

# Machine Intelligence 1

## 1.4 Additional Topics

Prof. Dr. Klaus Obermayer

Fachgebiet Neuronale Informationsverarbeitung (NI)

WS 2016/2017

## 1.4.1 Stochastic Approximation and Online Learning

# Online learning

$$\Delta \mathbf{w}_{ij}^{v'v} = -\eta \frac{\partial E_{[\mathbf{w}]}^T}{\partial \mathbf{w}_{ij}^{v'v}} = -\eta \frac{1}{p} \sum_{\alpha=1}^p \frac{\partial e_{[\mathbf{w}]}^{(\alpha)}}{\partial \mathbf{w}_{ij}^{v'v}}$$

# Online learning

$$\Delta \mathbf{w}_{ij}^{v'v} = -\eta \frac{\partial E_{[\mathbf{w}]}^T}{\partial \mathbf{w}_{ij}^{v'v}} = -\eta \frac{1}{p} \sum_{\alpha=1}^p \frac{\partial e_{[\mathbf{w}]}^{(\alpha)}}{\partial \mathbf{w}_{ij}^{v'v}}$$

**while** *convergence criterion not met* **do**

select the a data point:  $(\underline{\mathbf{x}}^{(\alpha)}, y_T^{(\alpha)})$

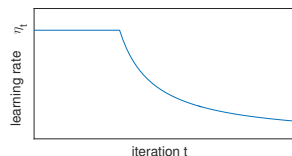
change weights according to:  $(\Delta \mathbf{w}_{ij}^{v'v})^{(t+1)} = -\eta_t \frac{\partial e_{[\mathbf{w}]}^{(\alpha)}}{\partial \mathbf{w}_{ij}^{v'v}}^{(t)}$

update learning rate  $\eta_t$

**end**

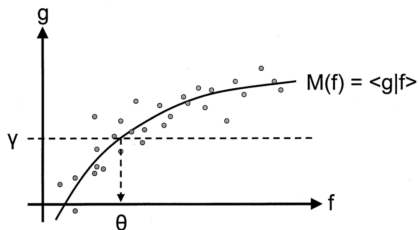
## ■ Adaptive learning rate

- first constant  $\eta_t = \eta_0$
- then decaying  $\eta_t = \frac{\eta_0}{t}$



# Convergence in online learning (1)

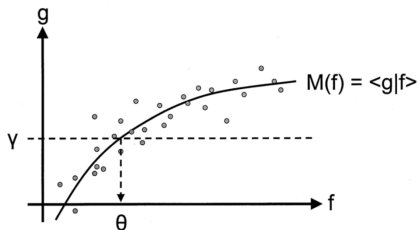
## Stochastic Approximation (Robbins & Munro, 1951)



$g, f$ : correlated random variables

# Convergence in online learning (1)

## Stochastic Approximation (Robbins & Munro, 1951)



$g, f$ : correlated random variables

## Application to online learning:

- $f$  are model parameters  $\underline{w}$
- $M(f)$  are batch-gradients  $\left. \frac{\partial E^T}{\partial \underline{w}} \right|_{[\underline{w}]}$
- $g$  are individual gradients  $\left. \frac{\partial e^\alpha}{\partial \underline{w}} \right|_{[\underline{w}]}$
- $\theta$  are optimal parameters  $\underline{w}^*$ , i.e.

$$\left. \frac{\partial E^T}{\partial \underline{w}} \right|_{[\underline{w}^*]} = M(\theta) \stackrel{!}{=} \gamma = 0$$

# Convergence in online learning (2)

Let  $g_t, f_t \in \mathbb{R}$ ,  $g_t|f_t$ , be correlated random variables and let the initial  $f_1$  be an arbitrary real number. If

- $g_t$  are bounded
- $M(f)$  is monotonously increasing
- $M(\theta) = \gamma$  and  $\frac{\partial M(\theta)}{\partial f} > 0$
- and for the learning rates  $\eta_t$  holds

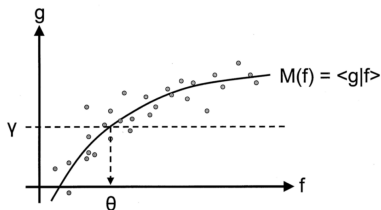
$$\sum_{t=1}^{\infty} \eta_t = \infty \quad \text{and} \quad \sum_{t=1}^{\infty} \eta_t^2 < \infty,$$

then the sequence

$$f_{t+1} = f_t + \eta_t(\gamma - g_t), \quad \eta_t > 0,$$

converges to  $\theta$  in the sense  $\lim_{t \rightarrow \infty} \langle (f_t - \theta)^2 \rangle = 0$ .

