

Machine Learning 1

Lecture 7 - Classification With Basis Functions
- Probabilistic Discriminative Models - Logistic
Regression

Erik Bekkers

$$y_1(x; \omega) > y_2(x; \omega) \rightarrow \text{class 1}$$

Catch-up: binary /
Classification with linear reg
The perceptron

Gaussians (Σ_1, Σ_2)

$$\begin{matrix} p(x|C) \\ p(C) \end{matrix} \left\{ \begin{matrix} p(C|x) \\ \end{matrix} \right\} \xrightarrow{\text{LDA}}$$



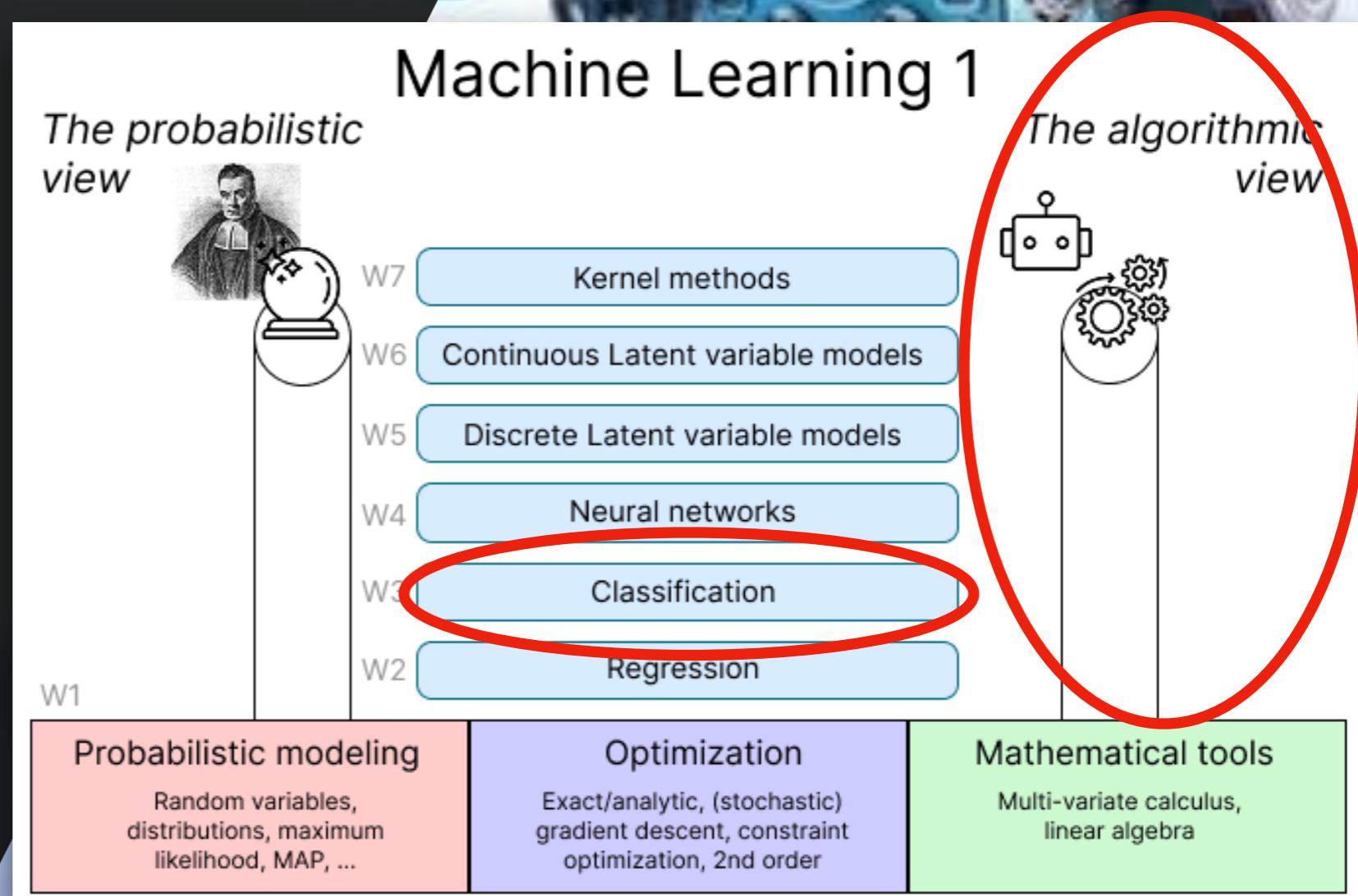
Machine Learning 1

Lecture 6.4 - Supervised Learning

Classification? - Discriminative Models - Least Squares Regression

Erik Bekkers

(Bishop 4.1.3)



Least Squares for Classification

- Each class C_k has its own linear model:

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

- Shorter notation: $y(\mathbf{x}) = \tilde{\mathbf{W}}^T \tilde{\mathbf{x}}$

- Matrix $\tilde{\mathbf{W}}$, column k contains:

$M \times K$
Input vector:

- Output/prediction vector:

$$\begin{pmatrix} y_1(\mathbf{x}) \\ y_2(\mathbf{x}) \\ \vdots \\ y_K(\mathbf{x}) \end{pmatrix} = \begin{pmatrix} -\tilde{w}_0 \\ -\tilde{w}_1 \\ \vdots \\ -\tilde{w}_K \end{pmatrix} \tilde{\mathbf{x}} = \tilde{\mathbf{w}}_k = \text{Concat}(w_{k0}, \mathbf{w}_k) \in \mathbb{R}^M$$
$$\tilde{\mathbf{x}} = \text{Concat}(1, \mathbf{x}) \in \mathbb{R}^M$$

$$\mathbf{y}(\mathbf{x}) = \begin{pmatrix} y_1(\mathbf{x}) \\ y_2(\mathbf{x}) \\ \vdots \\ y_K(\mathbf{x}) \end{pmatrix} = \tilde{\mathbf{W}}^T \tilde{\mathbf{x}} \in \mathbb{R}^K$$

- Classification:** Assign \mathbf{x} to class C_k if $k = \operatorname{argmax}_j y_j(\mathbf{x})$

Discriminative linear model

Least Squares for Classification (II)

- Data set: $N \times (D + 1)$ data matrix

$$\tilde{X} = \begin{pmatrix} \tilde{\mathbf{x}}_1^T \\ \vdots \\ \tilde{\mathbf{x}}_N^T \end{pmatrix}$$

- $N \times K$ target matrix

$$T = \begin{pmatrix} \mathbf{t}_1^T \\ \vdots \\ \mathbf{t}_N^T \end{pmatrix} =$$

one-hot encoding
for class C_3

| | | | | |
|--------------------|---|---|---|---|
| 0 | 0 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 | 0 |
| ⋮ | | | | |
| $y_k(\tilde{x}_n)$ | | | | |

- Use sum-of-squares regression error function

$$E_D(\tilde{\mathbf{W}}) = \frac{1}{2} \text{Tr} [(\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})^T(\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})] = \frac{1}{2} \sum_n \sum_k \left(\sum_m \tilde{X}_{nm} \tilde{W}_{mk} - t_{nk} \right)^2$$

- Minimize $E_D(\tilde{\mathbf{W}})$ as a function of $\tilde{\mathbf{W}}$:

$$\frac{\partial E}{\partial \mathbf{W}} = 0$$

- Solution: $\tilde{\mathbf{W}}_{LS} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{T} = \tilde{\mathbf{X}}^\dagger \mathbf{T}$

- Discriminant function: $\mathbf{y}_{LS}(\mathbf{x}) = \tilde{\mathbf{W}}_{LS}^T \tilde{\mathbf{x}} \in \mathbb{R}^K$

Matrix contains all products $\text{error}_i \cdot \text{error}_i$ but we are only interested in the trace (sum of diagonal) which is $\sum_k \text{error}_k^2$

for a matrix M with elements m_{ij} , i.e., $M = \begin{pmatrix} m_{11} & m_{12} & m_{13} & \dots \\ m_{21} & m_{22} & m_{23} & \dots \\ m_{31} & m_{32} & m_{33} & \dots \\ \vdots & \vdots & \vdots & \ddots \end{pmatrix}$

the trace $\text{Tr}[M]$ sums over the diagonal.

On the previous slide the matrix is $K \times K$ and contains the errors

$$m_{ij} = \sum_n (y_{ni} - t_{ni})(t_{nj} - t_{nj})$$

whereas we're only interested in the diagonal

$$m_{kk} = \sum_n (y_{nk} - t_{nk})^2$$

Least Squares for Classification: Problems

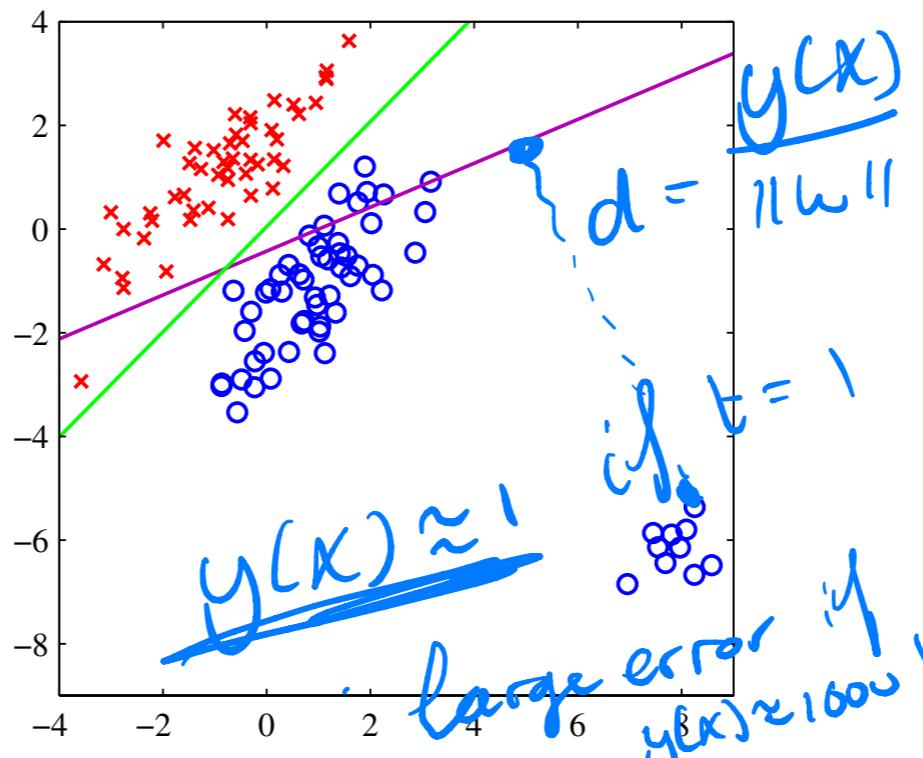
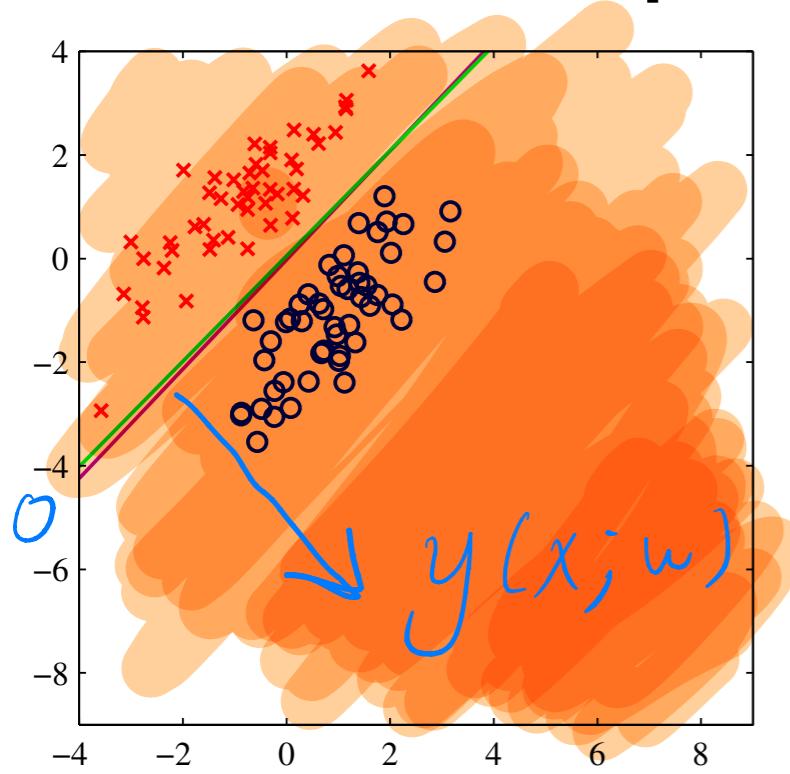


Figure: least squares is very sensitive to outliers (Bishop 4.4)

