

- Localized kernel

$$k(x', x) = \beta \phi(x')^T \mathbf{S}_N \phi(x)$$

- predictive mean

$$y(x', \mathbf{m}_N) = \sum_{n=1}^N k(x', x_n) t_n$$

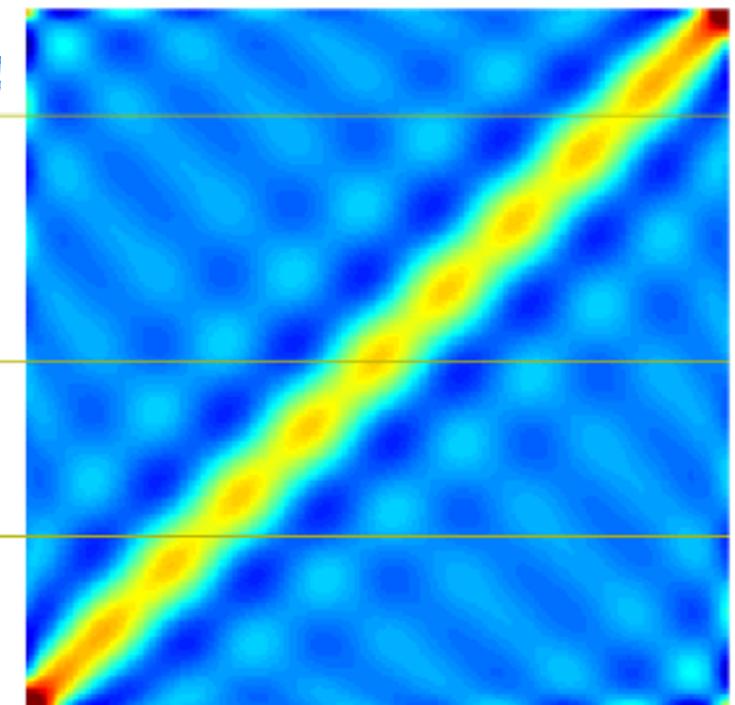
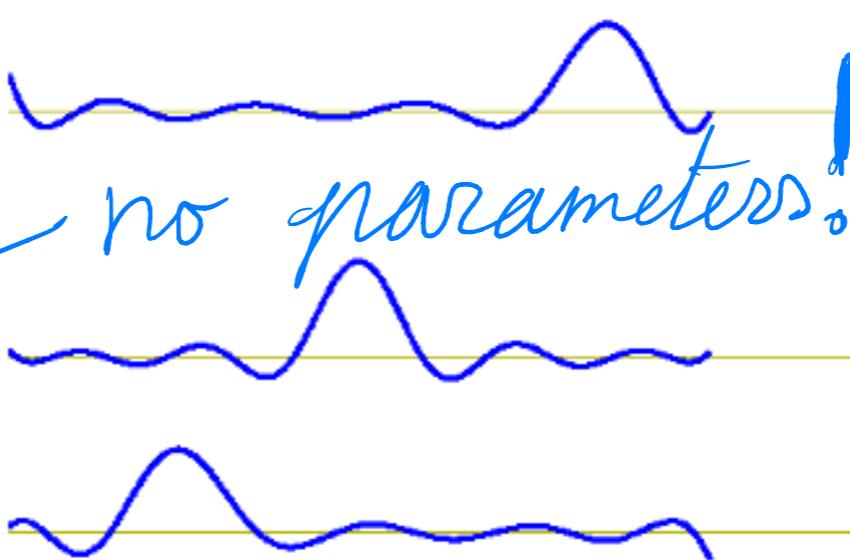
- Training points x_n close to x' contribute more!

- Covariance of between predictions:

$$\text{cov}[t_1, t_2 | x_1, x_2] = \text{cov}_{\mathbf{w}}[y(x_1, \mathbf{w}), y(x_2, \mathbf{w})] = \text{cov}_{\mathbf{w}}[\phi(x_1)^T \mathbf{w}, \mathbf{w}^T \phi(x_2)]$$

$$= \mathbb{E}_{\mathbf{w}}[\phi(x_1)^T \mathbf{w} \mathbf{w}^T \phi(x_2)] - \mathbb{E}_{\mathbf{w}}[\phi(x_1)^T \mathbf{w}] \mathbb{E}_{\mathbf{w}}[\mathbf{w}^T \phi(x_2)]$$

$$= \phi(x_1)^T (\mathbb{E}_{\mathbf{w}}[\mathbf{w} \mathbf{w}^T] - \mathbb{E}[\mathbf{w}] \mathbb{E}[\mathbf{w}]^T) \phi(x_2) = \phi(x_1)^T \text{cov}[\mathbf{w}, \mathbf{w}] \phi(x_2) = \phi(x_1)^T \mathbf{S}_N \phi(x_2)$$



l-g. covariance matrix
of posterior $p(\mathbf{w} | \mathbf{D})$
that was derived before
in the Gaussian case

Machine Learning 1

Lecture 5.2 - Supervised Learning
Bayesian Linear Regression - Bayesian Model Comparison

Erik Bekkers

(Bishop 3.4)



Machine Learning 1

Lecture 5.3 - Supervised Learning
Bayesian Linear Regression - Approximating
the Model Evidence

Erik Bekkers

(Bishop 3.5.0)

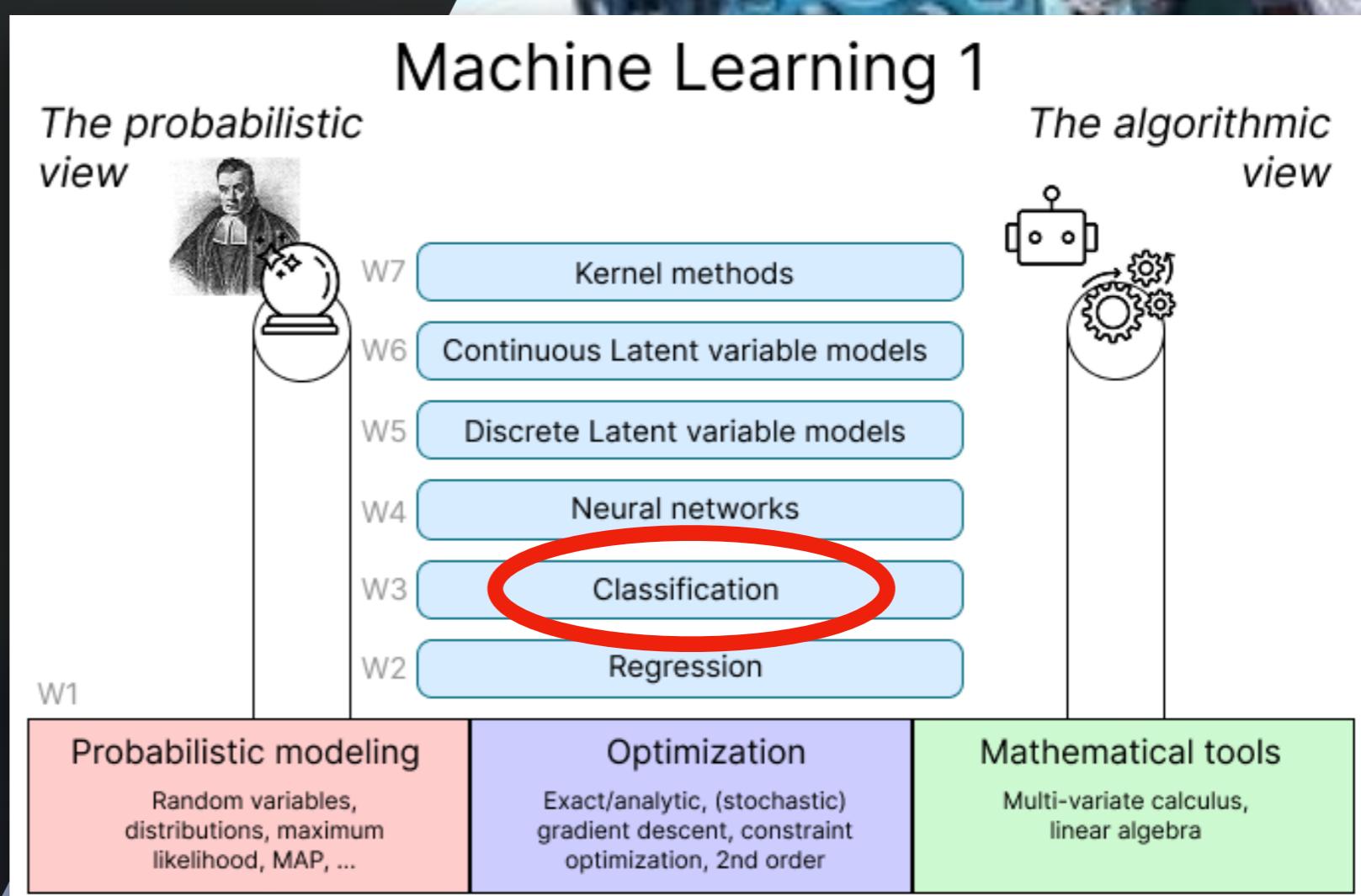


Machine Learning 1

Lecture 5.4 - Supervised Learning
Classification - Decision Regions

Erik Bekkers

(Bishop 1.5, 4.1)



Classification through decision regions

- Input: $\mathbf{x} = (x_1, \dots, x_D)^T \in \mathbb{R}^D$

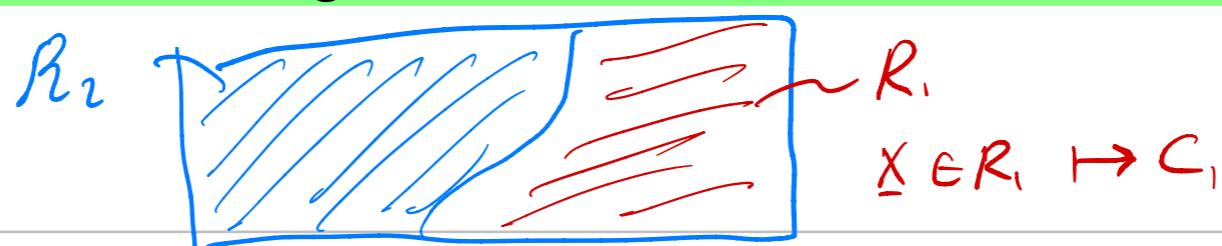
- Target:

- 2-class targets: $t \in \{C_1, C_2\} = \{"cat", "dog"\}$
- Multi-class targets $t \in \{C_1, C_2, \dots, C_K\}$