## Stochastic Approximation (Robbins & Munro, 1951)



$g, f$ : correlated random variables

# Convergence in online learning (2)

Let  $\frac{\partial e^\alpha}{\partial \underline{\mathbf{w}}} \big|_{[\underline{\mathbf{w}}_t]}, \underline{\mathbf{w}}_t \in \mathbb{R}^{|\mathcal{C}|}$  be correlated random variables and let the initial  $\underline{\mathbf{w}}_1$  be an arbitrary real vector. If

- $\frac{\partial e^\alpha}{\partial \underline{\mathbf{w}}} \big|_{[\underline{\mathbf{w}}_t]}$  are bounded
- $E^T$  is a **convex function**
- $\frac{\partial E^T}{\partial \underline{\mathbf{w}}} \big|_{[\underline{\mathbf{w}}^*]} = 0$  and  $\frac{\partial^2 E^T[\underline{\mathbf{w}}^*]}{\partial^2 \underline{\mathbf{w}}} > 0$  (min.  $\underline{\mathbf{w}}^*$ )
- and for the learning rates  $\eta_t$  holds

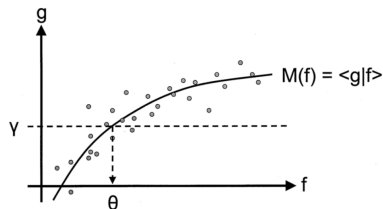
$$\sum_{t=1}^{\infty} \eta_t = \infty \quad \text{and} \quad \sum_{t=1}^{\infty} \eta_t^2 < \infty,$$

then the sequence

$$\underline{\mathbf{w}}_{t+1} = \underline{\mathbf{w}}_t - \eta_t \frac{\partial e^\alpha}{\partial \underline{\mathbf{w}}} \big|_{[\underline{\mathbf{w}}_t]}, \quad \eta_t > 0,$$

converges to  $\underline{\mathbf{w}}^*$  in the sense  $\lim_{t \rightarrow \infty} \langle (\underline{\mathbf{w}}_t - \underline{\mathbf{w}}^*)^2 \rangle = 0$ .

Stochastic Approximation  
(Robbins & Munro, 1951)