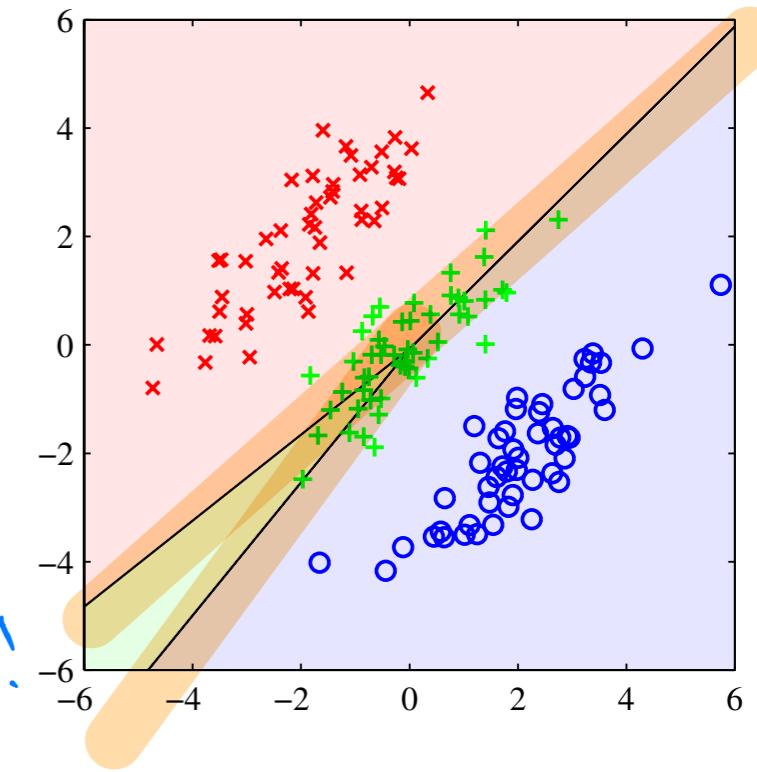


Figure: masking for least squares for $K > 2$ (Bishop 4.5)

1. The decision boundaries are very sensitive to outliers
↳ easy points but the least squares loss wants for every point $y(x) \approx 1$ when $t_n = 1$
2. For $K > 2$ some decision regions can become very small or are even completely ignored
3. The components of $\mathbf{y}_{LS}(\mathbf{x})$ are not probabilities!

- $y_k(\mathbf{x})$
- if $\sum_{k=1}^K t_k = 1$

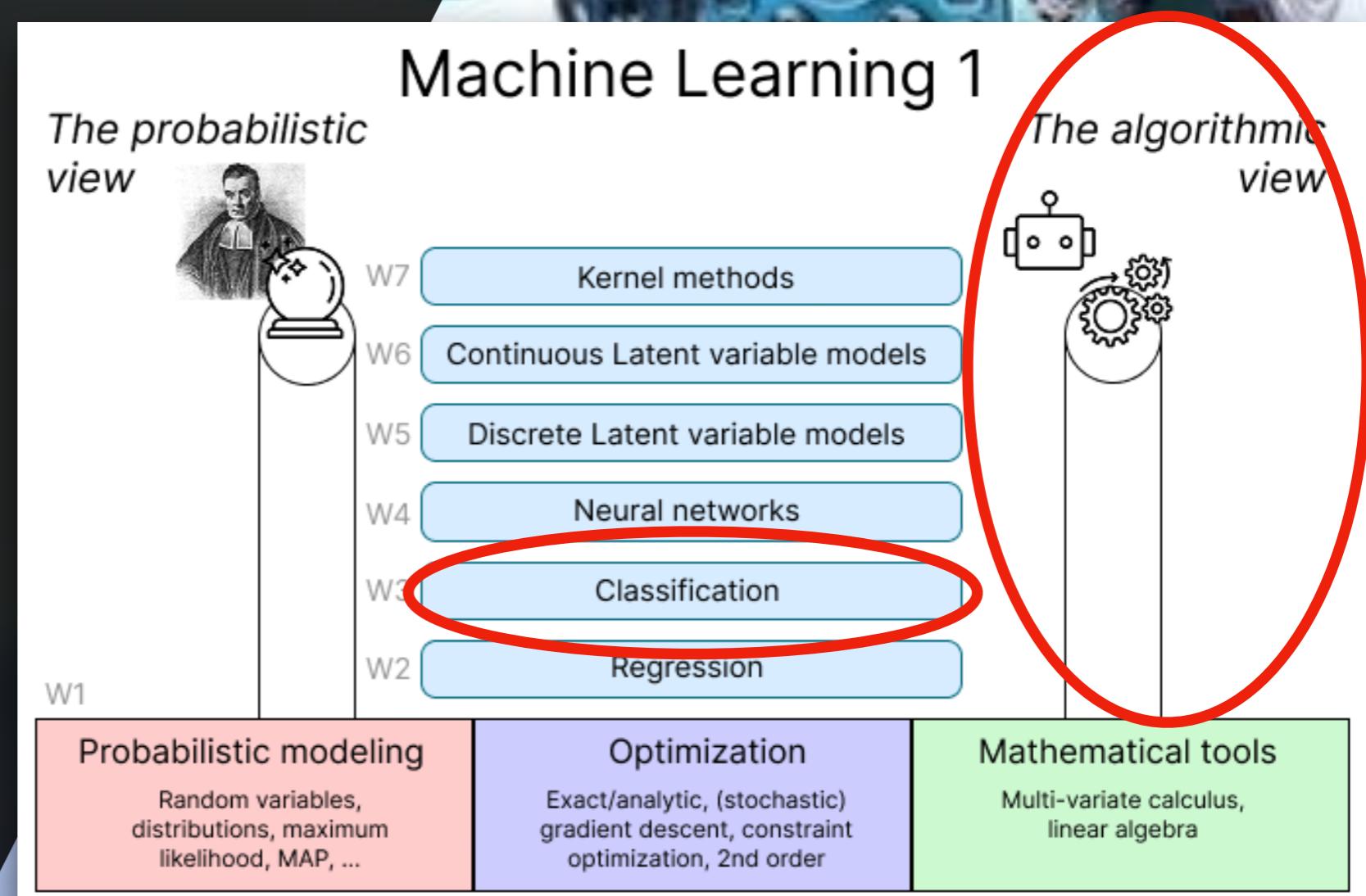
then $\rightarrow \sum_k y_k = 1$

Machine Learning 1

Lecture 6.5 - Supervised Learning
Classification - Discriminative Models -
The Perceptron

Erik Bekkers

(Bishop 4.1.7)



K 2

The Perceptron Algorithm

- Input: $\mathbf{x} \in \mathbb{R}^D$
- Targets: $t \in \{C_1, C_2\} = \{1, -1\}$

- Prediction: $y(\mathbf{x}) = f(\underbrace{\mathbf{w}^T \phi(\mathbf{x})}_a)$, with $f(a) = \begin{cases} 1 & , a \geq 0 \\ -1 & , a < 0 \end{cases}$

- Class decisions: assign \mathbf{x} to class C_1 if $\underbrace{\mathbf{w}^T \phi_n}_a \geq 0$
(and to C_2 if $\underbrace{\mathbf{w}^T \phi}_a < 0$)

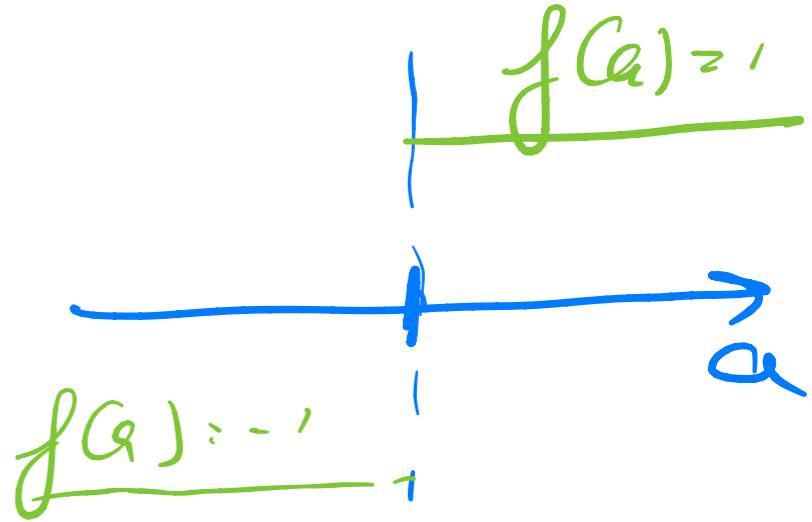
- For correct classification: find \mathbf{w} such that for all (\mathbf{x}_n, t_n) :

$$\underbrace{\mathbf{w}^T \phi_n}_a \cdot \underbrace{t_n}_a \geq 0$$

- Perceptron criterion: $E_P(\mathbf{w}) = - \sum_{\substack{n \in \mathcal{M}}} \mathbf{w}^T \phi(\mathbf{x}_n) t_n$

both positive: +
both negative: +
otherwise: -

with $\mathcal{M} = \{n : \mathbf{w}^T \phi_n t_n < 0\}$ ← set of incorrectly classified points



Perceptron: Stochastic Gradient Descent

- $$E_P(\mathbf{w}) = - \sum_{n \in \mathcal{M}} \mathbf{w}^T \phi(\mathbf{x}_n) t_n$$

$$= - \sum_{n \in \mathcal{M}} E_n(\mathbf{w})$$

- Stochastic Gradient Descent (SGD). For each misclassified \mathbf{x}_n :

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n(\mathbf{w})^T$$

$$= \mathbf{w}^{(\tau)} + \eta \phi_n t_n$$

(after each update we recompute the set of incorrect predictions and randomly select one point)

- If \mathbf{X} is linearly separable, then perceptron SGD will converge

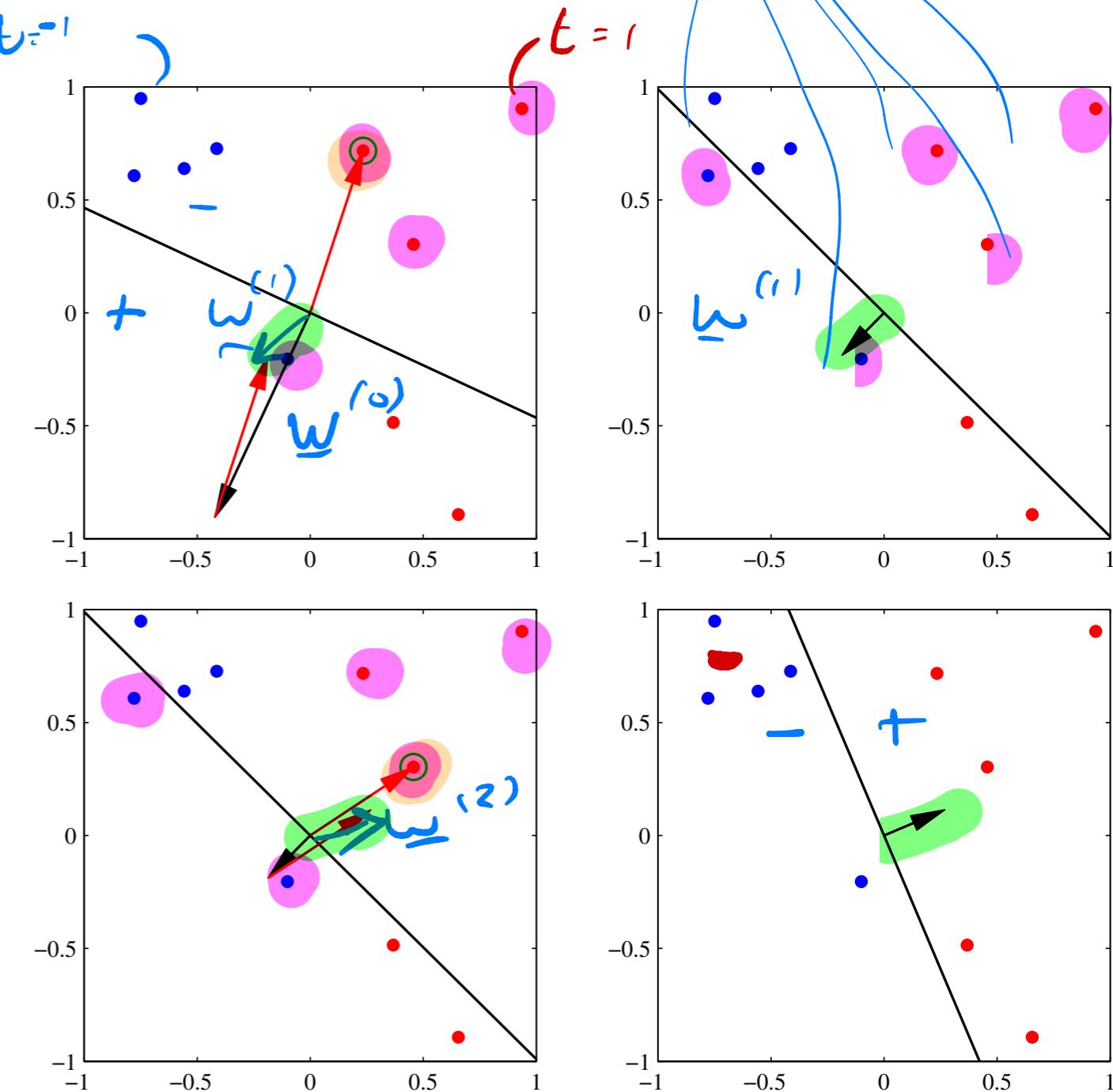


Figure: SGD for perceptron criterion (Bishop 4.7)
for \mathbf{x}_n in C_1 : add $\phi(\mathbf{x}_n)$ to \mathbf{w}
for \mathbf{x}_n in C_2 : subtract $\phi(\mathbf{x}_n)$ from \mathbf{w} .