Strategy: One-hot encoding
a vector with only a 1 at the class-index but all other elements zero

- Divide input space \mathbb{R}^D into K decision regions R_k
- Assign each decision region to a class C_k
- Boundaries of decision regions are called *decision boundaries/surfaces*.

if x falls in R_2
($x \in R_2$) then
classify as C_2 ($x \mapsto C_2$)
"maps to"



$$\underline{t} = \begin{pmatrix} t_1 \\ t_2 \\ \vdots \\ t_K \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \end{pmatrix} \in \{0, 1\}^K$$

is a one-hot encoding

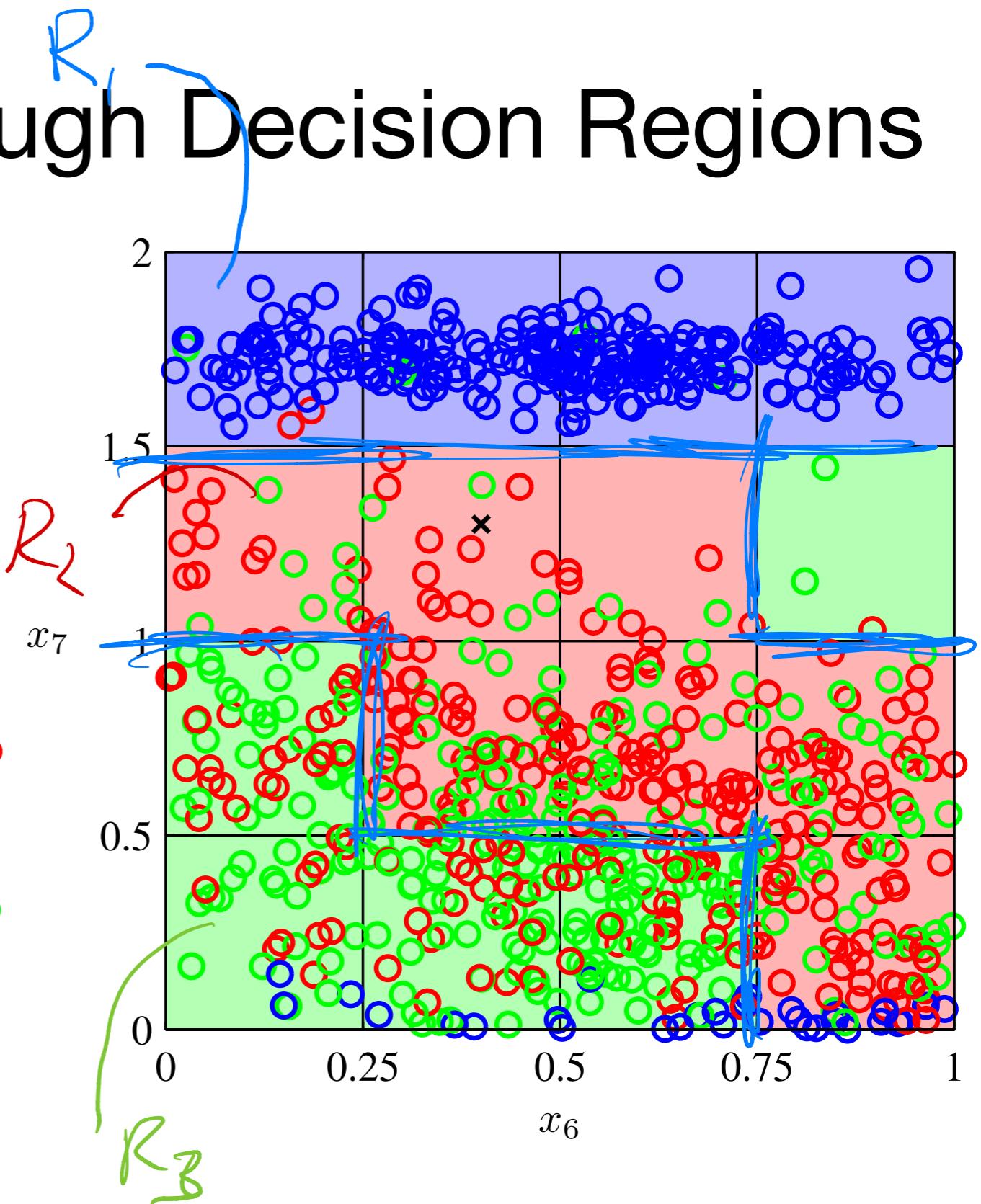
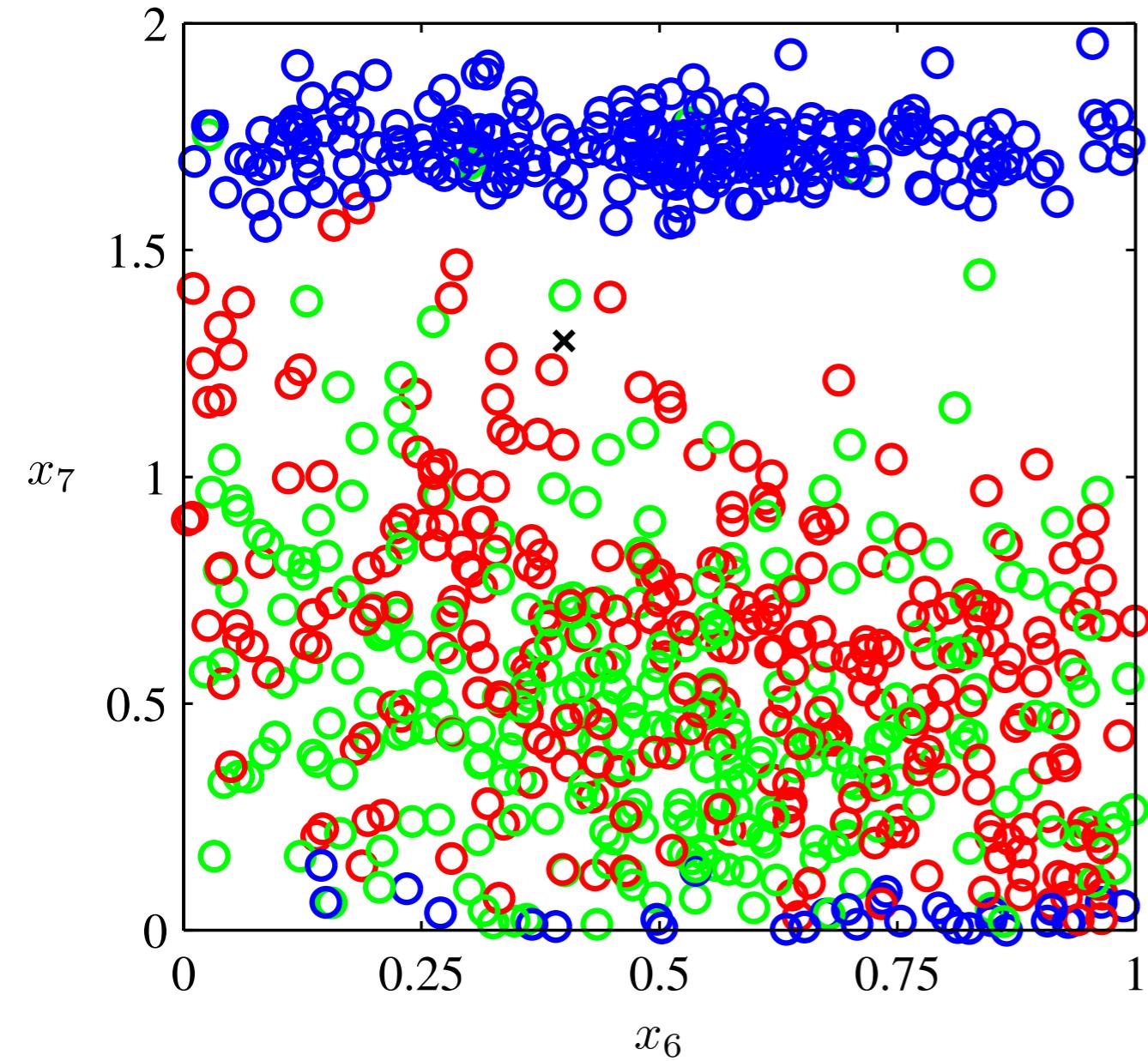
for class C_k if $t_i = 1$

if $i = k$ and zero otherwise.

In the above \underline{t} is a one-hot encoding for class C_2 since $t_2 = 1$ and all other entries are zero.

Such a one-hot encoding can be thought of as a categorical distribution assigning probability 1 to the corresponding class and zero to others $p(\text{class} = C_k) = 1$.

Classification through Decision Regions



Figures: 3 class problem with decision boundaries. (Bishop 1.19 & 1.20)

Linear Classification

- Linear Classification: consider only *linear* decision boundaries
- For D -dimensional input space: $\mathbf{x} \in \mathbb{R}^D$
decision surface is a $D - 1$ dimensional hyperplane
- Datasets whose classes can be separated exactly by linear decision surfaces are called **linearly separable**