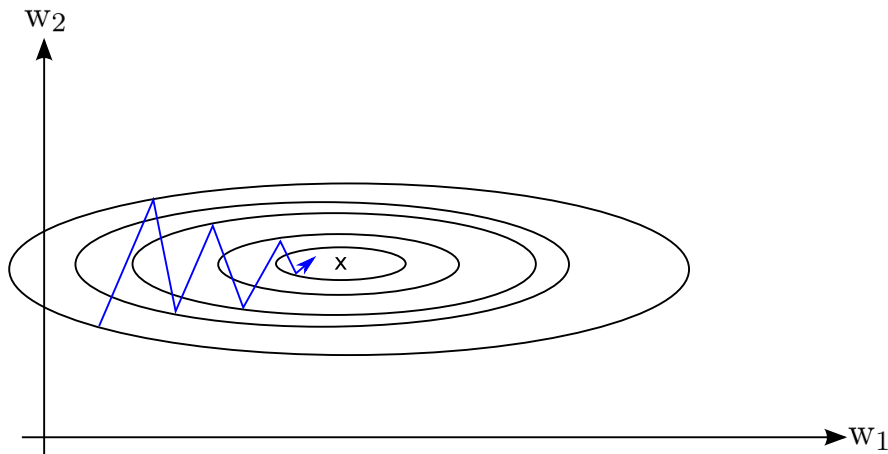


$g, f$ : correlated random variables

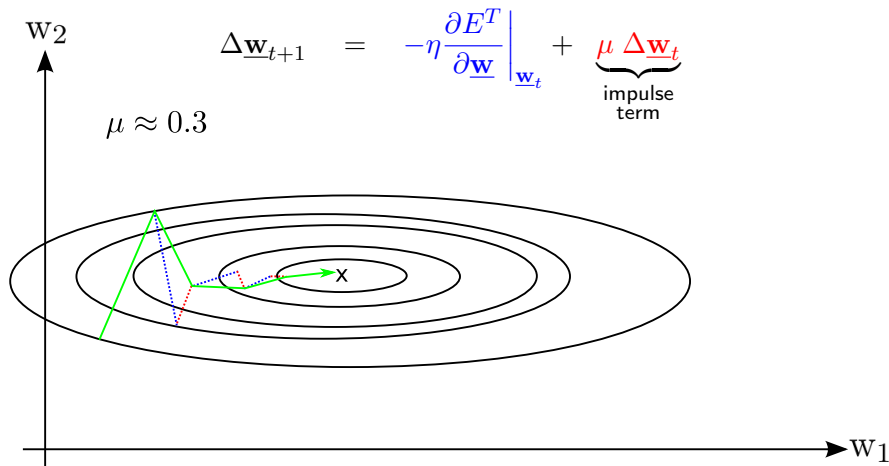


## 1.4.2 Improving Gradient-Descent Optimization

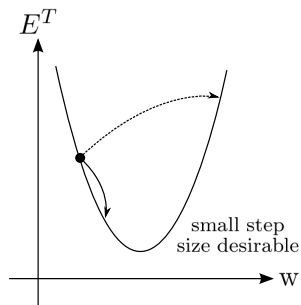
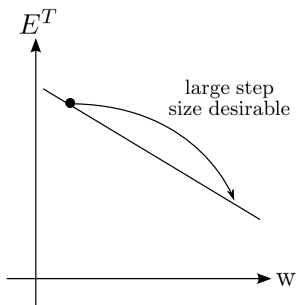
# Impulse terms



# Impulse terms



# Adaptive step size

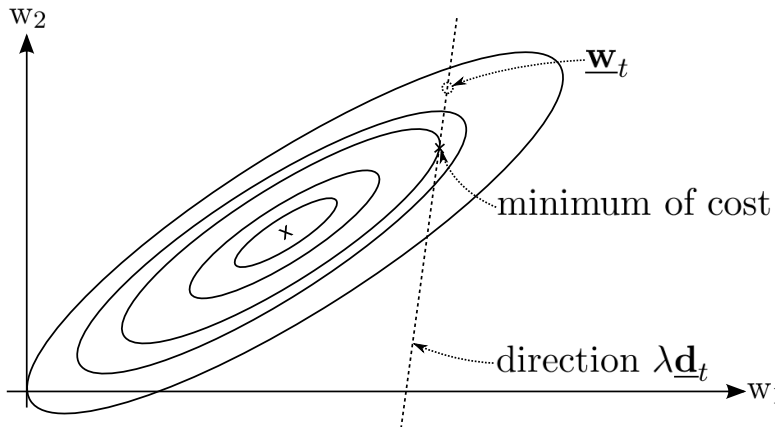


$$\eta_{t+1} = \begin{cases} \rho\eta_t, & \text{if } \Delta E^T < 0, \\ \delta\eta_t, & \text{if } \Delta E^T > 0, \end{cases} \quad \begin{array}{l} \text{increase step size, if } E^T \downarrow \\ \text{decrease step size, if } E^T \uparrow \end{array}$$

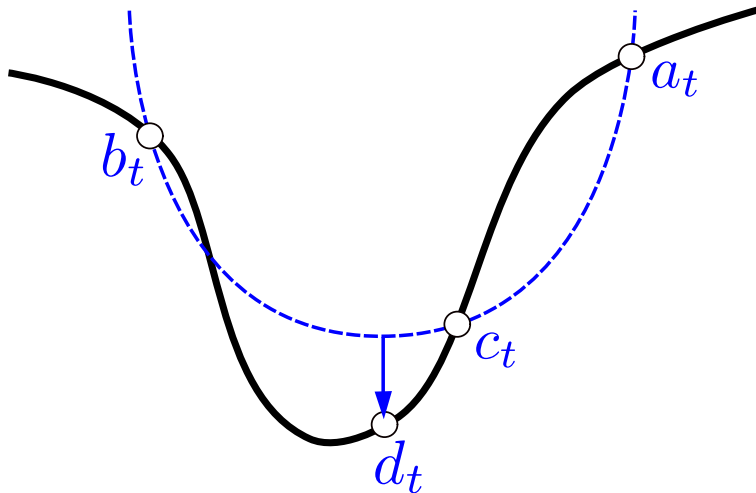
typical values:  $\rho = 1.1, \delta = 0.5$