Problems: Perceptron

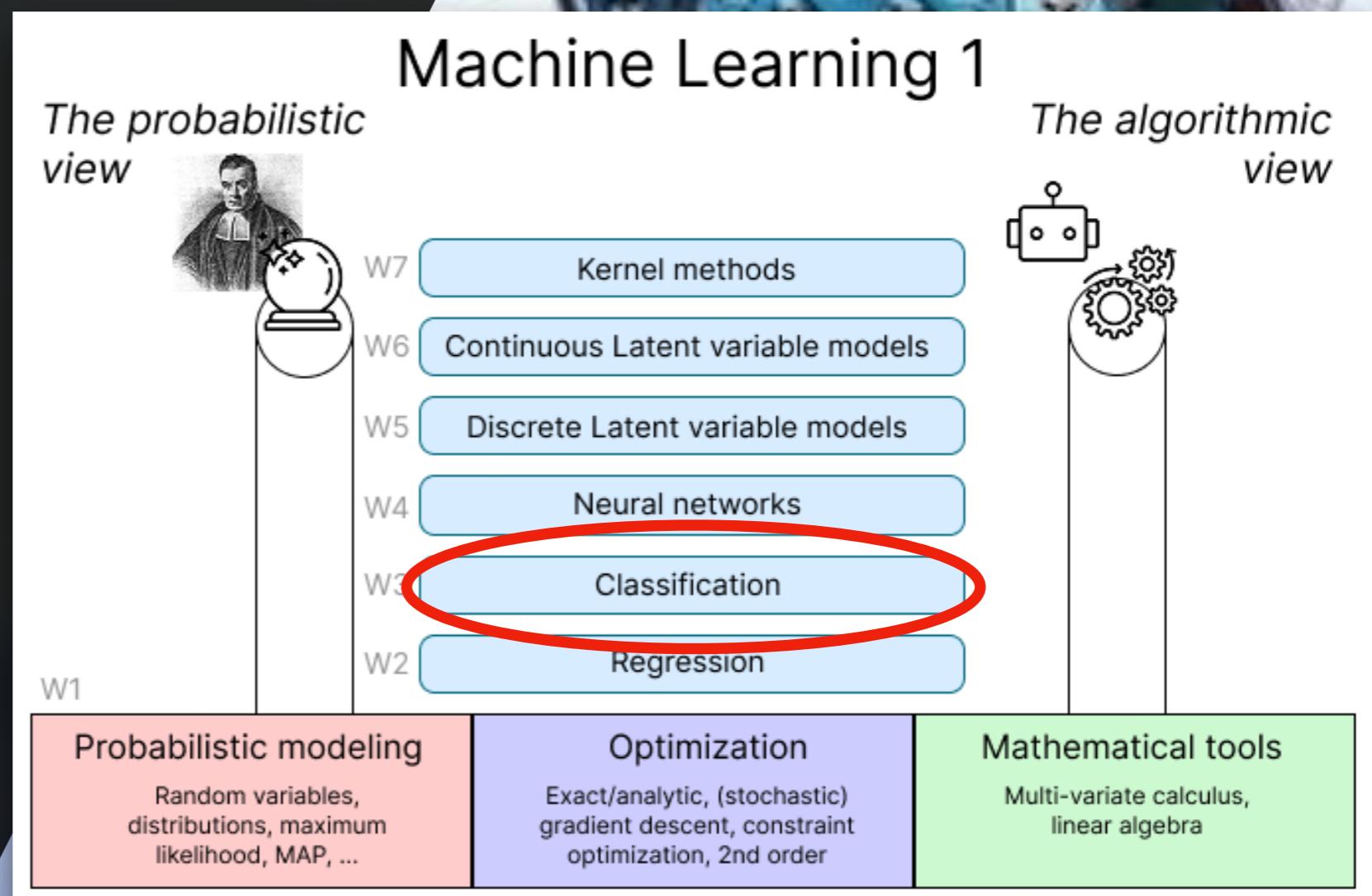
- ▶ Perceptron only works for 2 classes
 - ▶ There might be many solutions depending on the initialization of \mathbf{w} and on the order in which data is presented in SGD
 - ▶ If dataset is not linearly separable, the perceptron algorithm will not converge.
 - ▶ Based on linear combination of fixed basis functions.
 - ↳ learning basis
via function
via multi-layer perceptrons
- Random init*
- essentially determined by user*
- learning basis*
- via function*
- via multi-layer perceptrons*

Machine Learning 1

Lecture 7.1 - Supervised Learning
Classification - Classification With Basis Functions

Erik Bekkers

(Bishop 4.3.1)



Example: Use of Basis Functions

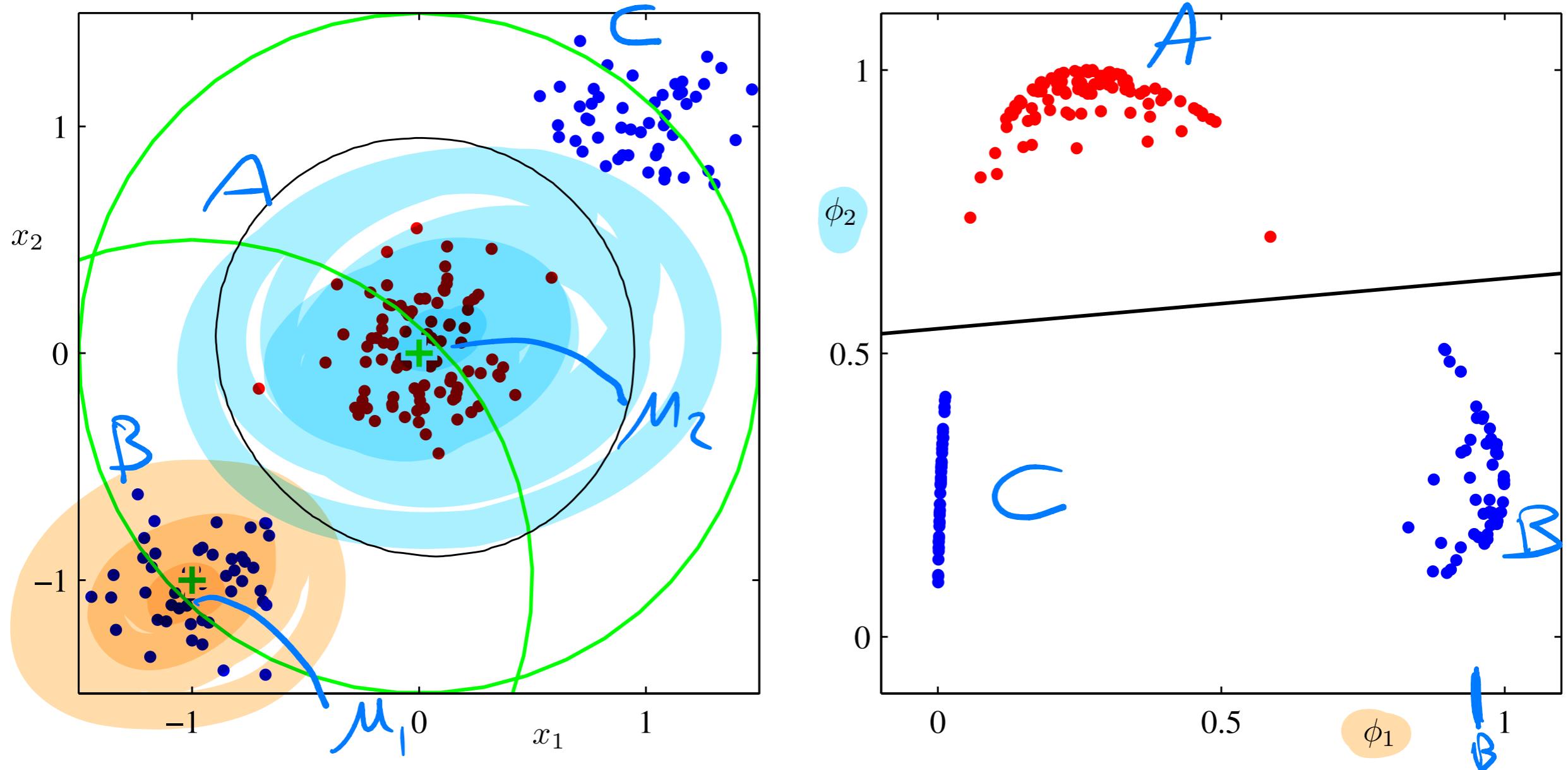


Figure: Left: original input space (x_1, x_2) , right: space of two gaussian basis functions with centres shown by the green crosses. (Bishop 4.12)

$$\phi_1(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^T(\mathbf{x} - \boldsymbol{\mu}_1)\right)$$

$$\phi_2(\mathbf{x}) = \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_2)^T(\mathbf{x} - \boldsymbol{\mu}_2)\right)$$

I

Limitations of Fixed Basis Functions

Advantages:

- ✓ ▶ Closed form solution for least-squares problem
- ✓ ▶ Tractable Bayesian treatment
- ✓ ▶ Nonlinear models mapping input variables to target variables through basis functions

neural networks

Limitations:

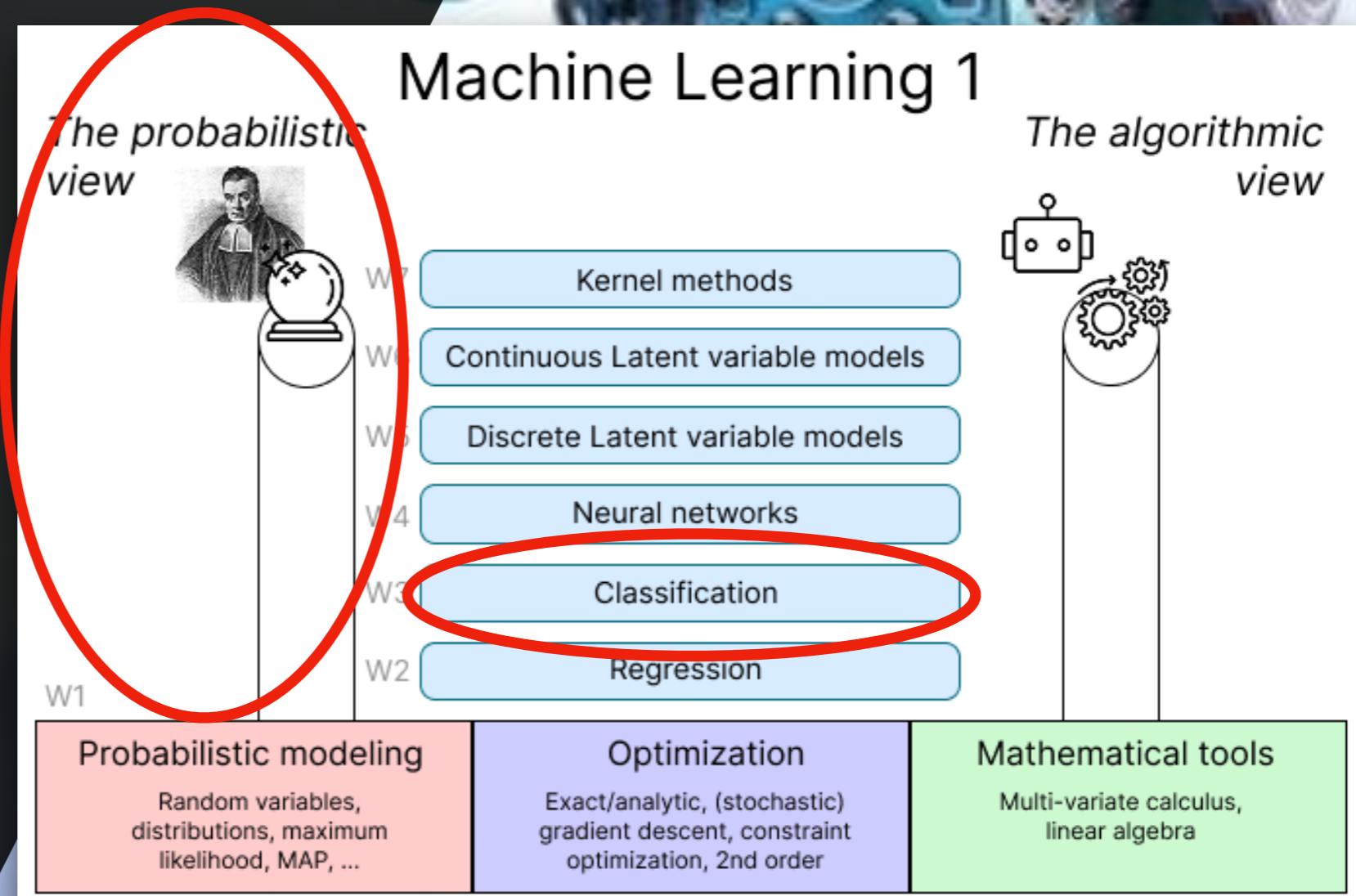
- ✗ ▶ Assumption: Basis functions $\phi_j(\mathbf{x})$ are fixed, not learned.
- ✗ ▶ **Curse of dimensionality:** to cover growing dimensions D of input vectors, the number of basis functions needs to grow rapidly / exponentially