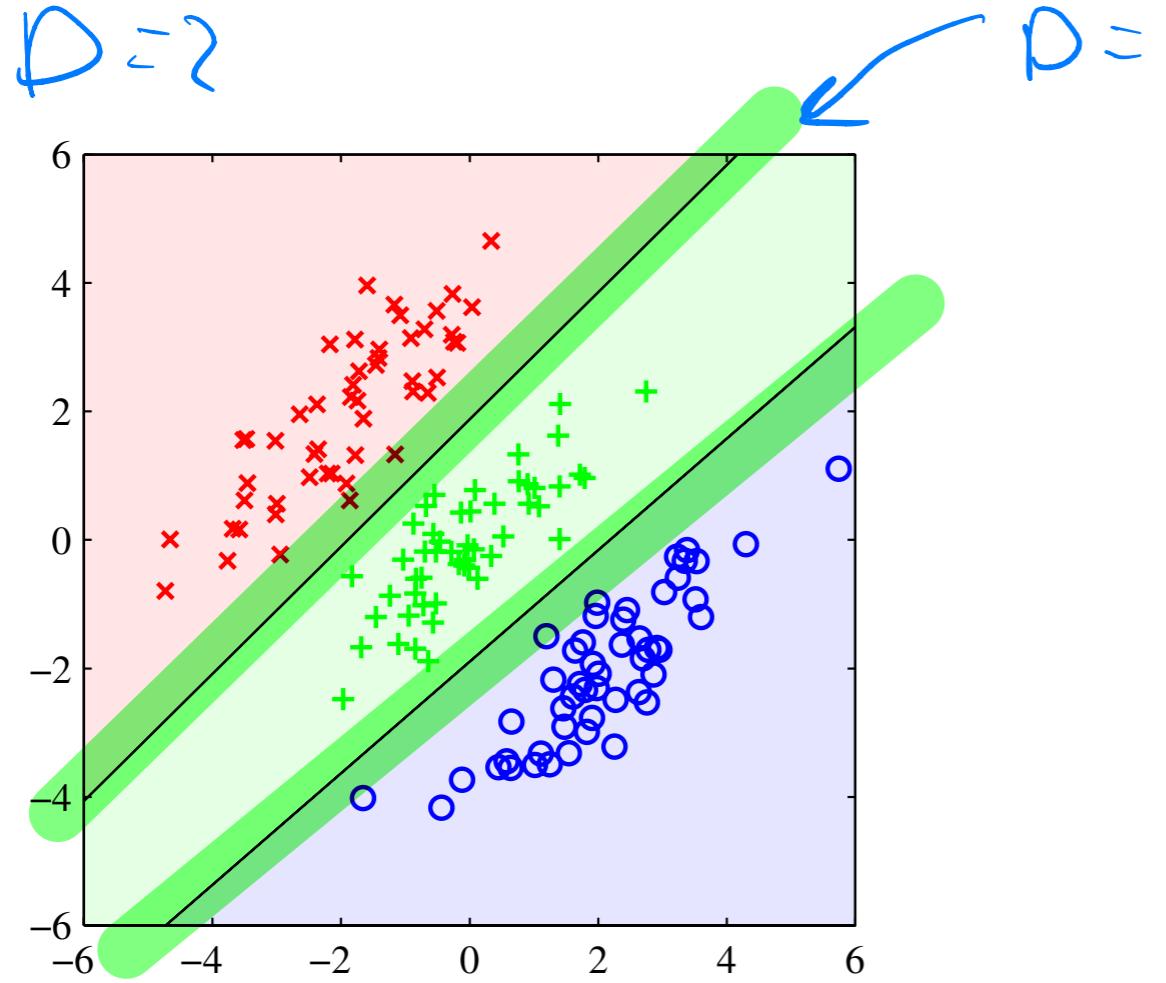


Figure: Linearly separable dataset (Bishop 4.5)

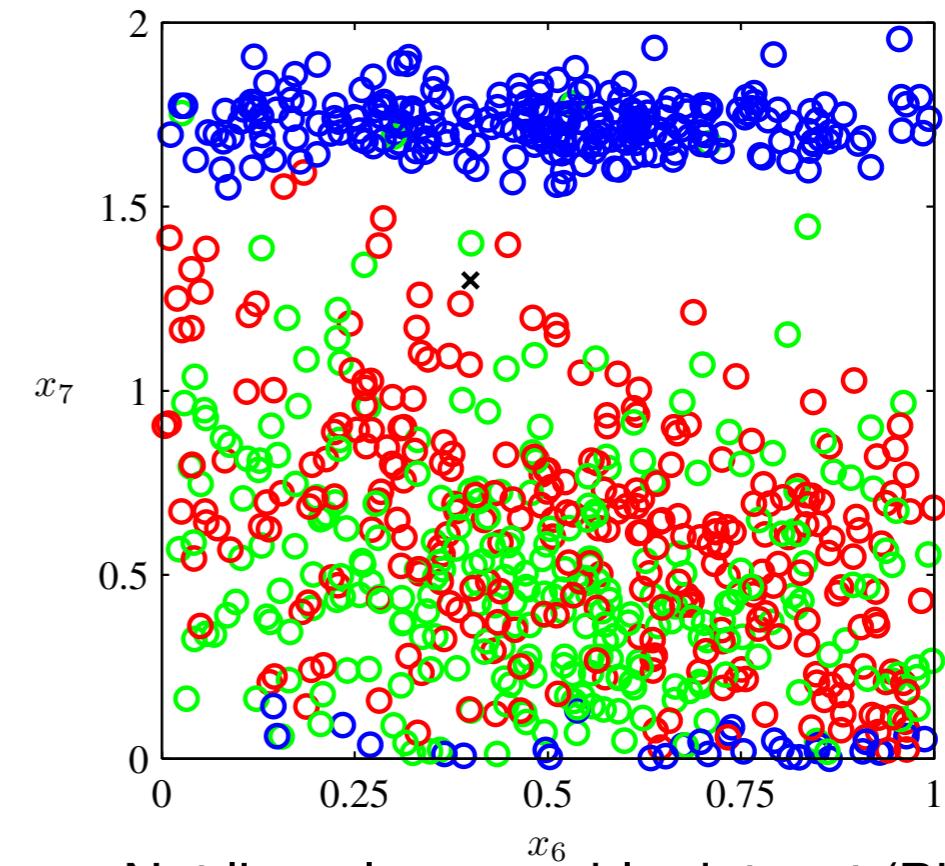


Figure: Not linearly separable dataset (Bishop 1.19)

Multiple Classes ($K > 2$)

- $K = 2$ classes:
 - 1 classifier determines
- Multiple classes: $K > 2$
 - $K - 1$ classifiers:
 - One-versus-the-rest
 - using multiple binary classifiers

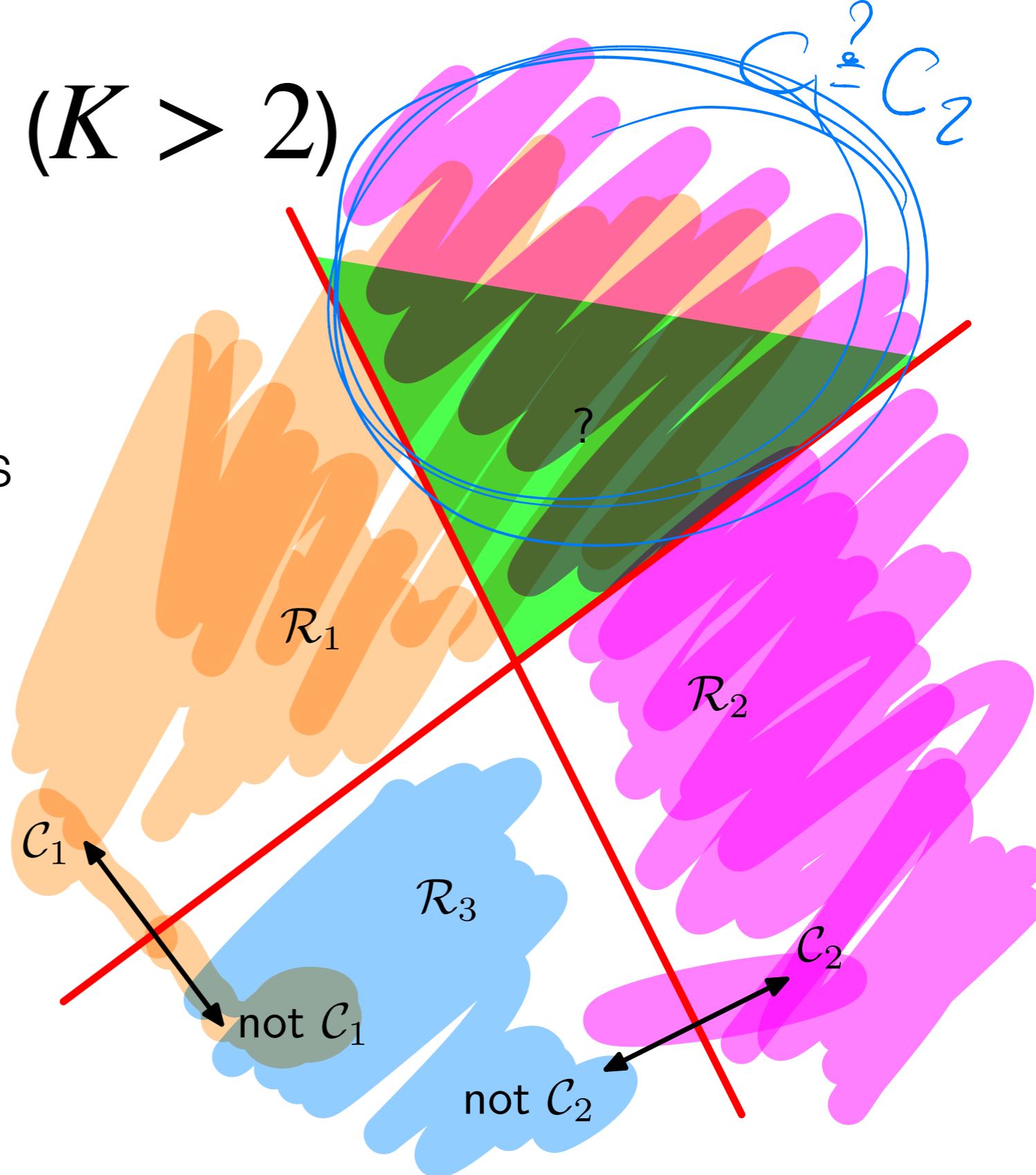
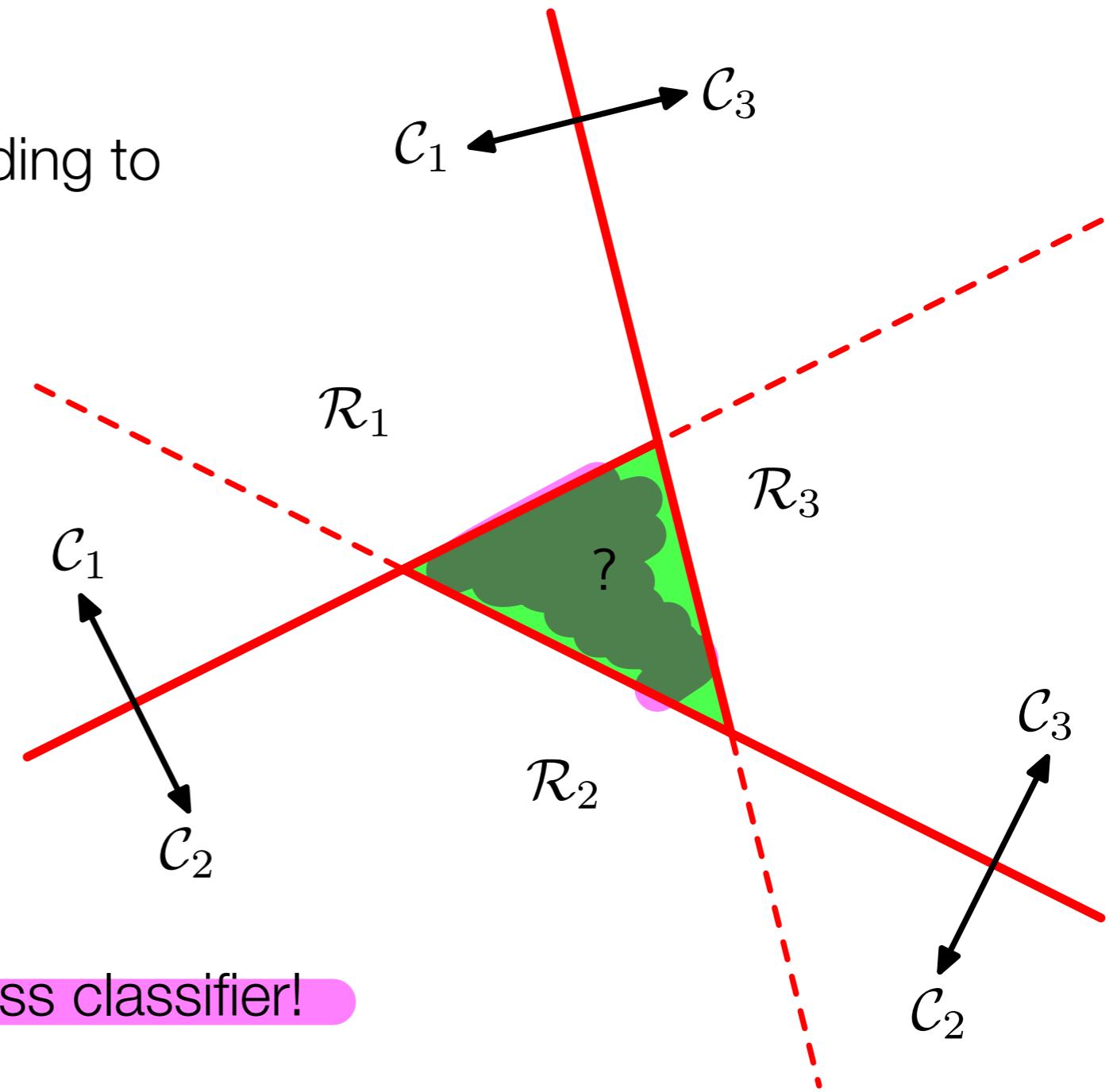