## 1.4.3 The Conjugate Gradient Method

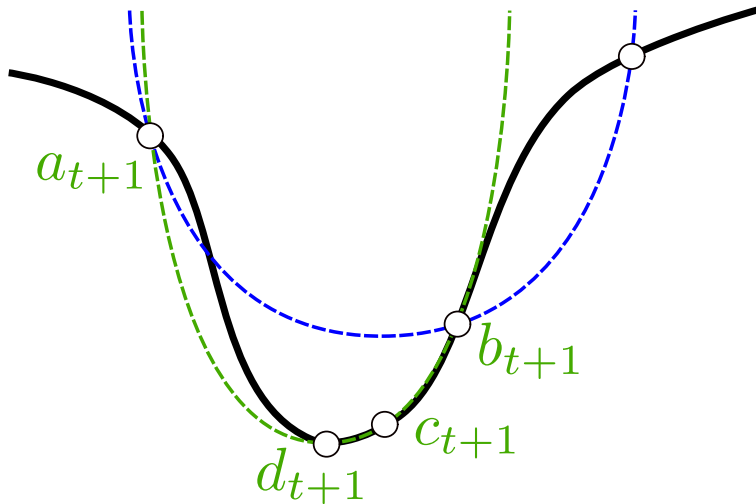
# Optimal step size



# Parabolic interpolation



## Parabolic interpolation





## Line search

### Line search using successive parabolic interpolation

**Initialization:**  $a_0, b_0, c_0$  (on  $\lambda \underline{d}_t$ );  $E_{(a_0)}^T, E_{(b_0)}^T > E_{(c_0)}^T$

**while** *stopping criterion not fulfilled* **do**

Fit a parabola through the three points  $a_t, b_t, c_t$

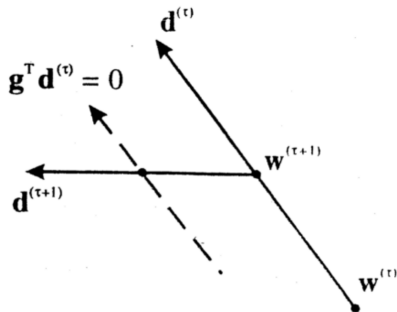
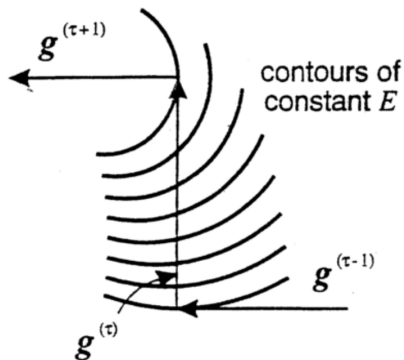
Calculate location  $d_t$  of its minimum

Set  $c_{t+1} = d_t, \quad b_{t+1} = c_t, \quad a_{t+1} = \begin{cases} a_t, & E_{(a_t)}^T < E_{(b_t)}^T \\ b_t, & \text{else} \end{cases}$

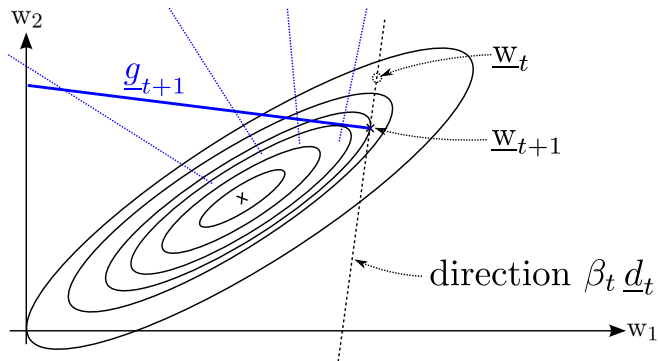
**end**

*For details and implementation see e.g. Numerical Recipes, 2nd edition, Chapter 10.2.*

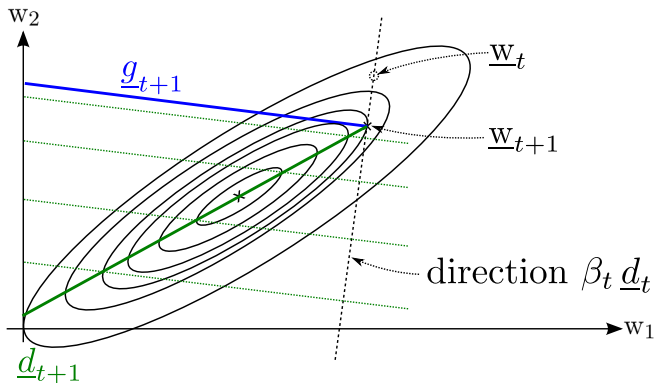
# The conjugate direction



# Parabolic cost function



# Parabolic cost function



$$\underline{d}_t^\top \underline{H} \underline{d}_{t+1} \stackrel{!}{=} 0, \quad H_{ij} := \frac{\partial^2 E^T}{\partial w_i \partial w_j} \quad (\text{Hebb Matrix})$$