Machine Learning 1

Lecture 7.2 - Supervised Learning
Classification - Probabilistic Discriminative
Models - Logistic Regression

Erik Bekkers

(Bishop 4.3.2)



Classification Strategies

- Discriminant functions

Direct mapping of input to target:

e.g. Perceptron
linear model

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

$\left\{ \begin{array}{l} t=0 \\ t>0 \\ t<0 \end{array} \right.$ $\left\{ \begin{array}{l} y=0 \\ y>0 \\ y<0 \end{array} \right.$

$w^T x$

- Probabilistic discriminative models

Posterior class probabilities:

$$p(C_k | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$$

- Probabilistic generative models

Class-conditional densities:

Prior class probabilities:

Gaussian:

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

LDA or QDA

gives posterior of the form

$$p(C_k) \quad p(C_k | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x})$$

in terms of $\{\mu_k\}, \Sigma$

Logistic Regression for Two Classes

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

From now on we write instead of t

$$\mathbf{t} = (t_1, \dots, t_N)^T \quad \text{with} \quad t_n \in \{C_1, C_2\} = \{1, 0\}$$

- Basis functions $\phi := \phi(\mathbf{x}) = (\phi_1(\mathbf{x}), \phi_2(\mathbf{x}), \dots)^T \in \mathbb{R}^M$

- Probabilistic Discriminative Linear Models: Posteriors $p(C_k | \phi)$ are modeled with nonlinear functions f that take a linear function of ϕ as input.

$$p(C_k | \phi, \mathbf{w}) = f(\mathbf{w}^T \phi)$$

$$C_1 \Leftrightarrow t=1, C_2 \Leftrightarrow t=0$$

- Logistic regression ($K=2$)

Combined into a single expression

$$\begin{cases} p(C_1 | \phi, \mathbf{w}) = y(\phi) = \sigma(\mathbf{w}^T \phi) \\ p(C_2 | \phi, \mathbf{w}) = 1 - y(\phi) = 1 - \sigma(\mathbf{w}^T \phi) \end{cases}$$
$$p(t | \phi, \mathbf{w}) = y(\phi)^t (1 - y(\phi))^{1-t}$$
$$p(C_1 | \phi, \mathbf{w}) \quad p(C_2 | \phi, \mathbf{w})$$

Logistic Regression for Two Classes

- Let $\phi \in \mathbb{R}^M$ and M -dimensional feature vector
- Probabilistic Discriminative Modeling* with Logistic Regression:

$$p(C_1 | \phi, w) = \sigma(w^T \phi) \quad \text{and} \quad p(C_2 | \phi, w) = 1 - \sigma(w^T \phi)$$

parameters: $w : M$ parameters \rightarrow Scales linearly

- Generative Modeling* with Gaussian conditional densities:

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

and class priors $p(C_k)$ + K params

parameters:

Scales quadratically with M $\left\{ \begin{array}{l} \mu_k : 2M \text{ parameters} \\ \Sigma : \sim 2M^2 \text{ parameters} \end{array} \right.$

param could be further reduced with Naive Bayes

Logistic Regression for Two Classes

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

$$\mathbf{t} = (t_1, \dots, t_N)^T \quad \text{with} \quad t_n \in \{C_1, C_2\} = \{0, 1\}$$

- Conditional likelihood function:

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

$p(\mathbf{x}, t | \mathbf{y})$

$$y_n = p(C_1 | \phi_n) = \sigma(\mathbf{w}^T \phi_n)$$

$$\phi_n := \phi(\mathbf{x}_n)$$

- Maximizing the conditional likelihood/minimizing the cross-entropy

$$E(\mathbf{w}) = -\ln p(\mathbf{t} | \mathbf{X}, \mathbf{w}) = -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln (1 - y_n)\}$$

information theory

- $E(\mathbf{w})$: convex, but no closed form solution!

$y_n = \sigma(\mathbf{w}^T \phi_n)$ is nonlinear in \mathbf{w}

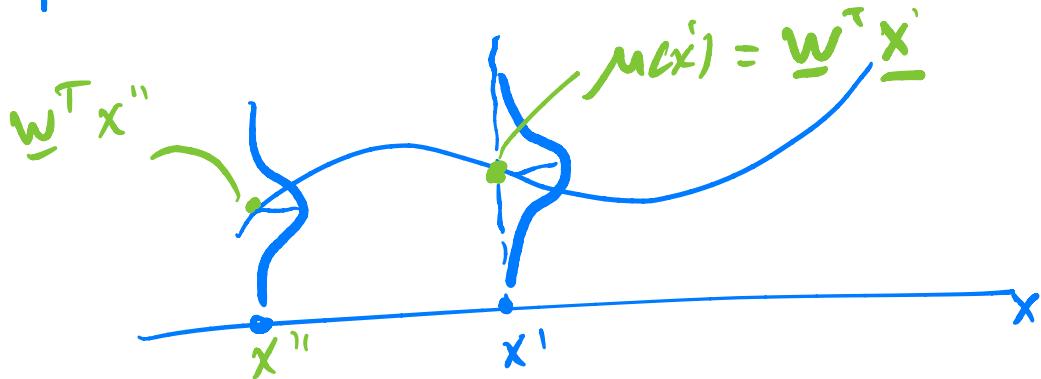
$$H_{\mathbf{w}} = \frac{\partial^2}{\partial \mathbf{w} \cdot \partial \mathbf{w}^T} E(\mathbf{w})$$

see lecture 7.4.

$$\forall \mathbf{x} : \mathbf{x}^T \mathbf{w} \geq 0$$

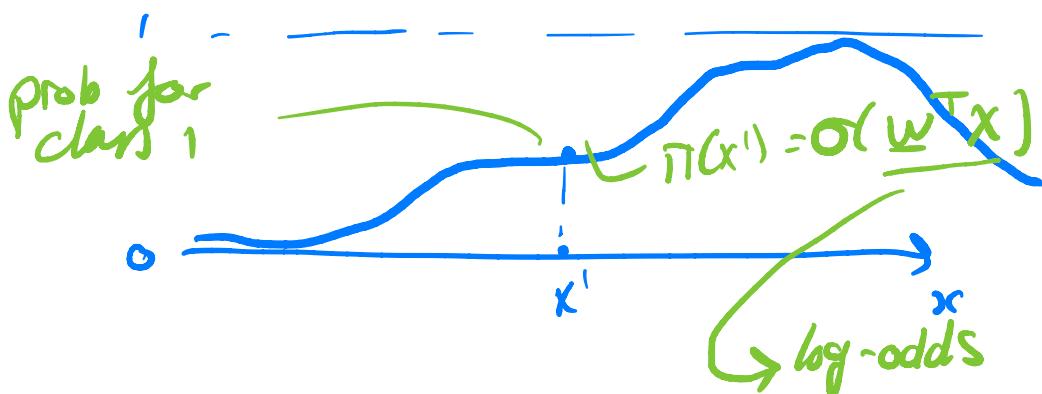
Regression w Gaussian predictor

$$p(t' | \underline{x}', \underline{w}) = N(t' | \mu(\underline{x}'), \sigma^2)$$



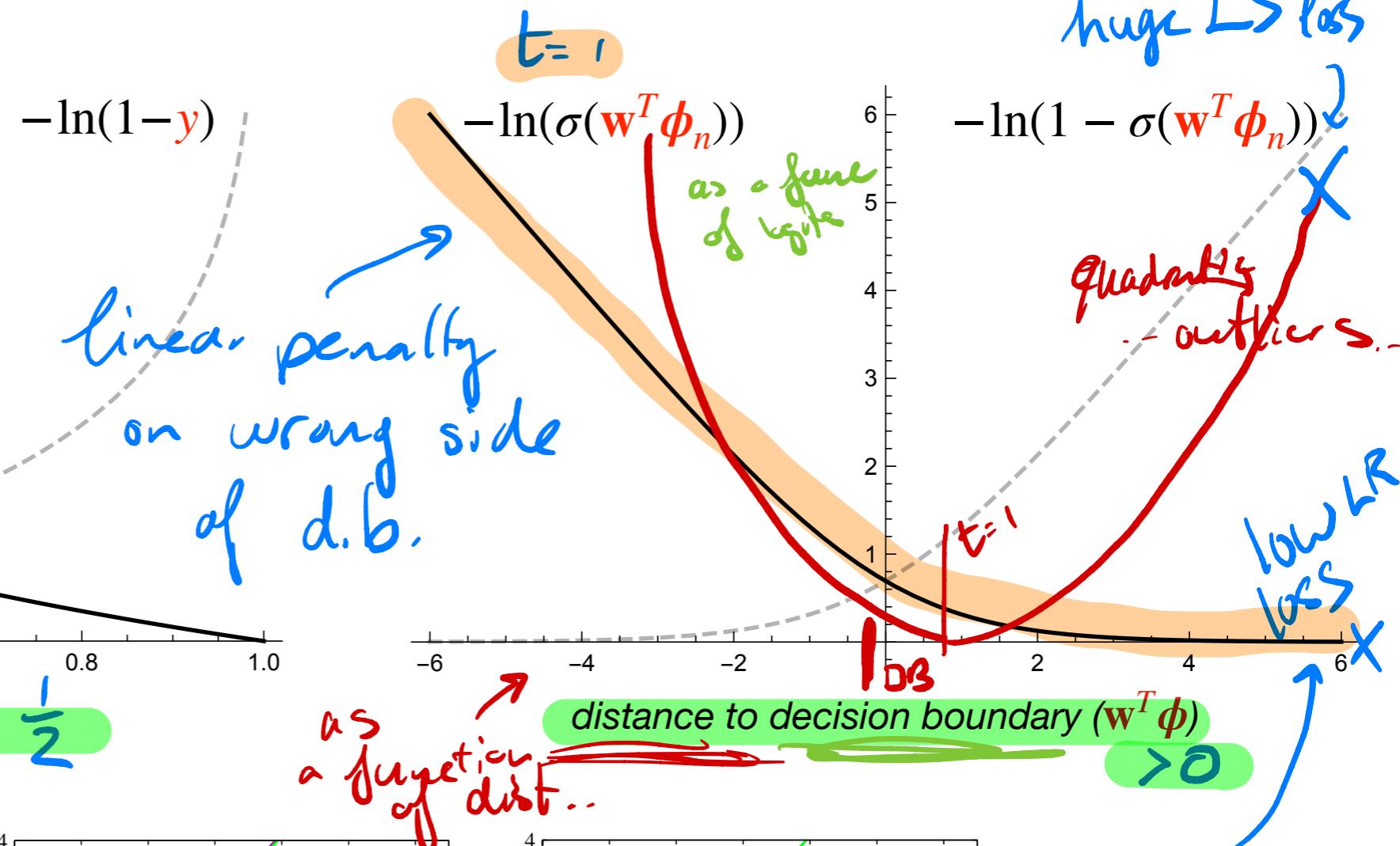
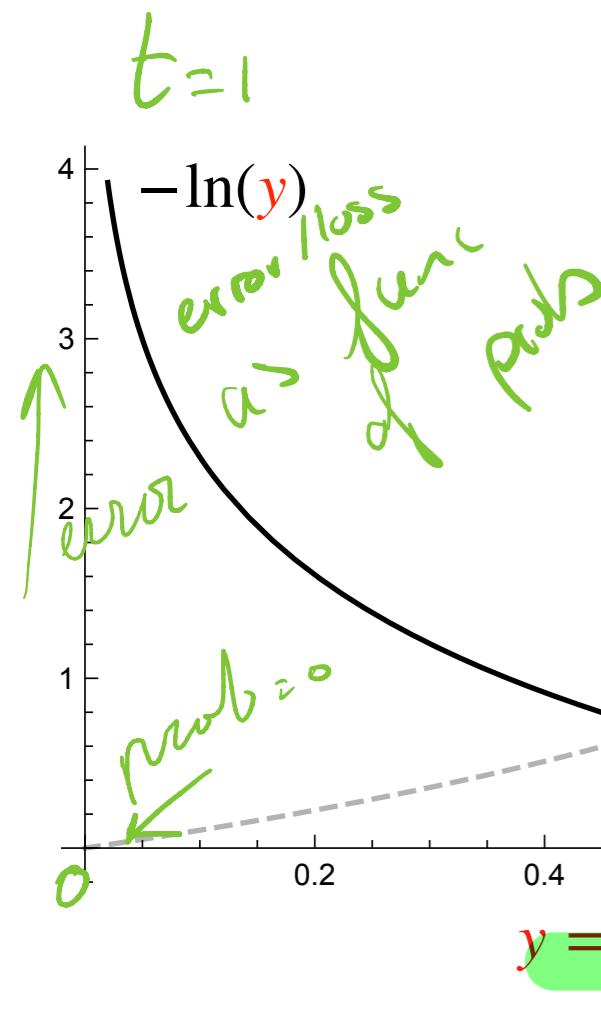
Classification w Bernoulli pred dist.