Figure: one-versus-the-rest classifiers (Bishop 4.2)

Multiple Classes ($K > 2$)

- $K(K - 1)/2$ classifiers:
- Points are classified according to majority vote of classifiers
- one-versus-one



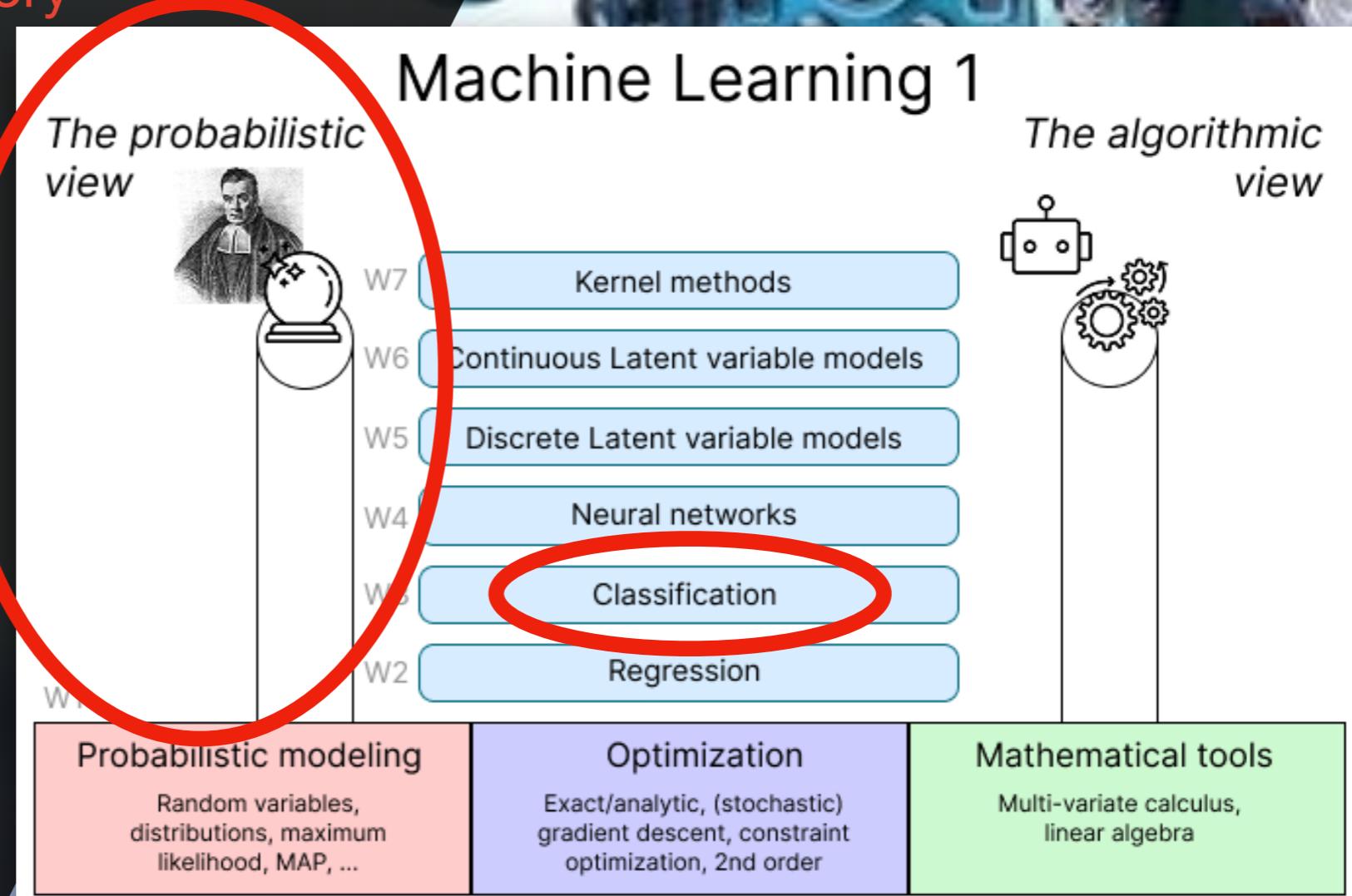
- **Solution:** Make one K -class classifier!
(See later)

Figure: one-versus-one classifiers (Bishop 4.2)

Machine Learning 1

Lecture 5.5 - Supervised Learning
Classification - **Decision Theory**

Erik Bekkers
(Bishop 1.5)



Decision theory

- Dataset: Input vectors $\mathbf{x} \in \mathbb{R}^D$, ground truth targets $t \in \{C_1, \dots, C_k\}$
- Divide input space \mathbb{R}^D into K decision regions R_k
- Every observed datapoint: $\begin{cases} \text{ground truth} \\ \text{prediction} \end{cases}$
- Confusion matrix: ground truth classes vs. predicted classes

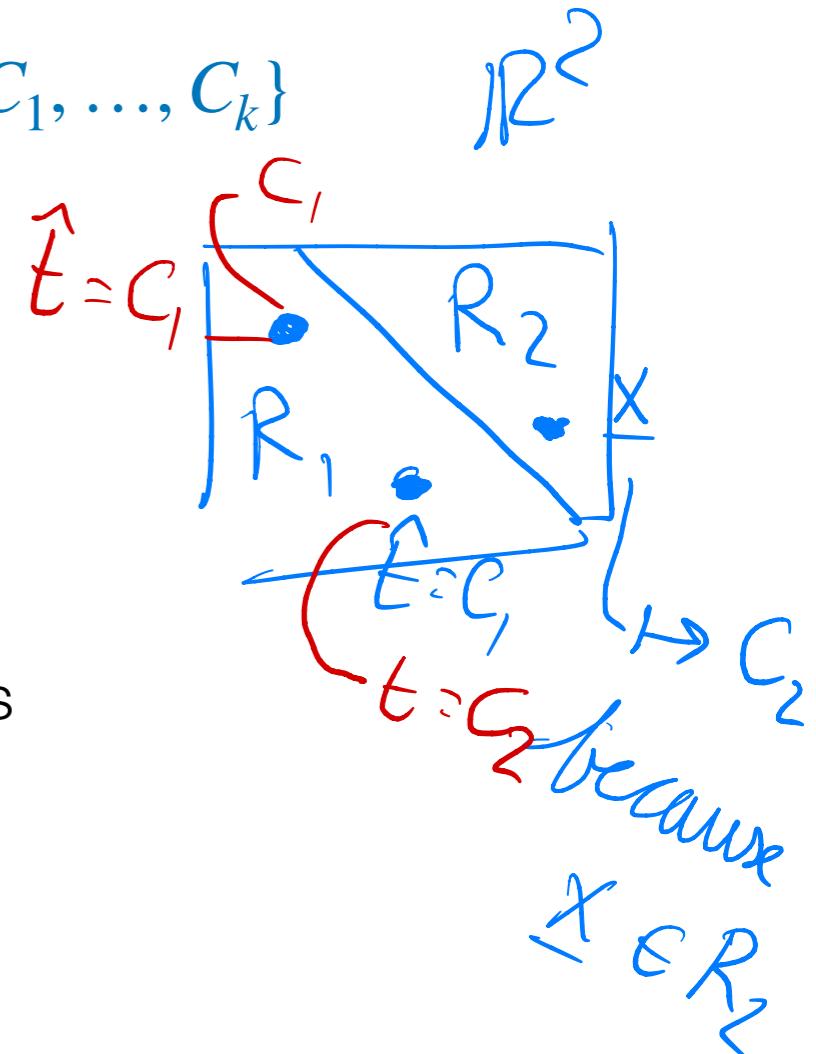
*matrix of counts or frequencies
for each type of prediction*

	R_1	R_2	\dots	R_K
C_1	6	1	\dots	0
C_2	5	3	\dots	1
\vdots	\vdots	\vdots	\ddots	\vdots
C_K	2	0	\dots	8

Pred C_1 pred C_2 ... pred C_K

ground truth

- Diagonal elements: correctly classified
- Off-diagonal elements: misclassified



Decision theory: Misclassification Rate

- Classification goal: Minimize the misclassification rate
- Assume observations are drawn from joint distribution $(\mathbf{x}, t) \sim p(\mathbf{x}, t)$
- Probability of a misclassification:

Note: For every one of K classes, we can make $K - 1$ mistakes (so $K(K - 1)$ mistakes)

(Only when predicted class C_i coincides with true class C_k)

$$p(\text{mistake}) = \sum_{i=1}^K \sum_{k \neq i} p(\mathbf{x} \in R_i, t = C_k) \xrightarrow{\substack{\text{assign to class } i \\ \text{should belong to class } k}} = 1 - p(\text{correct}) = 1 - \sum_{k=1}^K p(\mathbf{x} \in R_k, C_k)$$

Minimizing misclassification rate

- Assign \mathbf{x} to class C_k if

$$p(\mathbf{x}, t = C_k) > p(\mathbf{x}, t = C_j)$$

- Note: $p(\mathbf{x}, C_k) = p(C_k | \mathbf{x})p(\mathbf{x})$

→ instead maximize:

$$p(C_k | \mathbf{x}) > p(C_j | \mathbf{x})$$

Decision theory: Misclassification Rate

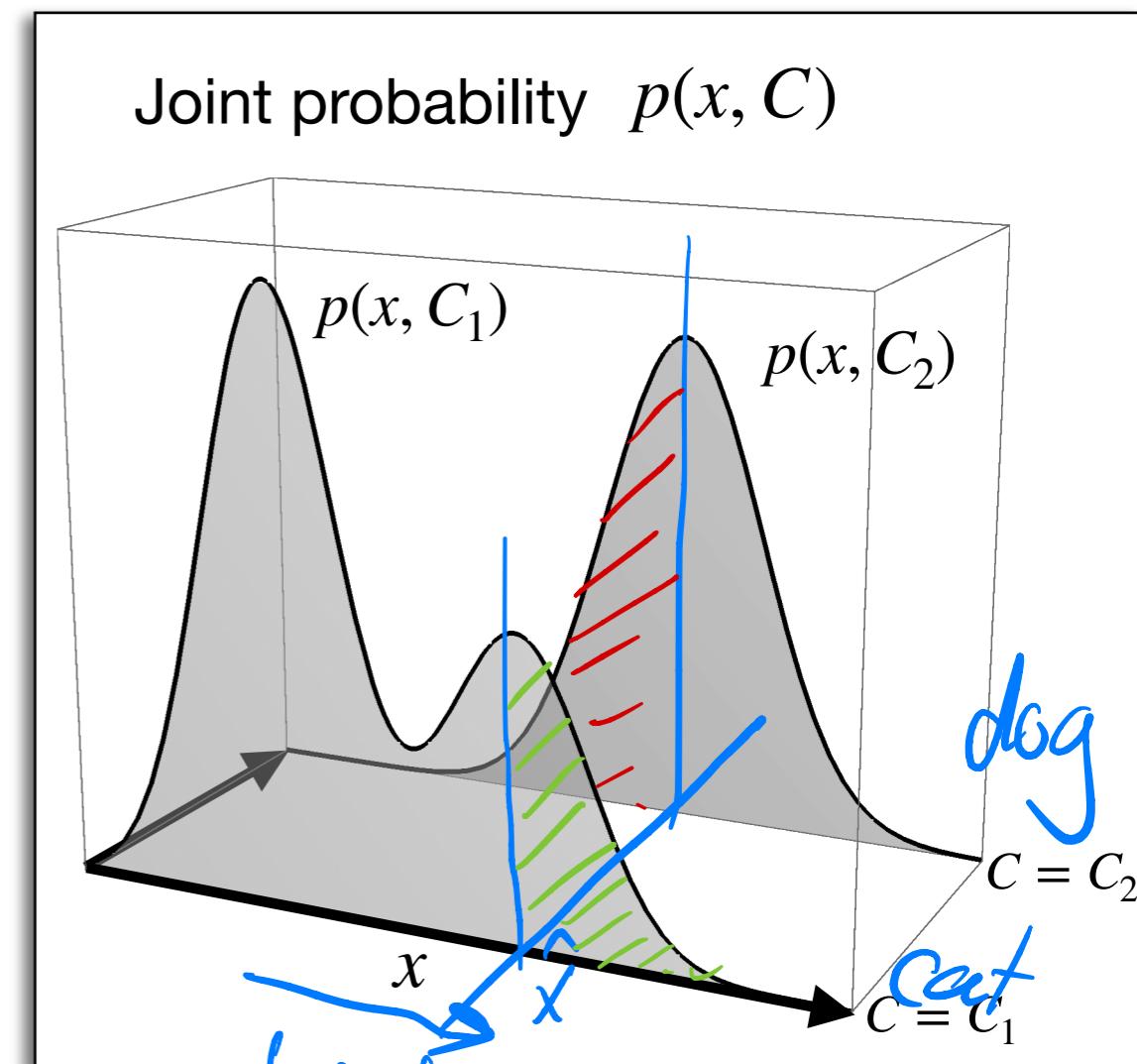
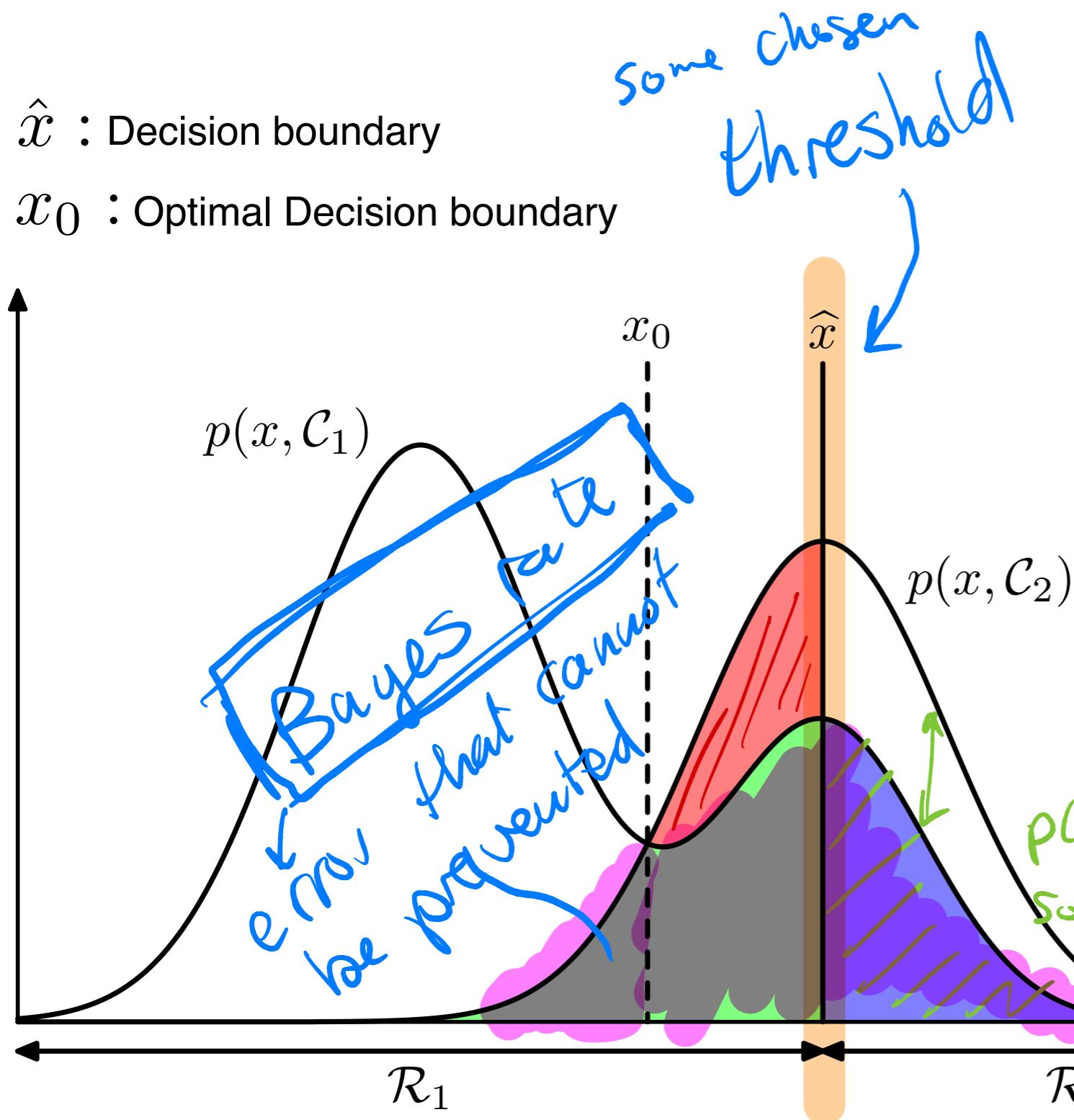


Figure: joint probability distributions and decision boundary (Bishop 1.24)

Minimizing the Misclassification Rate: Problems

- Not all errors have the same impact!

Example: Medical diagnosis of cancer

- Error 1: Label a healthy person as having cancer.