# Adaptive momentum: the Polak-Ribiere rule

$$\underline{\mathbf{g}}_{t+1} := \left. \frac{\partial E^T}{\partial \underline{\mathbf{w}}} \right|_{\underline{\mathbf{w}}_{t+1}} \quad (\text{gradient at } \underline{\mathbf{w}}_{t+1})$$

$$\underline{\mathbf{d}}_{t+1} = -\underline{\mathbf{g}}_{t+1} + \beta_t \underline{\mathbf{d}}_t \quad (\text{conjugate direction})$$

$$\beta_t = \frac{\underline{\mathbf{g}}_{t+1}^T (\underline{\mathbf{g}}_{t+1} - \underline{\mathbf{g}}_t)}{\underbrace{\underline{\mathbf{g}}_t^T \underline{\mathbf{g}}_t}_{\text{Polak-Ribiere rule}}} \quad (\text{"smart momentum"})$$

# Conjugate gradient descent algorithm

$$\mathbf{w}_{ij}(t+1) = \mathbf{w}_{ij}(t) + \eta_t \underline{\mathbf{d}}_{t-1}$$

**Initialization:**  $\underline{\mathbf{w}}$ ,  $\underline{\mathbf{d}} = -\underline{\mathbf{g}}$

**while** *stopping criterion not fulfilled* **do**

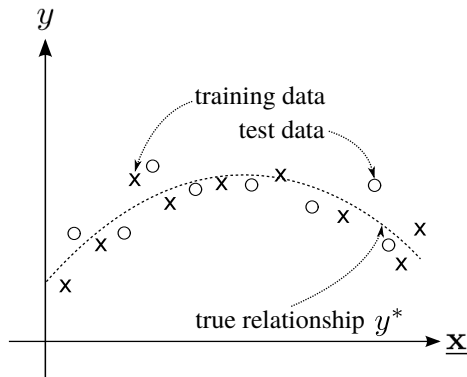
    Minimize  $E^T$  along  $\underline{\mathbf{d}}$  using line-search       $\rightarrow$  new  $\underline{\mathbf{w}}$

    Calculate the new conjugate direction       $\rightarrow$  new  $\underline{\mathbf{d}}$

**end**

## 1.4.4 Overfitting and Underfitting

# Overfitting and underfitting



data generation:

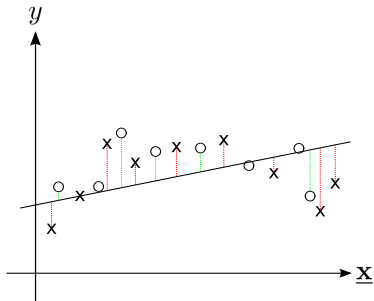
$$\text{e.g. } y_T = y^*_{(\underline{x})} + \underbrace{\eta}_{\text{noise}}$$

## Goal

Find a **good** model for  $y^*_{(.)}$  (explanation, prediction)



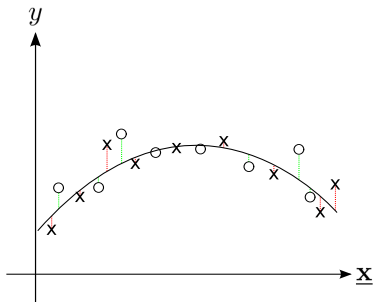
# Underfitting



## Diagnostics

$$\left. \begin{array}{l} E^T \text{ large} \\ E^G \text{ large} \end{array} \right\} E^T \approx E^G$$

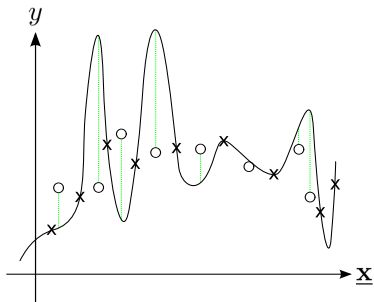
# Well fitted model



## Diagnostics

$$\left. \begin{array}{l} E^T \text{ small} \\ E^G \text{ small} \end{array} \right\} E^T \approx E^G$$