$$\begin{aligned} p(t' | \underline{x}', \underline{w}) &= \text{Bernoulli}(t' | \pi(x')) \\ &= \pi(x')^{t'} (1 - \pi(x'))^{1-t'} \end{aligned}$$

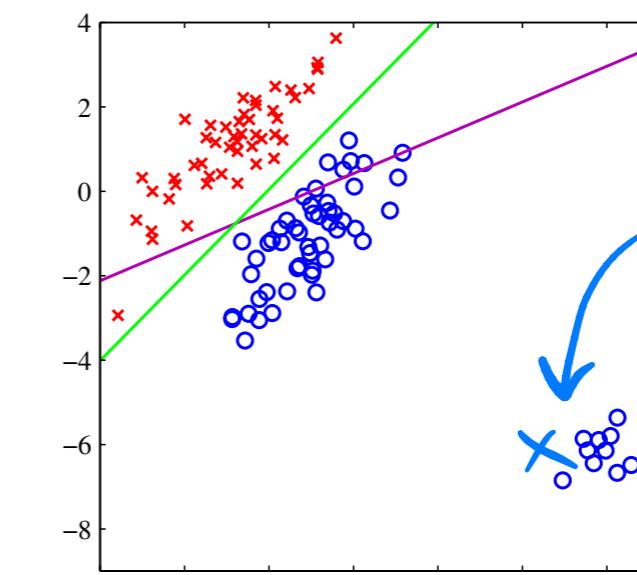
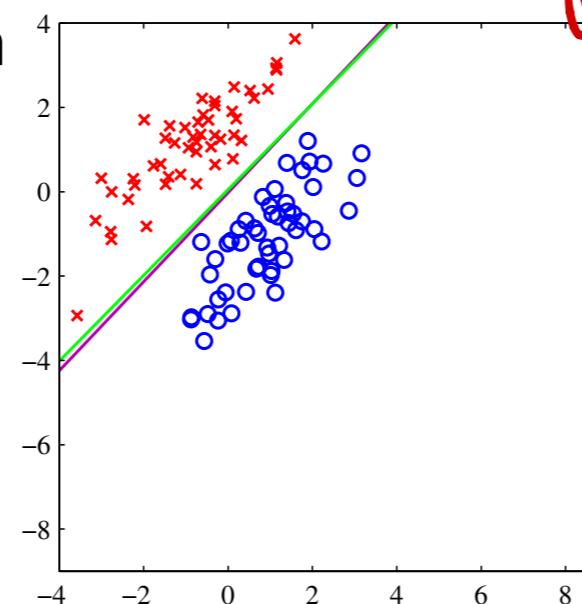


The cross-entropy loss

$$E_n(\mathbf{w}) = - \{ t_n \ln(y_n) + (1 - t_n) \ln(1 - y_n) \}$$



Least squares outlier problem does not take place with logistic regression



large dist to dec. b.
 $\phi_n^T \mathbf{w} \gg 0$
 $t=1$

Classification with Logistic Regression

- Given: Dataset $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N)^T$ with binary targets
 $\mathbf{t} = (t_1, \dots, t_N)^T$ with $t_n \in \{C_1, C_2\} = \{0,1\}$
- Basis functions: $\boldsymbol{\phi}(\mathbf{x}) = (\phi_0(\mathbf{x}), \dots, \phi_{M-1}(\mathbf{x}))^T$
- Posterior distributions (Prob. Discr. Model): $p(C_1 | \mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}))$

- Minimize

$$E(\mathbf{w}) = -\ln p(\mathbf{t} | \mathbf{X}, \mathbf{w}) = -\sum_{n=1}^N t_n \ln y_n + (1 - t_n) \ln(1 - y_n)$$

with *stochastic gradient descent* or *iterative reweighted least squares*, to find \mathbf{w}^*

Convex
so unique

- New datapoint \mathbf{x}' is assigned to class C_1 if

$$p(C_1 | \mathbf{x}', \mathbf{w}^*) = \sigma((\mathbf{w}^*)^T \boldsymbol{\phi}(\mathbf{x}')) > \frac{1}{2} , \quad (\mathbf{w}^*)^T \underline{\boldsymbol{\phi}} > 0$$

- Decision boundaries:

linear

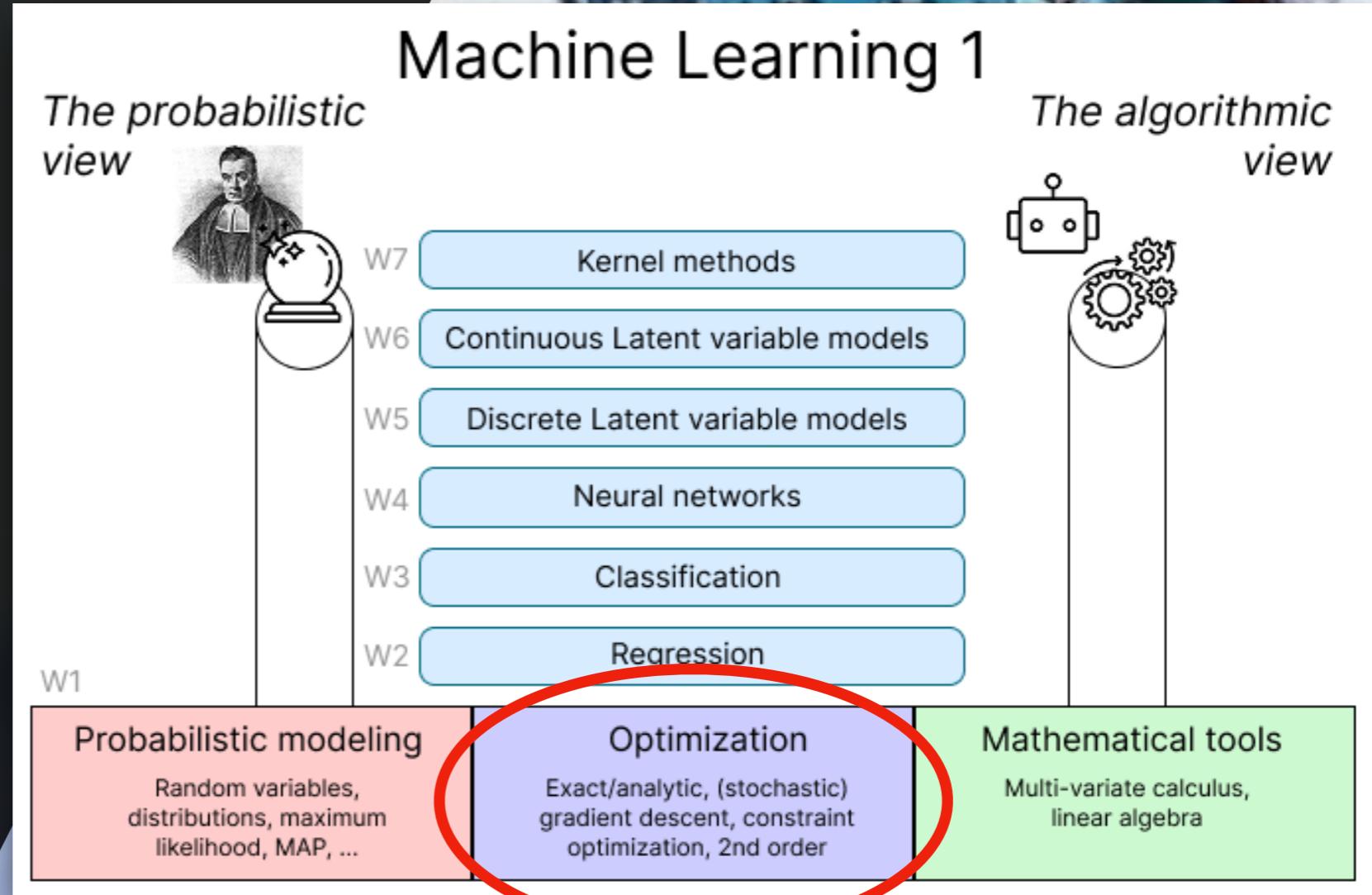
$$\mathbf{w}^{*\top} \underline{\boldsymbol{\phi}} = 0$$

Machine Learning 1

Lecture 7.3 - Supervised Learning
Classification - Logistic Regression:
Stochastic Gradient Descent

Erik Bekkers

(Bishop 4.3.2)



Has very simple gradient! Logistic Regression (K=2): SGD

- Stochastic Gradient Descent for cross-entropy:

$$E(\mathbf{w}) = - \sum_{n=1}^N t_n \ln y_n + (1 - t_n) \ln(1 - y_n) = \sum_{n=1}^N E_n(\mathbf{w})$$

$\hookrightarrow y_n = \sigma(\mathbf{w}^T \phi)$

- Update rule given a random data point (\mathbf{x}_n, t_n)

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n(\mathbf{w}^{(\tau)})^T$$

- Gradient: $\nabla E_n(\mathbf{w}) = \left(\frac{\partial E_n(\mathbf{w})}{\partial w_0}, \dots, \frac{\partial E_n(\mathbf{w})}{\partial w_{M-1}} \right)$

$$\frac{\partial E_n(\mathbf{w})}{\partial w_j} = \frac{\partial E_n(\mathbf{w})}{\partial y_n} \cdot \frac{\partial y_n}{\partial w_j} = \left(-\frac{t_n}{y_n} + \frac{1-t_n}{1-y_n} \right) \cdot \frac{\partial y_n}{\partial w_j}$$

$$\frac{\partial y_n}{\partial w_j} = \frac{\partial}{\partial w_j} \sigma(\mathbf{w}^T \phi_n)$$



Logistic Regression (K=2): SGD

Let's derive $\frac{\partial y_n}{\partial w_j} = \frac{\partial}{\partial w_j} \sigma(\mathbf{w}^T \phi_n)$ using $\frac{\partial \sigma(a)}{\partial a} = \sigma(a)(1 - \sigma(a))$

$$\begin{aligned} \frac{\partial}{\partial w_j} \sigma(\mathbf{w}^T \phi_n) &= \cancel{\sigma(\mathbf{w}^T \phi_n)}(1 - \sigma(\mathbf{w}^T \phi_n)) \frac{\partial(\mathbf{w}^T \phi_n)}{\partial w_j} \\ &= y_n (1 - y_n) \phi_{nj} \end{aligned}$$

$\frac{\partial E_n(\mathbf{w})}{\partial w_j} = -\frac{t_n}{y_n} \frac{\partial y_n}{\partial w_j} + \frac{1 - t_n}{1 - y_n} \frac{\partial y_n}{\partial w_j}$

$$\begin{aligned} &= -\frac{t_n}{y_n} y_n (1 - y_n) \phi_{nj} + \cancel{\frac{1 - t_n}{1 - y_n} y_n (1 - y_n) \phi_{nj}} \\ &= -t_n \phi_{nj} + \cancel{t_n y_n \phi_{nj}} + \cancel{y_n \phi_{nj}} - \cancel{t_n y_n \phi_{nj}} \\ &= (y_n - t_n) \phi_{nj} \end{aligned}$$

Logistic Regression (K=2): SGD

- Stochastic Gradient Descent for cross-entropy:

$$E(\mathbf{w}) = - \sum_{n=1}^N t_n \ln y_n + (1 - t_n) \ln(1 - y_n)$$

- Update rule given a random data point (\mathbf{x}_n, t_n)

Simple form b/c activation function f is the **link function!** (Bishop 4.3.6)

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n(\mathbf{w}^{(\tau)})^T$$

$$\nabla E_n(\mathbf{w}) = \left(\frac{\partial E_n(\mathbf{w})}{\partial w_0}, \dots, \frac{\partial E_n(\mathbf{w})}{\partial w_{M-1}} \right) = (y_n - t_n) \boldsymbol{\phi}_n$$

error + feature vector
 \sim

- Gradient Descent for LR has a perceptron-like update rule:

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta(y_n - t_n) \boldsymbol{\phi}_n$$

↳ .. or ..