- Error 2: Label a sick person as healthy. Lack of treatment!

severity of error

- If cancer only occurs in 1% of all patients, a classifier which labels everyone as healthy has a misclassification rate of 1%!

class imbalance

should also
be corrected for

Expected Loss

$$p(x \in (a, b)) = \int_a^b p(x) dx$$

$$R, ; a \leq x \leq b$$



- Possible solution: use different weights for different error types

$$L = \begin{pmatrix} & \text{label cancer} & & \text{label healthy} \\ & 0 & & 1000 \\ & 1 & & 0 \end{pmatrix} \begin{matrix} \text{true cancer} \\ \text{true healthy} \end{matrix}$$

- Expected loss:

$$\mathbb{E}_{(\mathbf{x}, t) \sim p(\mathbf{x}, t)}[L] = \sum_{k,j} L_{kj} p(x \in R_j, t = C_k) = \sum_{k,j} L_{kj} \int_{R_j} p(x, C_k) dx$$

- Objective: Minimize the expected loss:

Assign \mathbf{x} to C_k if $\sum_{j=1}^K L_{jk} p(\mathbf{x}, C_j)$ is minimal

Classification Strategies

- Discriminant functions

Direct mapping of input to target:

$$\hat{y} = y(x, w)$$

- Probabilistic discriminative models

Posterior class probabilities:

$$p(C_k | x)$$

- Probabilistic generative models

Class-conditional densities:

$$p(x|C_k) \quad p(x, C_k)$$
$$p(C_k) \quad p(C_k|x)$$

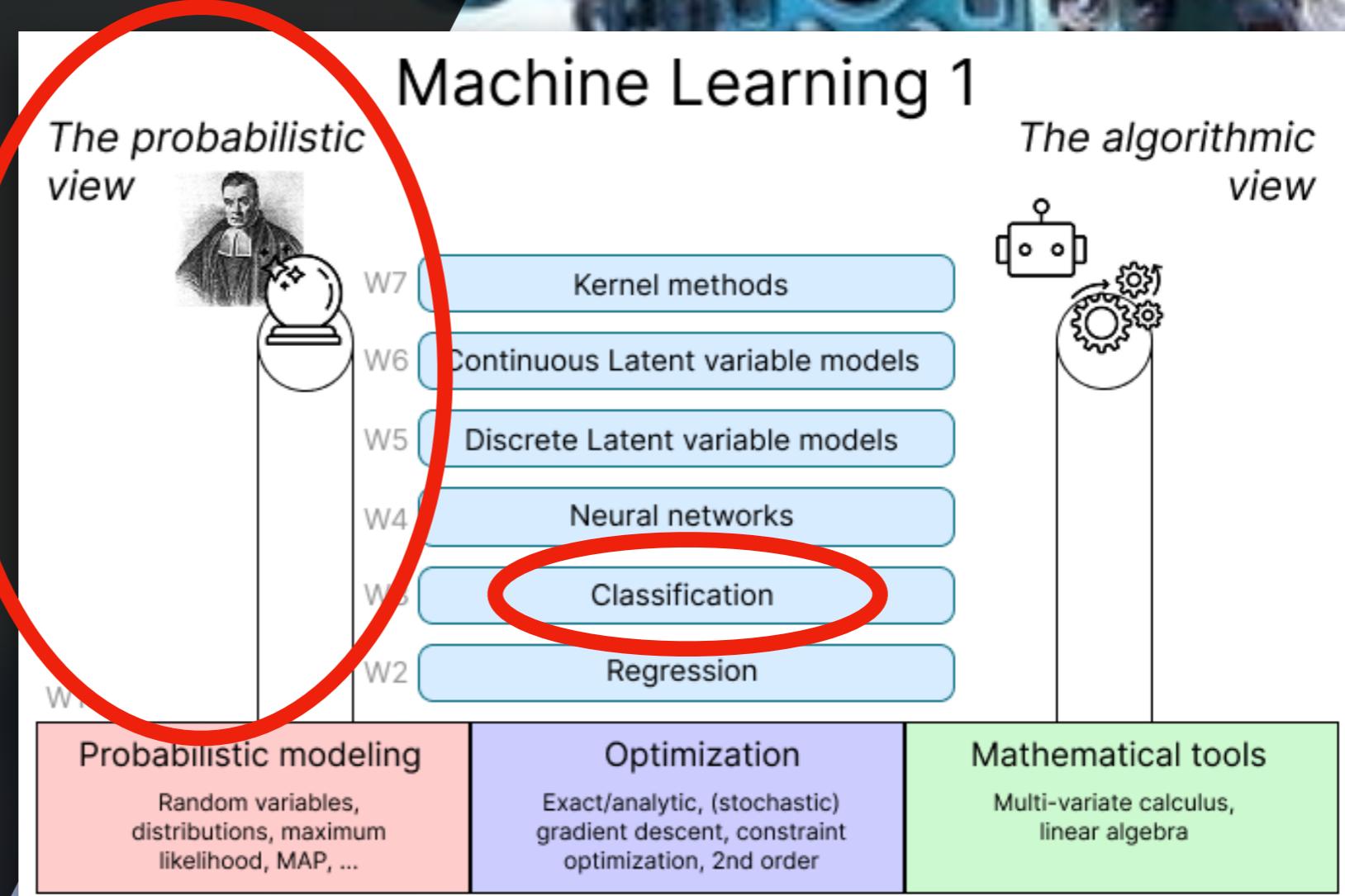
Prior class probabilities:

Machine Learning 1

Lecture 5.6 - Supervised Learning
Classification - **Probabilistic Generative Models**

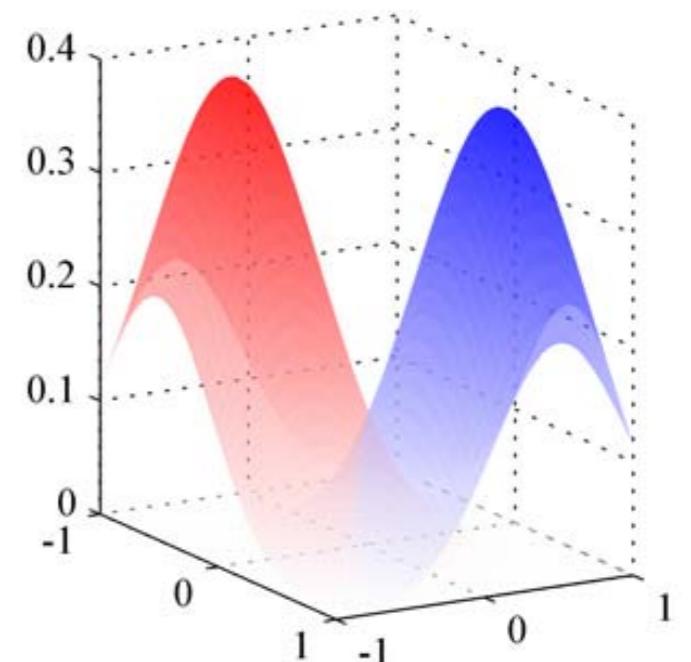
Erik Bekkers

(Bishop 1.5)



Probabilistic Generative Models: general K

- ▶ Model:
 - ▶ Class-conditional densities: $p(\mathbf{x} | C_k)$
 - ▶ Prior class probabilities: $p(C_k)$
- ▶ This gives access to joint and posterior:
 - ▶ $p(\mathbf{x}, C_k)$ and $p(C_k | \mathbf{x})$
- ▶ Why are we interested in joint or posterior?
 - ▶ Decision theory tells us that the best prediction for input \mathbf{x} , is to choose the class with highest joint $p(\mathbf{x}, C_k)$
 - ▶ Or equivalently: choose class with highest posterior $p(C_k | \mathbf{x})$
 - ▶ Decision boundary between C_k and C_j are at $p(C_k | \mathbf{x}) = p(C_j | \mathbf{x})$



Probabilistic Generative Models: $K = 2$

- ▶ Class-conditional densities:
- ▶ Prior class probabilities:
- ▶ Joint distribution:
- ▶ Posterior distribution: $K = 2$

$$p(C_1 | \mathbf{x}) = \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_1)p(C_1) + p(\mathbf{x} | C_2)p(C_2)} =$$

- ▶ Log odds $a = \ln \frac{\sigma}{1 - \sigma} =$

Logistic Sigmoid Function

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

$$\sigma(-a) = 1 - \sigma(a)$$

$$\sigma'(a) = \sigma(a)(1 - \sigma(a))$$

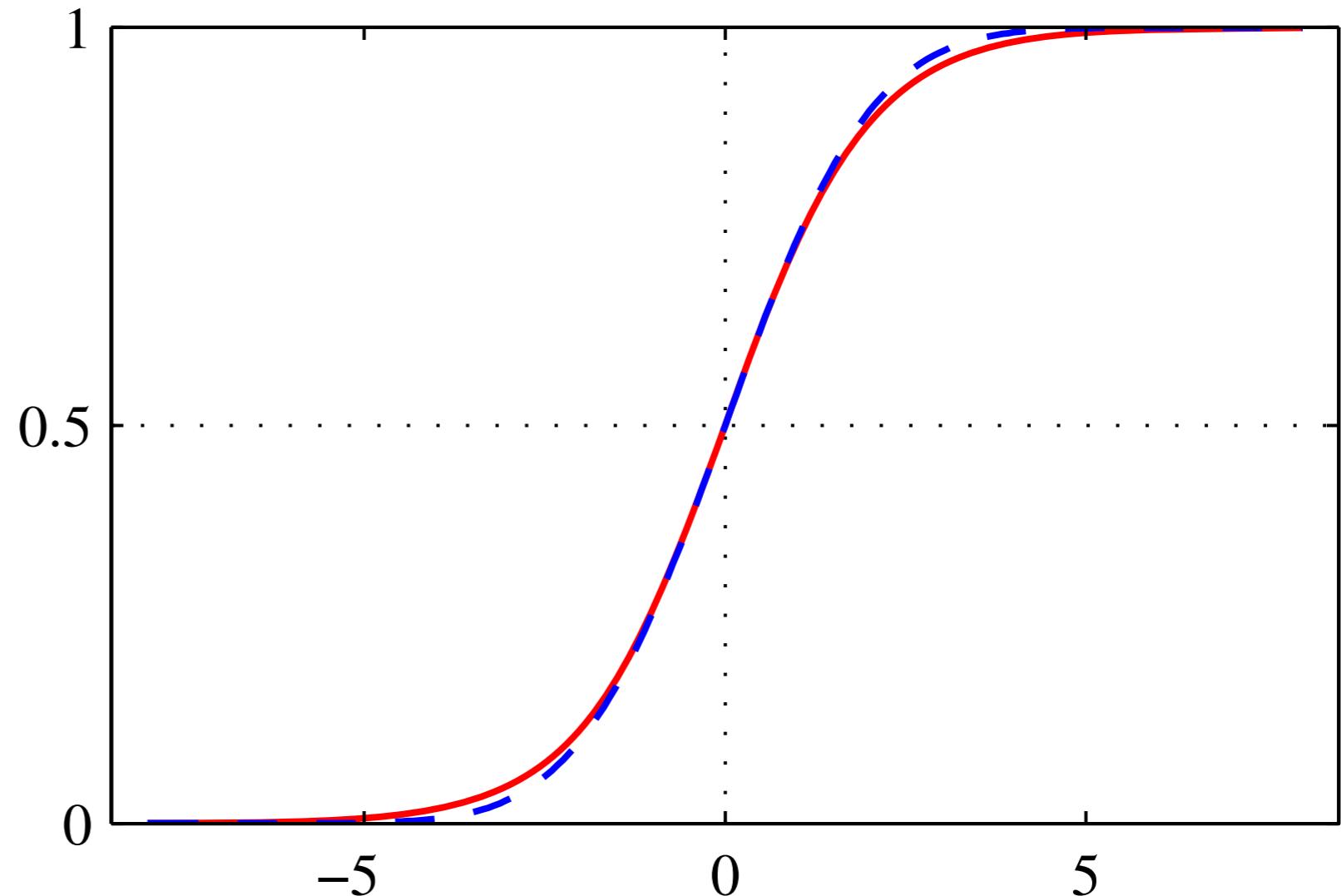


Figure: Logistic Sigmoid function (red) (Bishop 4.9)