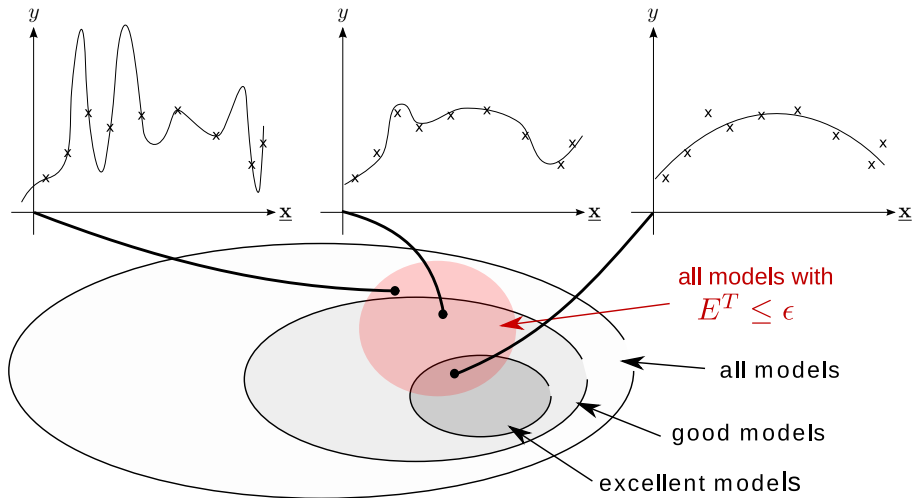
# Overfitting



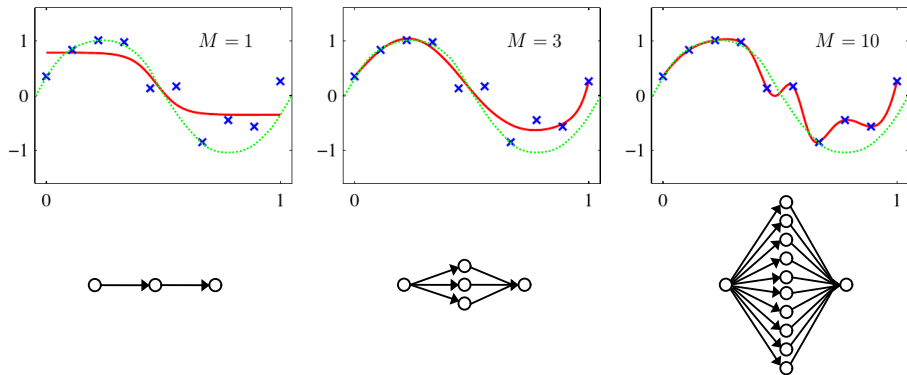
## Diagnostics

$$\left. \begin{array}{l} E^T \text{ small} \\ E^G \text{ large} \end{array} \right\} E^T \ll E^G$$

# Consequences



# Example with MLPs



MLP with  $M$  sigmoid hidden neurons  $f_j^1(x) = \tanh(x)$  and linear output neuron  $f_1^2(x) = x$ , from Bishop (2006)

## 1.4.5 Bias and Variance

# Bias and variance

## Example scenario

Observations:  $y_T = \underbrace{y^*_{(\underline{x})}}_{\text{true relationship}} + \underbrace{\eta}_{\text{noise}}, \quad \underline{x} \in \mathbb{R}^N, \quad y_T \in \mathbb{R}, \quad \eta \sim \mathcal{N}(0, \sigma_\eta)$

- $\underline{x}, y_T$  are random variables
- Many iid. datasets of equal length from  $P_{(y_T, \underline{x})} = P_{(y_T | \underline{x})} P_{(\underline{x})}$

# Bias and variance

## Example scenario

Observations:  $y_T = \underbrace{y^*_{(\underline{x})}}_{\text{true relationship}} + \underbrace{\eta}_{\text{noise}}, \quad \underline{x} \in \mathbb{R}^N, \quad y_T \in \mathbb{R}, \quad \eta \sim \mathcal{N}(0, \sigma_\eta)$

- $\underline{x}, y_T$  are random variables
- Many iid. datasets of equal length from  $P_{(y_T, \underline{x})} = P_{(y_T | \underline{x})} P_{(\underline{x})}$
- fitting one MLP to every dataset
- $\underline{w}$  (model parameters) are random variables
- $y_{(\underline{x}; \underline{w})}$  (predicted values) are random variables