Stochastic Gradient Descent

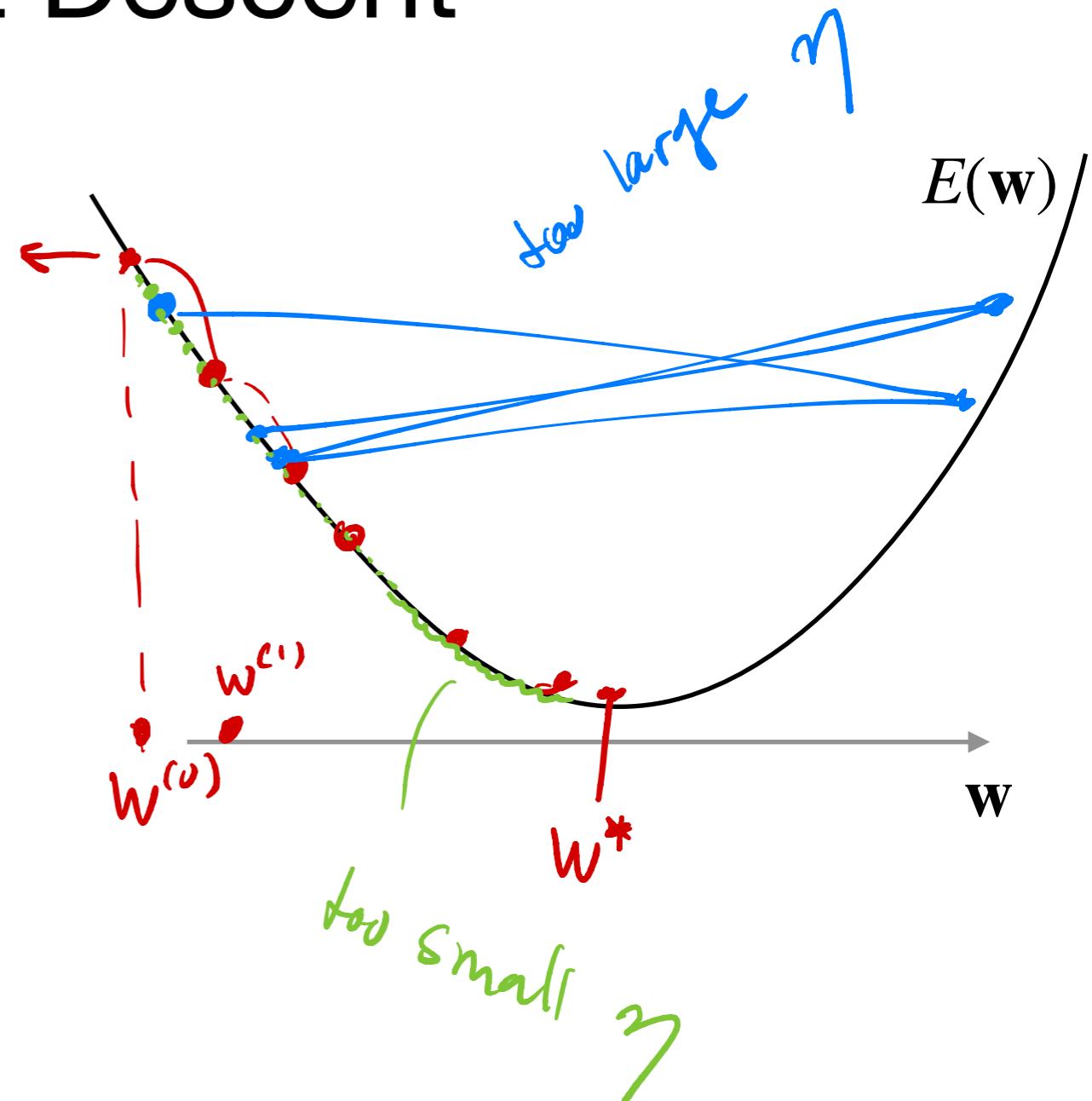
1. Initialize $\mathbf{w}^{(0)}$
2. Choose a learning rate η
3. While $\|\mathbf{w}^{(\tau+1)} - \mathbf{w}^{(\tau)}\| > \epsilon$

I. Choose a random data point (\mathbf{x}_n, t_n)

II. Update \mathbf{w} :

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta(y_n - t_n)\phi_n$$

- ➊ If η too large: no convergence
- ➋ If η too small: very slow convergence
- Converged \mathbf{w}^* : estimate of minimizer of $E(\mathbf{w})$!



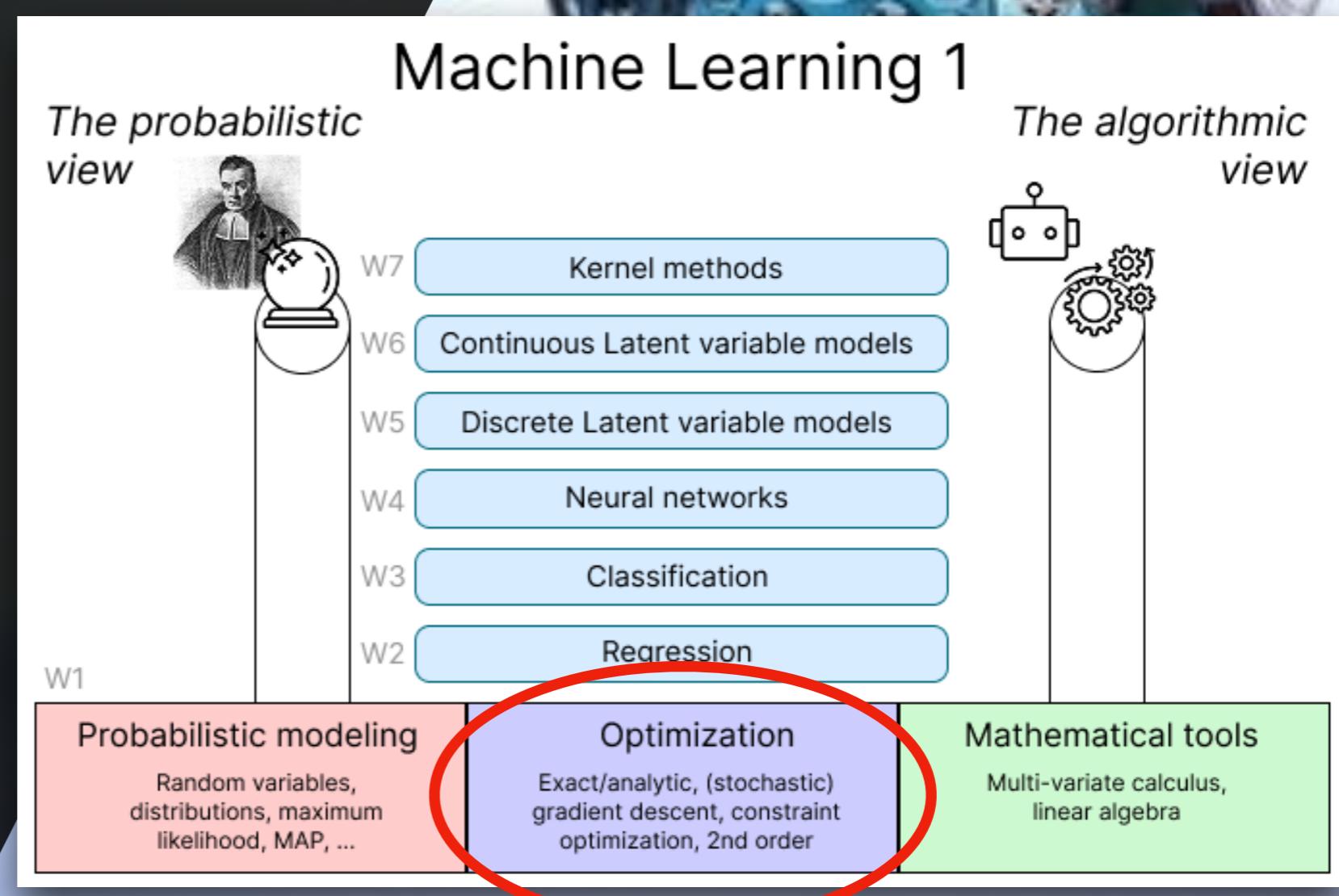
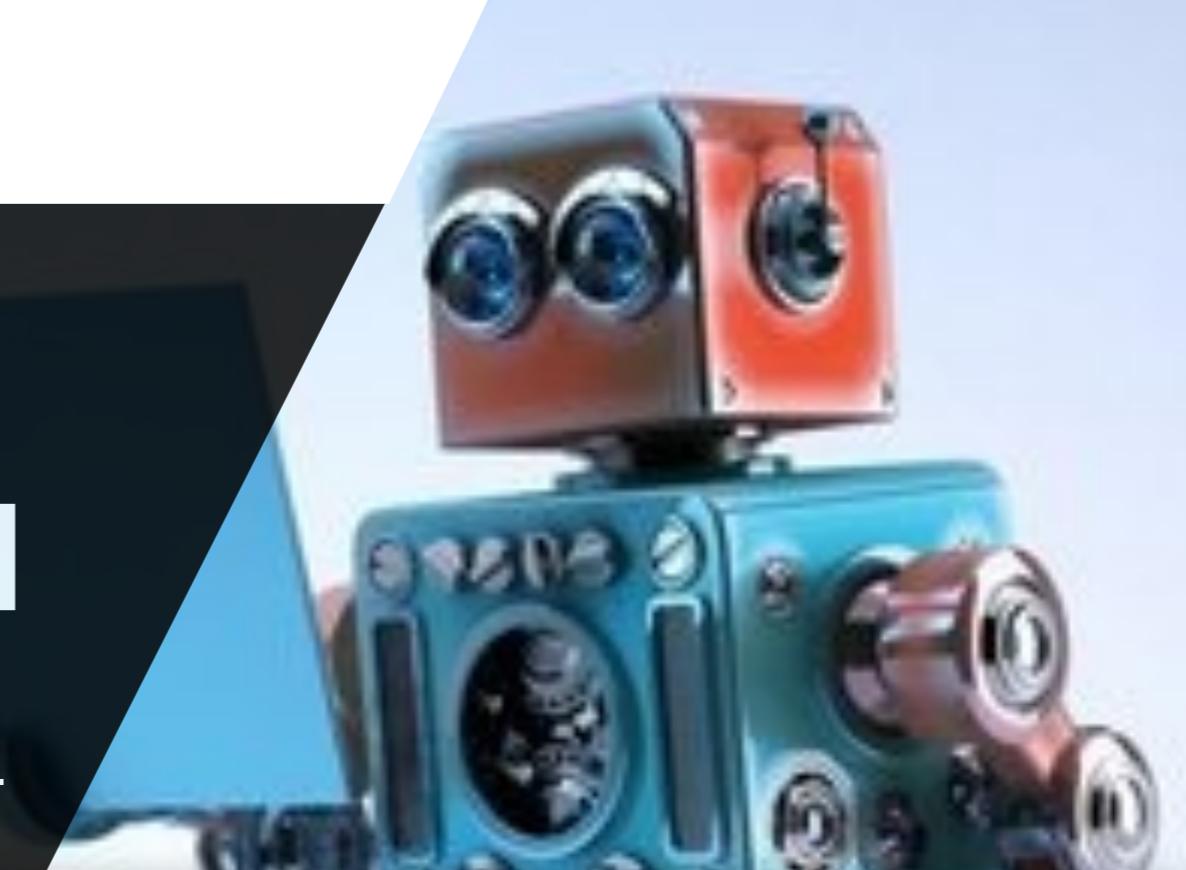
Machine Learning 1

Lecture 7.4 - Supervised Learning
Classification - Logistic Regression: Newton-Raphson Optimization

Erik Bekkers

(Bishop 4.3.3)
Not material for mid-term. But please watch video by the end of the course, or try to follow my annotations