Probabilistic Generative Models: general K

- For multiple classes (general K):

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k)p(C_k)}{\sum_{j=1}^K p(\mathbf{x} | C_j)p(C_j)} =$$

$$a_k = \ln(p(\mathbf{x} | C_k)p(C_k))$$

- Softmax: if $a_k \gg a_j$ for all $j \neq k$:

- Note: for $K = 2$:

$$p(C_1 | \mathbf{x}) = \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_1)p(C_1) + p(\mathbf{x} | C_2)p(C_2)} = \frac{1}{1 + \frac{p(\mathbf{x} | C_2)p(C_2)}{p(\mathbf{x} | C_1)p(C_1)}}$$

=

$$a = \frac{p(\mathbf{x} | C_2)p(C_2)}{p(\mathbf{x} | C_1)p(C_1)}$$

Class Conditional Densities: Continuous Inputs

- ▶ Gaussian Class-conditional densities:

$$p(\mathbf{x} | C_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)\right\}$$

- ▶ Class posteriors simplify with the assumption of shared covariance matrix: $\Sigma_k = \Sigma$

- ▶ For $K = 2$ classes they are given as $p(C_1 | \mathbf{x}) = \frac{1}{1 + \exp(-a)} = \sigma(a)$ with

$$a = \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)} =$$

Class Conditional Densities: Continuous Inputs

- ▶ Gaussian Class-conditional densities:

$$p(\mathbf{x} | C_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)\right\}$$

- ▶ Class posteriors simplify with the assumption of shared covariance matrix: $\Sigma_k = \Sigma$

- ▶ For $K = 2$ classes they are given as $p(C_1 | \mathbf{x}) = \frac{1}{1 + \exp(-a)} = \sigma(a)$ with

$$a = \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)} = \ln \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_1, \Sigma) - \ln \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_2, \Sigma) + \ln \frac{p(C_1)}{p(C_2)}$$

Class Conditional Densities: Continuous Inputs

- ▶ Gaussian Class-conditional densities:

$$p(\mathbf{x} | C_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)\right\}$$

- ▶ Class posteriors simplify with the assumption of shared covariance matrix: $\Sigma_k = \Sigma$

- ▶ For $K = 2$ classes they are given as $p(C_1 | \mathbf{x}) = \frac{1}{1 + \exp(-a)} = \sigma(a)$ with

$$\begin{aligned} a &= \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)} = \ln \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_1, \Sigma) - \ln \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_2, \Sigma) + \ln \frac{p(C_1)}{p(C_2)} \\ &= -\frac{1}{2} \ln |\Sigma| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2} \ln |\Sigma| + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_2)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) + \ln \frac{p(C_1)}{p(C_2)} \end{aligned}$$

Class Conditional Densities: Continuous Inputs

- ▶ Gaussian Class-conditional densities:

$$p(\mathbf{x} | C_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)\right\}$$

- ▶ Class posteriors simplify with the assumption of shared covariance matrix: $\Sigma_k = \Sigma$

- ▶ For $K = 2$ classes they are given as $p(C_1 | \mathbf{x}) = \frac{1}{1 + \exp(-a)} = \sigma(a)$ with

$$\begin{aligned} a &= \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)} = \ln \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_1, \Sigma) - \ln \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_2, \Sigma) + \ln \frac{p(C_1)}{p(C_2)} \\ &= -\frac{1}{2} \ln |\Sigma| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2} \ln |\Sigma| + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_2)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) + \ln \frac{p(C_1)}{p(C_2)} \\ &= (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \Sigma^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(C_1)}{p(C_2)} \end{aligned}$$

Class Conditional Densities: Continuous Inputs

- ▶ Gaussian Class-conditional densities:

$$p(\mathbf{x} | C_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)\right\}$$

- ▶ Class posteriors simplify with the assumption of shared covariance matrix: $\Sigma_k = \Sigma$

- ▶ For $K = 2$ classes they are given as $p(C_1 | \mathbf{x}) = \frac{1}{1 + \exp(-a)} = \sigma(a)$ with

$$\begin{aligned} a &= \ln \frac{p(\mathbf{x} | C_1)p(C_1)}{p(\mathbf{x} | C_2)p(C_2)} = \ln \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_1, \Sigma) - \ln \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_2, \Sigma) + \ln \frac{p(C_1)}{p(C_2)} \\ &= -\frac{1}{2} \ln |\Sigma| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2} \ln |\Sigma| + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_2)^T \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) + \ln \frac{p(C_1)}{p(C_2)} \\ &= (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \Sigma^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(C_1)}{p(C_2)} \end{aligned}$$

- ▶ Generalized Linear Model: $p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$

$$\mathbf{w} = \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

$$w_0 = -\frac{1}{2} \boldsymbol{\mu}_1^T \Sigma^{-1} \boldsymbol{\mu}_1 + \frac{1}{2} \boldsymbol{\mu}_2^T \Sigma^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(C_1)}{p(C_2)}$$

Example: Linear Discriminant Analysis for K=2

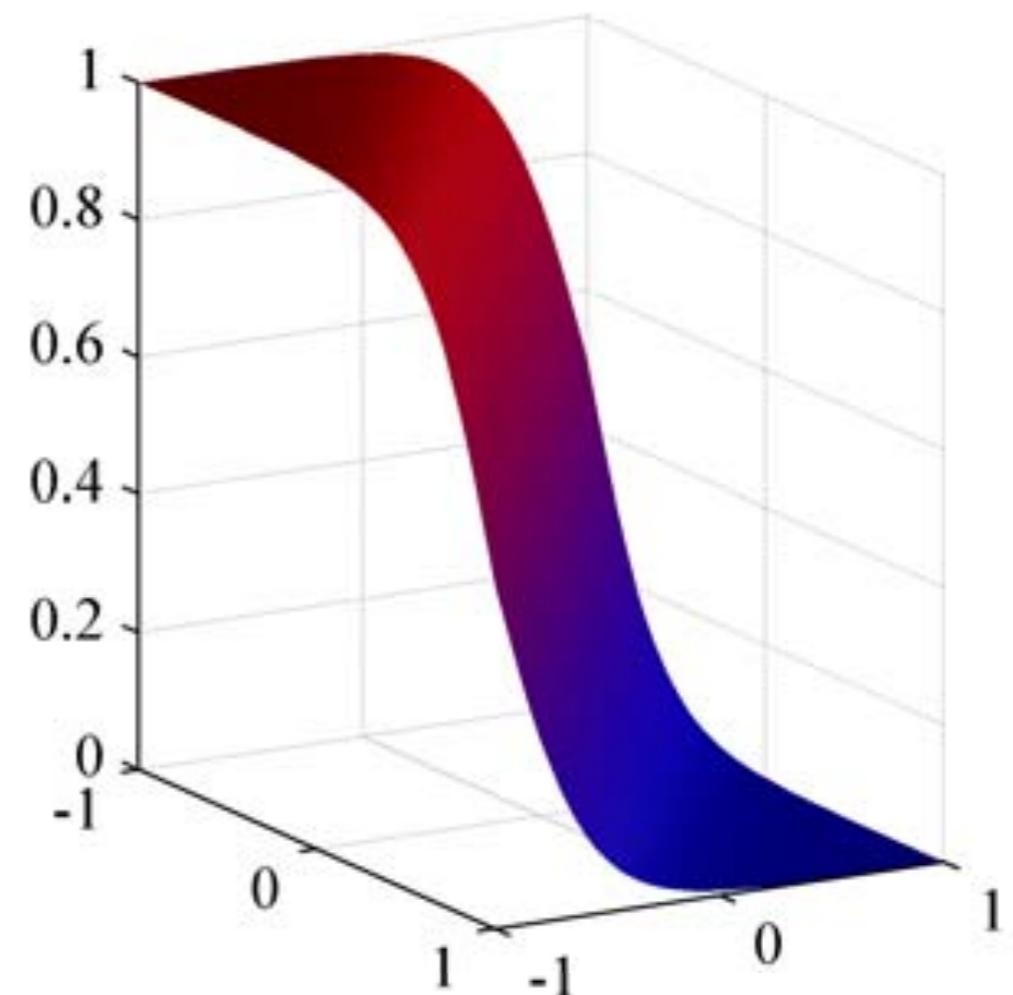
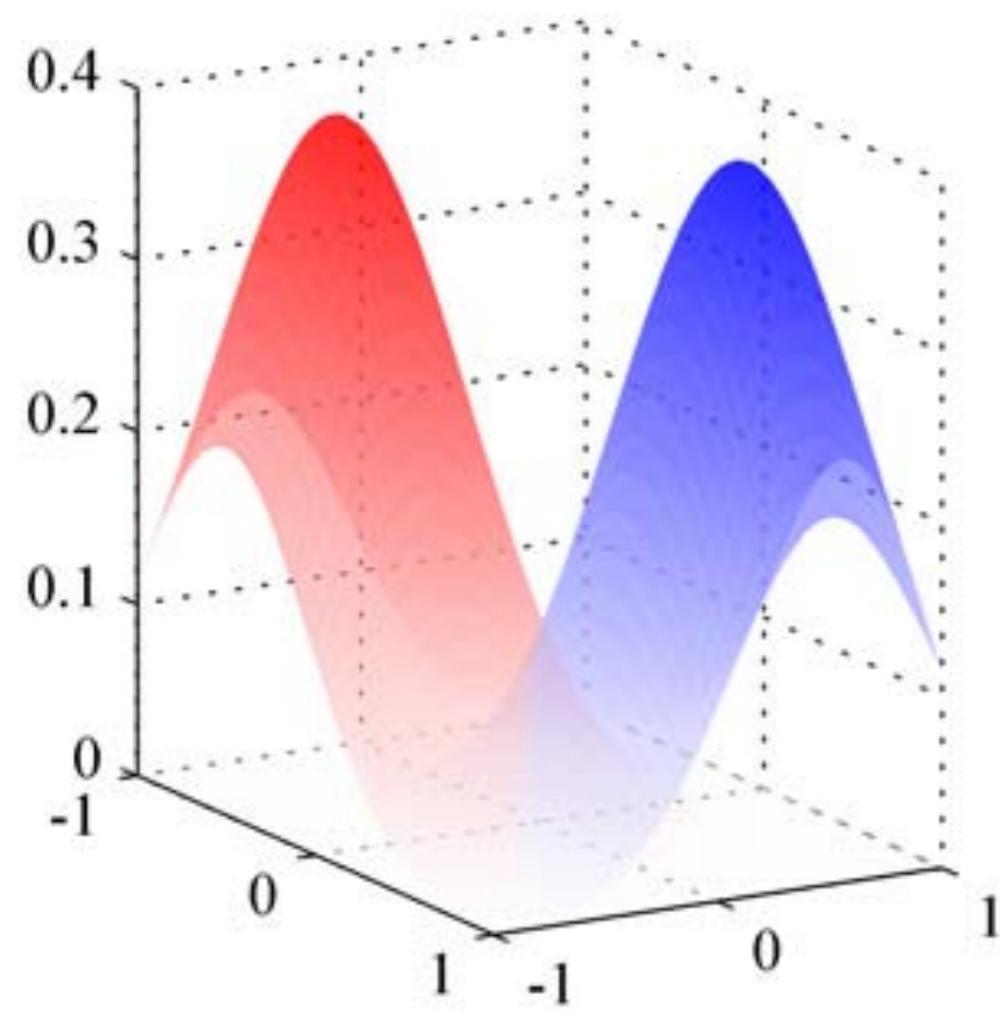