# Bias and variance

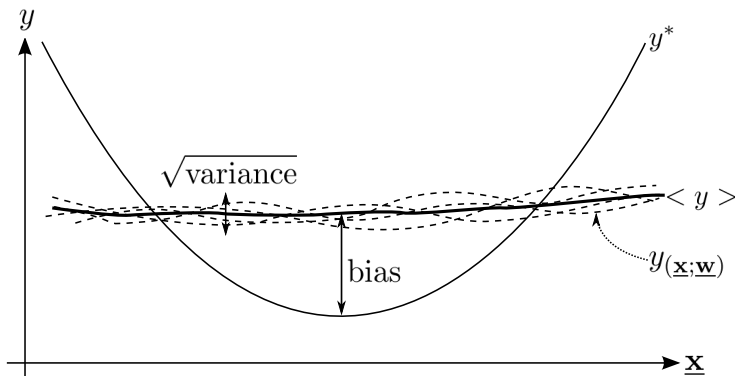
Example: squared error cost

$$\left\langle \left( \underbrace{y(\underline{\mathbf{x}}; \underline{\mathbf{w}})}_y - \underbrace{y^*_{(\underline{\mathbf{x}})}}_{y^*} \right)^2 \right\rangle_{\text{all datasets}}$$

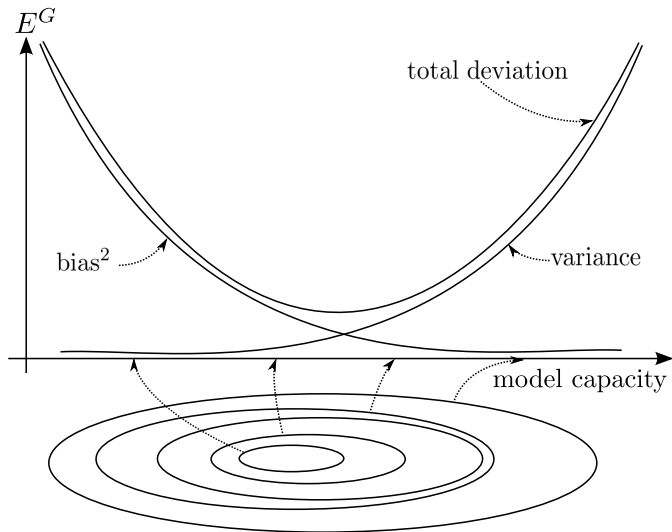
with

$$\begin{aligned} \langle (y - y^*)^2 \rangle &= \langle (y - \langle y \rangle + \langle y \rangle - y^*)^2 \rangle \\ &= \underbrace{(y^* - \langle y \rangle)^2}_{\text{bias}^2} + \underbrace{\langle (y - \langle y \rangle)^2 \rangle}_{\text{variance}} + 2 \underbrace{\langle y - \langle y \rangle \rangle (\langle y \rangle - y^*)}_{=0} \end{aligned}$$

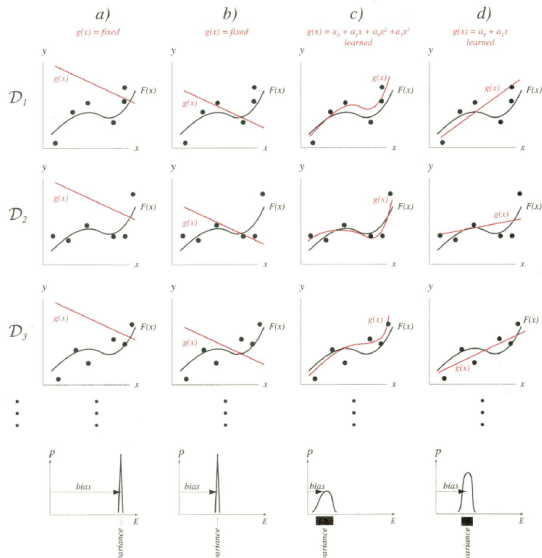
# Bias and variance



# Bias and variance



# Regression example



from Duda &amp; Hart 467

# Bias, variance and generalization error

Relation to generalization performance (for previous example)

$$E^G = \langle e^G \rangle_{P(\underline{\mathbf{x}})}, \text{ where } e^G = \langle (y(\underline{\mathbf{x}}; \underline{\mathbf{w}}) - y_T)^2 \rangle_{P(y_T | \underline{\mathbf{x}})}$$

# Bias, variance and generalization error

Relation to generalization performance (for previous example)

$$E^G = \langle e^G \rangle_{P(\underline{\mathbf{x}})}, \text{ where } e^G = \langle (y(\underline{\mathbf{x}}; \underline{\mathbf{w}}) - y_T)^2 \rangle_{P(y_T | \underline{\mathbf{x}})}$$

$$\begin{aligned} e^G &= \langle (y - y_T)^2 \rangle \\ &= \langle (y - y^* + y^* - y_T)^2 \rangle \\ &= (y - y^*)^2 + \langle (y^* - y_T)^2 \rangle + 2(y - y^*) \underbrace{\langle (y^* - y_T) \rangle}_{=0} \\ &= (y - y^*)^2 + \sigma_\eta^2 \end{aligned}$$

$$\langle e^G \rangle_{\text{ensemble}} = \underbrace{\langle (y - y^*)^2 \rangle}_{\text{bias}^2 + \text{variance}}_{\text{ensemble}} + \sigma_\eta^2$$