Slide credits: Patrick Forré and Rianne van den Berg



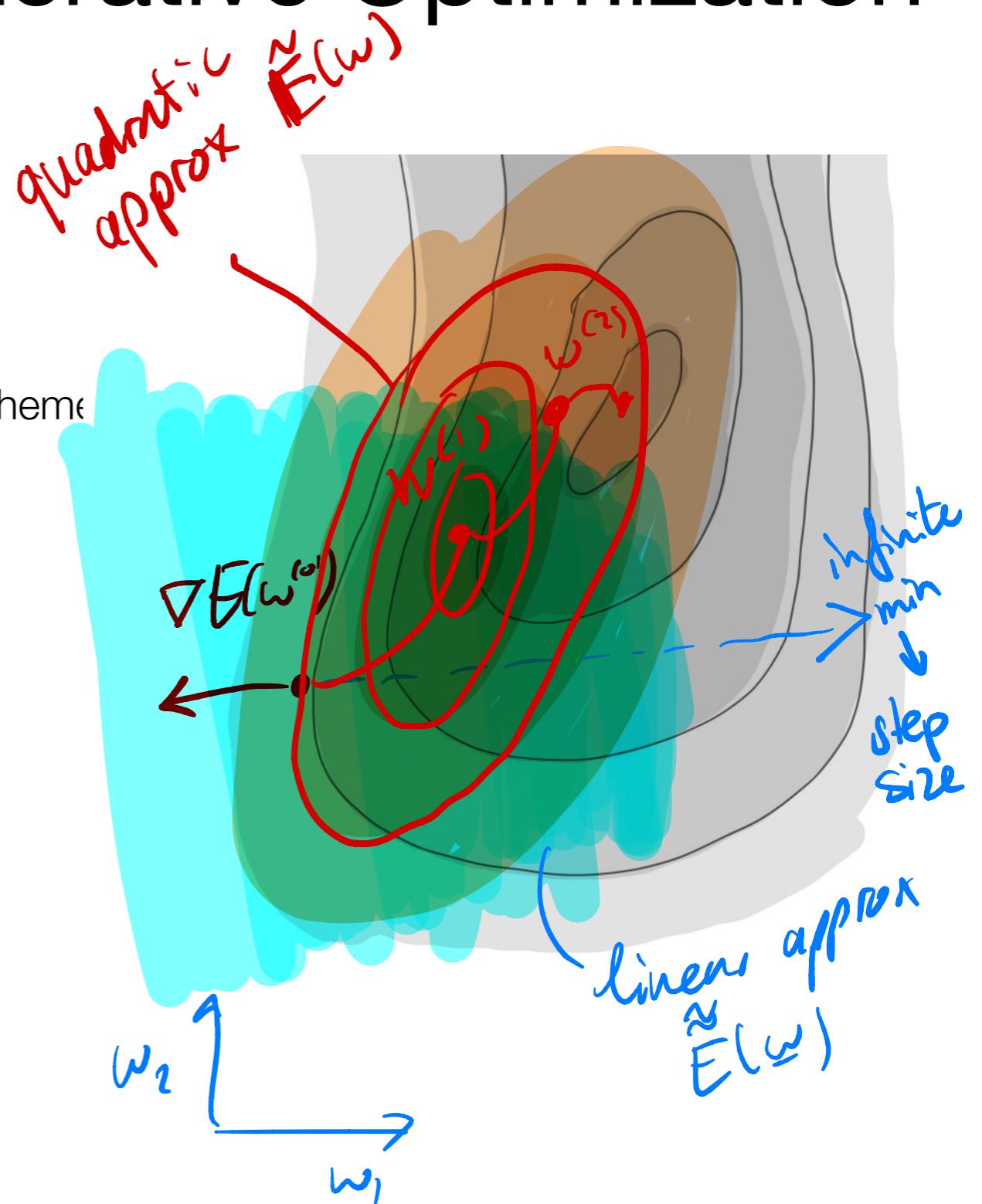
Newton-Raphson Iterative Optimization

- Goal: minimize

$$E(\mathbf{w}) = - \sum_{n=1}^N t_n \ln y_n + (1 - t_n) \ln(1 - y_n)$$

- Newton-Raphson iterative second order optimization scheme

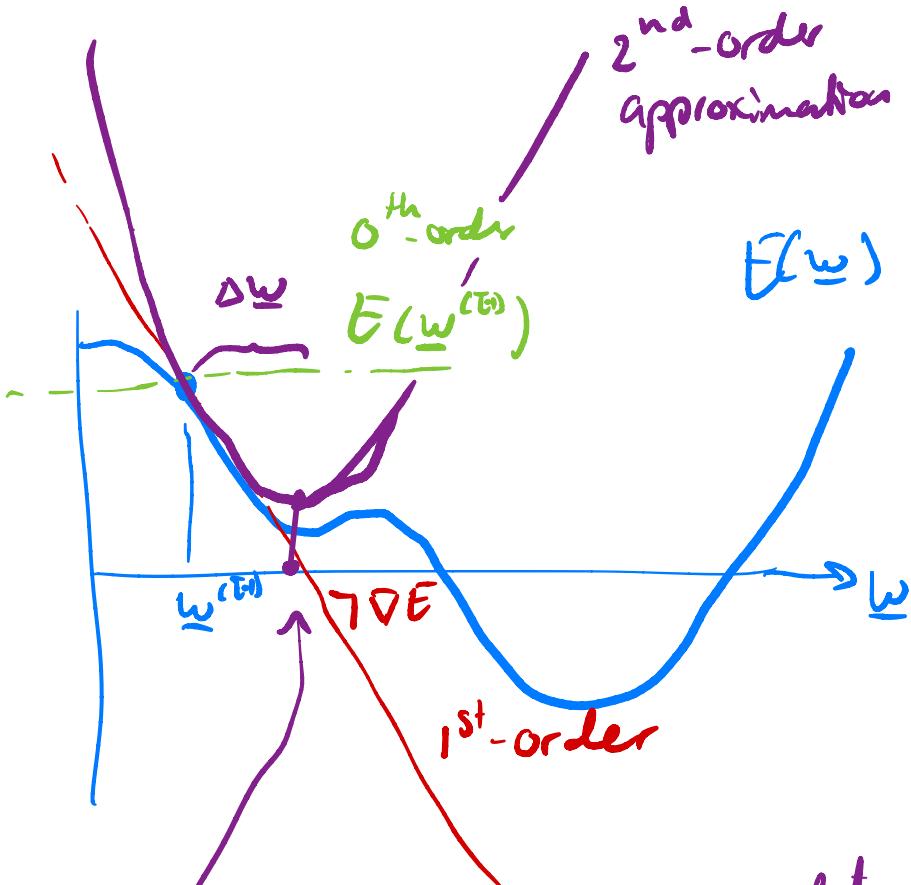
- Initial guess \mathbf{w}^0
- For $\tau = 1, \dots$:
 - Approximate $E(\mathbf{w})$ with a quadratic function $\tilde{E}_{\tau-1}(\mathbf{w})$ around $\mathbf{w}^{(\tau-1)}$
 - Set $\mathbf{w}^{(\tau)}$ such that it minimizes $\tilde{E}(\mathbf{w})$, i.e.,
$$\mathbf{w}^{(\tau)} = \underset{\mathbf{w}}{\operatorname{argmin}} \tilde{E}_{\tau-1}(\mathbf{w})$$
 - Stop when $\|\mathbf{w}^{(\tau+1)} - \mathbf{w}^{(\tau)}\| = 0$
- You have found \mathbf{w}^* such that $\frac{\partial}{\partial \mathbf{w}} E(\mathbf{w}^*) = 0$



Taylor expansion:

$$f(x_0 + \Delta x) \approx f(x_0) + \Delta x \frac{\partial f}{\partial x} \Big|_{x_0} + \frac{1}{2} \Delta x^2 \frac{\partial^2 f}{\partial x^2} \Big|_{x_0} + \dots$$

higher order derivatives



well-defined minimum so let

$$\begin{aligned} w^{(t)} &= \arg \min_{w^{(t-1)} + \Delta w} E_{w^{(t)}}(w) \\ &= w^{(t-1)} + \Delta w \end{aligned}$$

Newton-Raphson Iterative Optimization

- Given your old estimate $\mathbf{w}^{(\tau-1)}$, approximate $E(\mathbf{w})$ with a second order Taylor expansion around $\mathbf{w}^{(\tau-1)}$

$$E(\mathbf{w}) \approx \tilde{E}(\mathbf{w}^{(\tau-1)} + \Delta\mathbf{w}) = \underbrace{E(\mathbf{w}^{(\tau-1)})}_{0^{\text{th}}\text{-order}} + \underbrace{(\Delta\mathbf{w})^T \nabla E(\mathbf{w}^{(\tau-1)})}_{\text{linear}} + \frac{1}{2} (\Delta\mathbf{w})^T \mathbf{H} \Delta\mathbf{w} \quad \underbrace{+}_{\text{quadratic}}$$

- minimizing*
- Gradient (first order derivatives):

$$\nabla E_n(\mathbf{w}) = \left(\frac{\partial E_n(\mathbf{w})}{\partial w_0}, \dots, \frac{\partial E_n(\mathbf{w})}{\partial w_{M-1}} \right)$$

- Hessian Matrix (second order derivatives):

$$H_{ij} = \frac{\partial^2 E(\mathbf{w})}{\partial w_i \partial w_j}$$

- Choose $\Delta\mathbf{w}$ such that next estimate $\mathbf{w}^{(\tau)} = \mathbf{w}^{(\tau-1)} + \Delta\mathbf{w}$ minimizes $\tilde{E}(\mathbf{w})$:

$$\frac{\partial}{\partial \Delta\mathbf{w}} \tilde{E}(\mathbf{w}^{(\tau-1)} + \Delta\mathbf{w}) = \nabla E(\mathbf{w}^{(\tau-1)}) + (\Delta\mathbf{w})^T H = 0 \rightarrow H \Delta\mathbf{w} = -\nabla E$$

$$\Delta\mathbf{w} = -H^{-1} \nabla E$$

- Update rule:

$$\mathbf{w}^{(\tau)} = \mathbf{w}^{(\tau-1)} - H^{-1} \nabla E(\mathbf{w}^{(\tau-1)})$$

Optimal step size!

Newton-Raphson Iterative Optimization

- Update rule: $\mathbf{w}^{(\tau)} = \mathbf{w}^{(\tau-1)} - \mathbf{H}^{-1} \nabla^T E(\mathbf{w}^{(\tau-1)})$
- Gradient $\nabla E_n(\mathbf{w}) = \left(\frac{\partial E_n(\mathbf{w})}{\partial w_0}, \dots, \frac{\partial E_n(\mathbf{w})}{\partial w_{M-1}} \right) = (y_n - t_n) \boldsymbol{\phi}_n$

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \boldsymbol{\phi}_n =$$

$$\left(E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w}) \right)$$

Newton-Raphson Iterative Optimization

- Update rule: $\mathbf{w}^{(\tau)} = \mathbf{w}^{(\tau-1)} - \mathbf{H}^{-1} \nabla^T E(\mathbf{w}^{(\tau-1)})$
- Gradient $\nabla E_n(\mathbf{w}) = \left(\frac{\partial E_n(\mathbf{w})}{\partial w_0}, \dots, \frac{\partial E_n(\mathbf{w})}{\partial w_{M-1}} \right) = (y_n - t_n) \boldsymbol{\phi}_n$

all data $\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \boldsymbol{\phi}_n = \Phi^T (\mathbf{y} - \mathbf{t})$ $\left(E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w}) \right)$

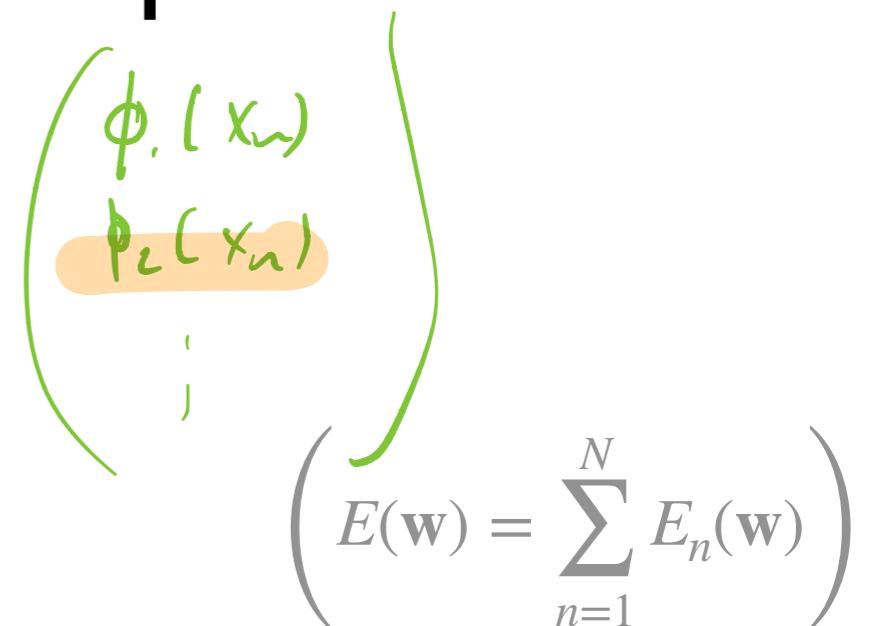
feat vec error

- Hessian $H_{ij} = \frac{\partial^2 E(\mathbf{w}^{(\tau-1)})}{\partial w_i \partial w_j}$

Newton-Raphson Iterative Optimization

- Update rule: $\mathbf{w}^{(\tau)} = \mathbf{w}^{(\tau-1)} - \mathbf{H}^{-1} \nabla^T E(\mathbf{w}^{(\tau-1)})$
- Gradient $\nabla E_n(\mathbf{w}) = \left(\frac{\partial E_n(\mathbf{w})}{\partial w_0}, \dots, \frac{\partial E_n(\mathbf{w})}{\partial w_{M-1}} \right) = (y_n - t_n) \boldsymbol{\phi}_n^T$

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \boldsymbol{\phi}_n = \boldsymbol{\Phi}^T (\mathbf{y} - \mathbf{t})$$


$$E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w})$$

- Hessian $H_{ij} = \frac{\partial^2 E(\mathbf{w}^{(\tau-1)})}{\partial w_i \partial w_j} = \frac{\partial}{\partial w_i} \sum_{n=1}^N (y_n - t_n) \boldsymbol{\phi}_j(\mathbf{x}_n)$

$\underline{\mathbf{w}}^T \underline{\boldsymbol{\phi}}_n$



Newton-Raphson Iterative Optimization

- Update rule: $\mathbf{w}^{(\tau)} = \mathbf{w}^{(\tau-1)} - \mathbf{H}^{-1} \nabla^T E(\mathbf{w}^{(\tau-1)})$

- Gradient $\nabla E_n(\mathbf{w}) = \left(\frac{\partial E_n(\mathbf{w})}{\partial w_0}, \dots, \frac{\partial E_n(\mathbf{w})}{\partial w_{M-1}} \right) = (y_n - t_n) \boldsymbol{\phi}_n$