Figure: Left: class conditional densities $p(\mathbf{x} | C_k)$. Right: posterior $p(C_1 | \mathbf{x})$ as sigmoid of linear function of \mathbf{x} . (Bishop 4.9)

Linear Discriminant Analysis: General K

- ▶ Gaussian Class-conditional densities & fixed covariance:

$$p(\mathbf{x} | C_k) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \Sigma_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)\right\}$$

- ▶ Posterior distributions:

$$p(C_k | \mathbf{x}) = \frac{\exp(a_k(\mathbf{x}))}{\sum_{j=1}^K \exp(a_j(\mathbf{x}))}$$

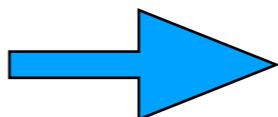
- ▶ $a_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$

$$\mathbf{w}_k = \Sigma_k^{-1} \boldsymbol{\mu}_k$$

$$w_{k0} = -\frac{1}{2} \boldsymbol{\mu}_k^T \Sigma_k^{-1} \boldsymbol{\mu}_k + \ln p(C_k)$$

- ▶ Decision boundary:

$$p(C_k | \mathbf{x}) = p(C_j | \mathbf{x})$$



- ▶ If all covariance matrices are different $\Sigma_k \neq \Sigma_j$ then $a_k(\mathbf{x})$ also contains quadratic terms in \mathbf{x}

Example: LDA and QDA

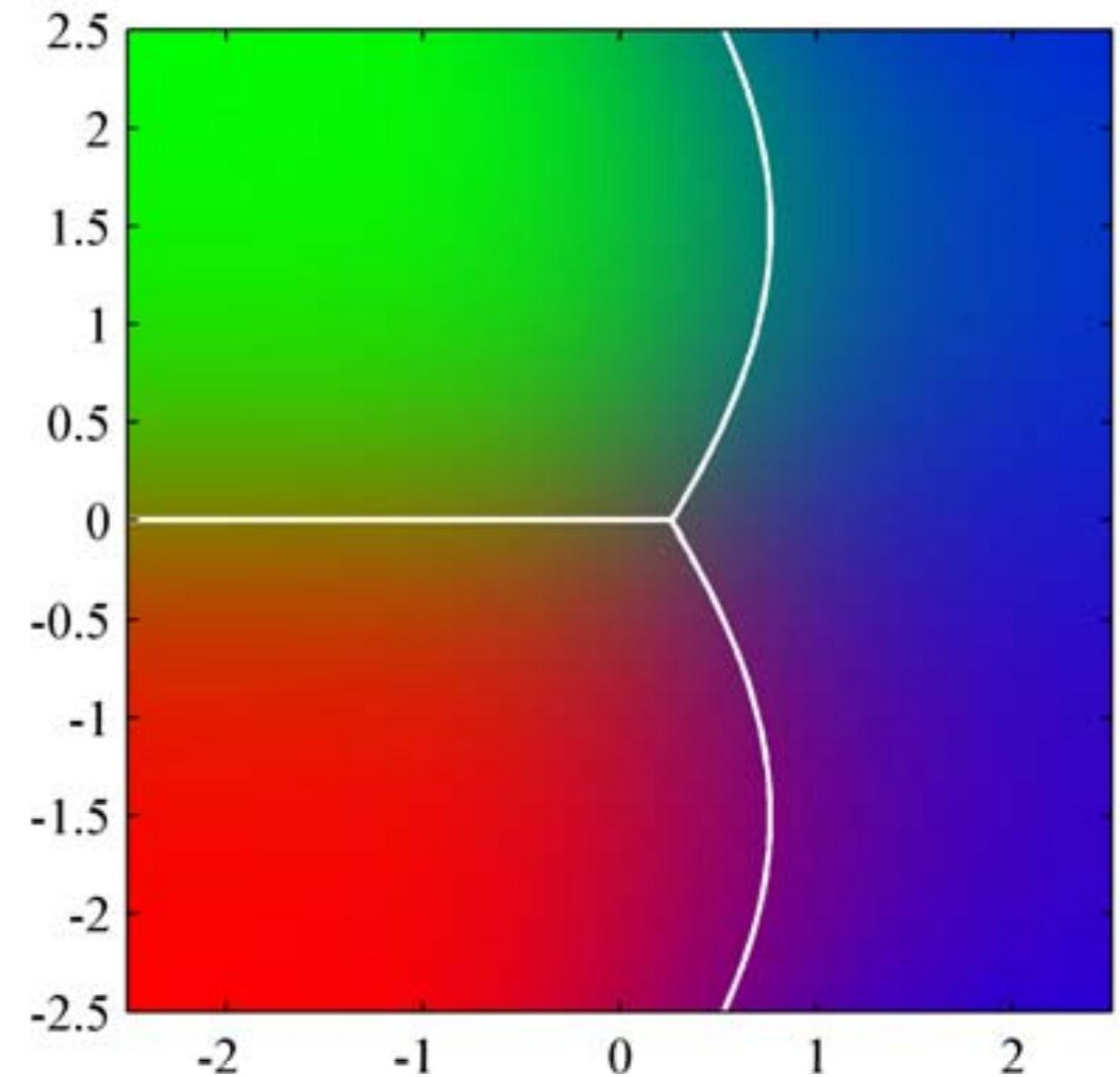
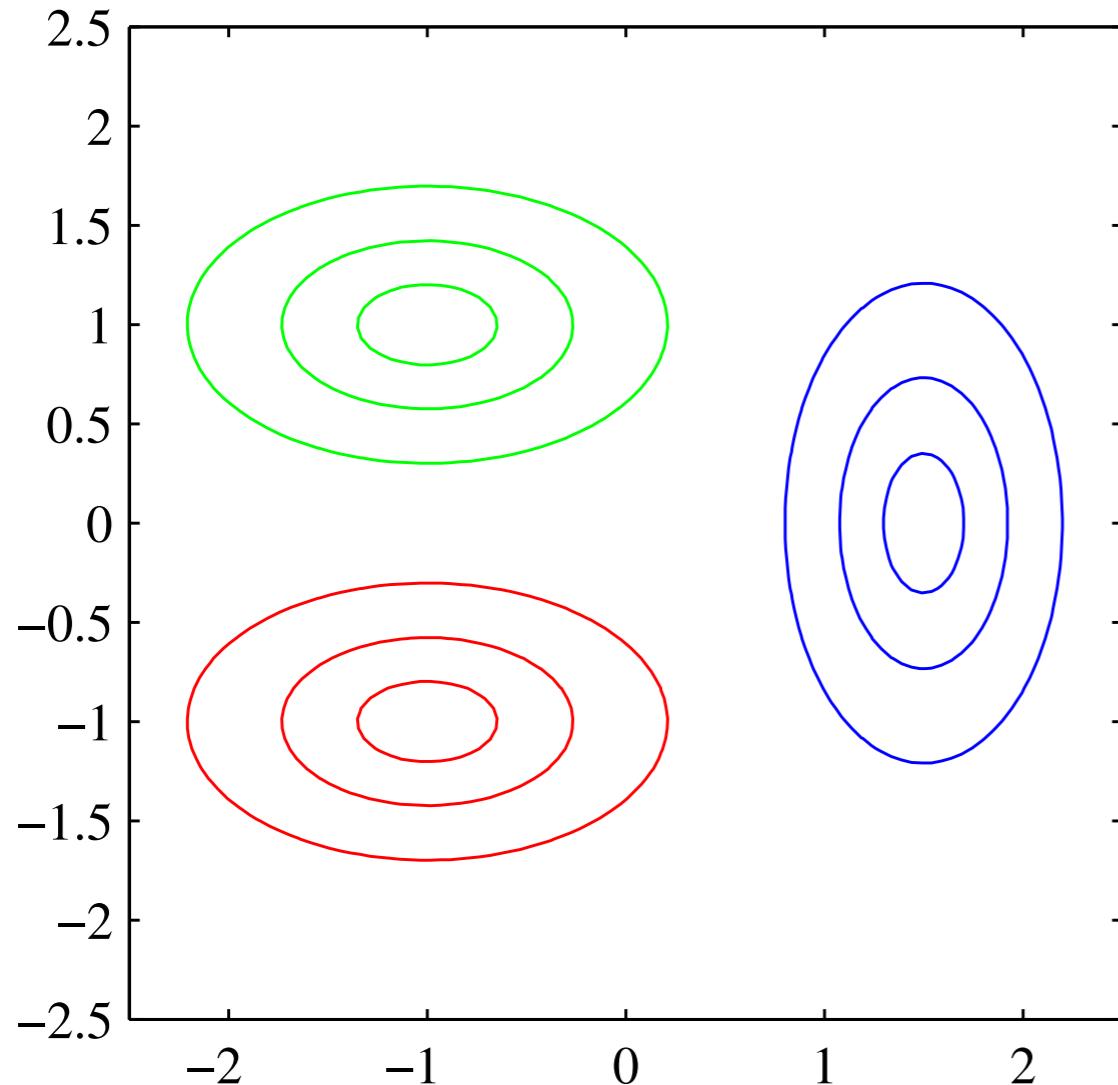


Figure: Left: Gaussian class conditional densities $p(\mathbf{x} | C_k)$, red and green have same covariance matrix. Right: posterior $p(C_k | \mathbf{x})$ distributions (RGB vectors) and decision boundaries. (Bishop 4.9)