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \boldsymbol{\phi}_n = \mathbf{\Phi}^T (\mathbf{y} - \mathbf{t})$$

$$\left(E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w}) \right)$$

- Hessian $H_{ij} = \frac{\partial^2 E(\mathbf{w}^{(\tau-1)})}{\partial w_i \partial w_j} = \frac{\partial}{\partial w_i} \sum_{n=1}^N (y_n - t_n) \boldsymbol{\phi}_j(\mathbf{x}_n) = \sum_{n=1}^N \boldsymbol{\phi}_j(\mathbf{x}_n) \frac{\partial y_n}{\partial w_i}$

Newton-Raphson Iterative Optimization

- Update rule: $\mathbf{w}^{(\tau)} = \mathbf{w}^{(\tau-1)} - \mathbf{H}^{-1} \nabla^T E(\mathbf{w}^{(\tau-1)})$

- Gradient $\nabla E_n(\mathbf{w}) = \left(\frac{\partial E_n(\mathbf{w})}{\partial w_0}, \dots, \frac{\partial E_n(\mathbf{w})}{\partial w_{M-1}} \right) = (y_n - t_n) \boldsymbol{\phi}_n$

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \boldsymbol{\phi}_n = \boldsymbol{\Phi}^T (\mathbf{y} - \mathbf{t})$$

$$\left(E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w}) \right)$$

- Hessian $H_{ij} = \frac{\partial^2 E(\mathbf{w}^{(\tau-1)})}{\partial w_i \partial w_j} = \frac{\partial}{\partial w_i} \sum_{n=1}^N (y_n - t_n) \boldsymbol{\phi}_j(\mathbf{x}_n) = \sum_{n=1}^N \boldsymbol{\phi}_j(\mathbf{x}_n) \frac{\partial y_n}{\partial w_i}$

$$= \sum_{n=1}^N y_n (1 - y_n) \boldsymbol{\phi}_i \boldsymbol{\phi}_j$$

"weight per data point"

$$y_n = \sigma(\boldsymbol{\Phi}_n^T \mathbf{w})$$

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

Newton-Raphson Iterative Optimization

- Update rule: $\mathbf{w}^{(\tau)} = \mathbf{w}^{(\tau-1)} - \mathbf{H}^{-1} \nabla^T E(\mathbf{w}^{(\tau-1)})$
- Gradient $\nabla E_n(\mathbf{w}) = \left(\frac{\partial E_n(\mathbf{w})}{\partial w_0}, \dots, \frac{\partial E_n(\mathbf{w})}{\partial w_{M-1}} \right) = (y_n - t_n) \boldsymbol{\phi}_n$

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \boldsymbol{\phi}_n = \boldsymbol{\Phi}^T (\mathbf{y} - \mathbf{t}) \quad \left(E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w}) \right)$$

Hessian $H_{ij} = \frac{\partial^2 E(\mathbf{w}^{(\tau-1)})}{\partial w_i \partial w_j} = \frac{\partial}{\partial w_i} \sum_{n=1}^N (y_n - t_n) \phi_j(\mathbf{x}_n) = \sum_{n=1}^N \phi_j(\mathbf{x}_n) \frac{\partial y_n}{\partial w_i}$

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

$R = \text{diag}(\{y_n(1-y_n)\})$

$\mathbf{H} = \sum_{n=1}^N y_n(1-y_n) \boldsymbol{\phi}_n \boldsymbol{\phi}_n^T = \boldsymbol{\Phi}^T R \boldsymbol{\Phi} \quad \begin{cases} R_{nn} = y_n(1-y_n) \\ R_{nm} = 0 \quad \forall n \neq m \end{cases}$

Newton-Raphson Iterative Optimization

- Update rule: $\mathbf{w}^{(\tau)} = \mathbf{w}^{(\tau-1)} - \mathbf{H}^{-1} \nabla^T E(\mathbf{w}^{(\tau-1)})$
- Gradient $\nabla E_n(\mathbf{w}) = \left(\frac{\partial E_n(\mathbf{w})}{\partial w_0}, \dots, \frac{\partial E_n(\mathbf{w})}{\partial w_{M-1}} \right) = (y_n - t_n) \boldsymbol{\phi}_n$

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \boldsymbol{\phi}_n = \boldsymbol{\Phi}^T (\mathbf{y} - \mathbf{t})$$

$$\left(E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w}) \right)$$

- Hessian $H_{ij} = \frac{\partial E(\mathbf{w}^{(\tau-1)})}{\partial w_i \partial w_j} = \frac{\partial}{\partial w_i} \sum_{n=1}^N (y_n - t_n) \boldsymbol{\phi}_j(\mathbf{x}_n) = \sum_{n=1}^N \boldsymbol{\phi}_j(\mathbf{x}_n) \frac{\partial y_n}{\partial w_i}$

$$= \sum_{n=1}^N y_n(1 - y_n) \boldsymbol{\phi}_i \boldsymbol{\phi}_j$$

$$\mathbf{H} = \sum_{n=1}^N y_n(1 - y_n) \boldsymbol{\phi}_n \boldsymbol{\phi}_n^T = \boldsymbol{\Phi}^T \mathbf{R} \boldsymbol{\Phi}$$

$$\Delta \mathbf{w} = -\mathbf{H}^{-1} \underline{\nabla E^T}$$

Iterative Reweighted Least Squares

- Update rule: $\mathbf{w}^{(\tau)} = \mathbf{w}^{(\tau-1)} - \mathbf{H}^{-1} \nabla^T E(\mathbf{w}^{(\tau-1)})$ optimal $\Delta \mathbf{w}$ give the approx.
- With Gradient: $\nabla E(\mathbf{w}) = \Phi^T(\mathbf{y} - \mathbf{t})$
- With Hessian: $\mathbf{H} = \Phi^T \mathbf{R} \Phi$, with $\mathbf{R} = \text{diag}_{N \times N}\{y_n(1 - y_n)\}$
- Newton-Raphson updates solutions to a weighted least squares problem \rightarrow iterative reweighted least squares:

$$\begin{aligned}
 \mathbf{w}^{(\tau)} &= \mathbf{w}^{(\tau-1)} - (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T (\mathbf{y} - \mathbf{t}) \\
 &= (\Phi^T \mathbf{R} \Phi)^{-1} \{ (\Phi^T \mathbf{R} \Phi) \mathbf{w}^{(\tau-1)} - \Phi^T (\mathbf{y} - \mathbf{t}) \} \\
 &= (\Phi^T \mathbf{R} \Phi)^{-1} \{ \Phi^T \mathbf{R} (\mathbf{z} + \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t})) - \Phi^T (\mathbf{y} - \mathbf{t}) \} \\
 &= (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{z} \quad \text{with } \mathbf{z} = \Phi \mathbf{w}^{(\tau-1)} - \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t})
 \end{aligned}$$

- Note similarity with ML solution to linear regression:

$$\begin{aligned}
 \mathbf{w}_{ML} &= (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} && \text{(solution to least squares regression } \mathbf{w}_{ML} = \underset{\mathbf{w}}{\text{argmin}} \sum_{n=1}^N (\mathbf{w}^T \boldsymbol{\phi}_n - t_n)^2 \text{)} \\
 \mathbf{w}^{(\tau)} &= (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{z} && \text{(solution to a } \text{weighted} \text{ least squares problem } \mathbf{w}^{(\tau)} = \underset{\mathbf{w}}{\text{argmin}} \sum_{n=1}^N r_i (\mathbf{w}^T \boldsymbol{\phi}_n - z_n)^2
 \end{aligned}$$

SGD vs Newton-Raphson

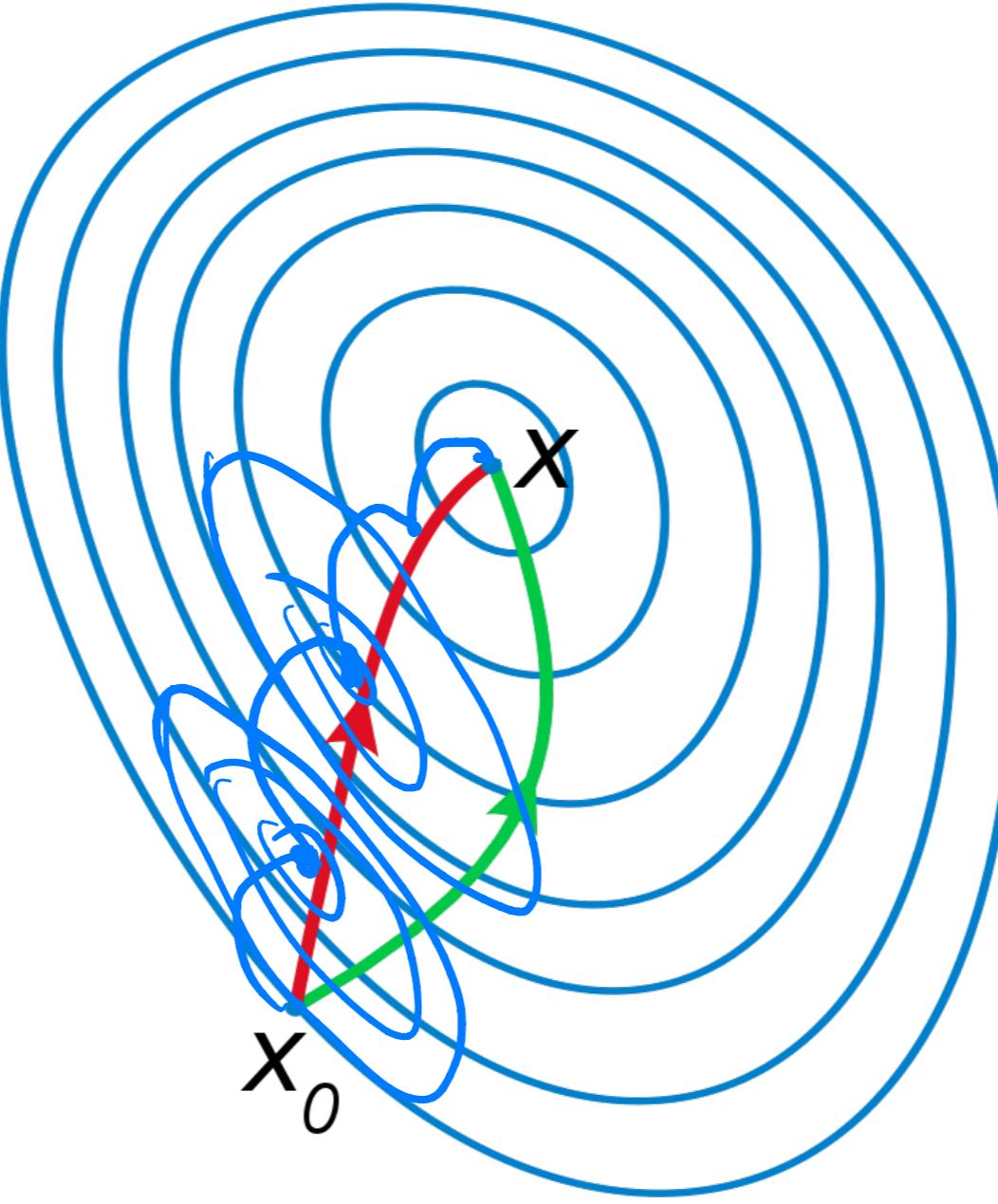


Figure: Green: gradient descent, which always goes in the direction of steepest descent. Red: Newton-Raphson's procedure, which takes into account curvature to take a more direct path. (Wikipedia)

The cross entropy loss is convex

- ▶ The error function $E(\mathbf{w})$ is convex since the Hessian is positive definite
- ▶ Proof Hessian is convex:

The Hessian of $E(\mathbf{w})$ is given by $\mathbf{H} = \Phi^T \mathbf{R} \Phi$, with $\mathbf{R} = \text{diag}_{N \times N}\{y_n(1 - y_n)\}$.

To show: $\forall_{\mathbf{w} \neq \mathbf{0} \in \mathbb{R}^M} : \mathbf{w}^T \mathbf{H} \mathbf{w} > 0$.

$$\mathbf{w}^T \mathbf{H} \mathbf{w} = \mathbf{w}^T \Phi^T \mathbf{R} \Phi \mathbf{w}$$

$$\left(\mathbf{R} = \sqrt{\mathbf{R}} \sqrt{\mathbf{R}} \right)$$

$$= \mathbf{w}^T \Phi^T \sqrt{\mathbf{R}} \sqrt{\mathbf{R}} \Phi \mathbf{w}$$

$$\left(\sqrt{\mathbf{R}} = \text{diag}(\sqrt{y_n(1 - y_n)}) \right)$$

$$= (\sqrt{\mathbf{R}} \Phi \mathbf{w})^T (\sqrt{\mathbf{R}} \Phi \mathbf{w})$$

$$= \|\sqrt{\mathbf{R}} \Phi \mathbf{w}\|^2 > 0$$

(since $\mathbf{w} \neq \mathbf{0}$)