$\Rightarrow$  bias-variance trade-off applies to wide range of inductive learning problems

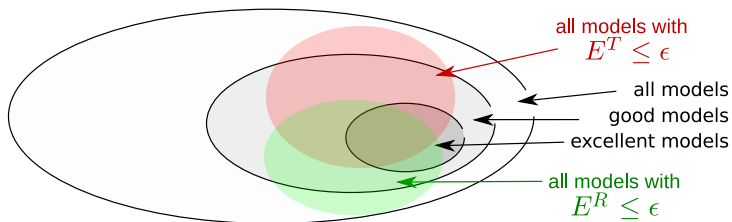
## 1.4.6 Regularization

# Regularization

## Risk function

$$R_{[\mathbf{w}]} = \underbrace{E_{[\mathbf{w}]}^T}_{\text{training error}} + \underbrace{\lambda E_{[\mathbf{w}]}^R}_{\text{regularization term}} \stackrel{!}{=} \min$$

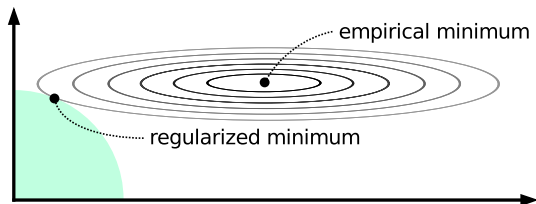
- $E^R$  : prior knowledge of solution
- $\lambda$  : regularization parameter





# Regularization example: weight decay

$$E_{[\mathbf{w}]}^R = \frac{1}{2} \sum_{i,j,v,v'} (w_{ij}^{vv'})^2$$

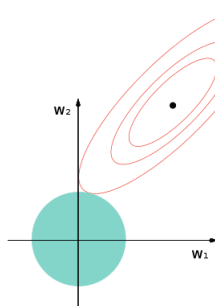
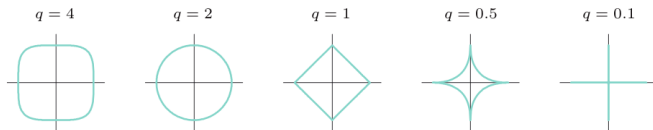


Minimization of  $R$  through gradient descent

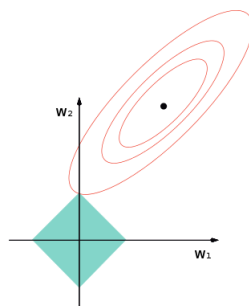
$$\Delta w_{ij}^{vv'} \sim -\frac{\partial R}{\partial w_{ij}^{vv'}} = -\underbrace{\frac{\partial E^T}{\partial w_{ij}^{vv'}}}_{\text{e.g. via backprop}} - \underbrace{\lambda w_{ij}^{vv'}}_{\text{decay term}}$$

# Other forms of regularization

- general form of regularization:  $E^R = \sum_{ijv'v} |w_{ij}^{v'v}|^q$



weight decay with  $q = 2$



lasso with  $q = 1$

## Regularization example: symmetries

Odd vs. even function

$$E_{[\mathbf{w}]}^R = \frac{1}{2p} \sum_{\alpha=1}^p \left( y_{(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}})} \pm y_{(-\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}})} \right)^2$$

## Regularization example: symmetries

Odd vs. even function

$$E_{[\mathbf{w}]}^R = \frac{1}{2p} \sum_{\alpha=1}^p \left( y_{(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}})} \pm y_{(-\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}})} \right)^2$$

Invariance under translation:

$$E_{[\mathbf{w}]}^R = \frac{1}{2p} \sum_{\alpha=1}^p \left( y_{(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}})} - y_{(\underline{\mathbf{x}}^{(\alpha)} - \underline{\mathbf{t}}; \underline{\mathbf{w}})} \right)^2$$

## Regularization example: symmetries

Odd vs. even function

$$E_{[\mathbf{w}]}^R = \frac{1}{2p} \sum_{\alpha=1}^p \left( y_{(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}})} \pm y_{(-\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}})} \right)^2$$

Invariance under translation:

$$E_{[\mathbf{w}]}^R = \frac{1}{2p} \sum_{\alpha=1}^p \left( y_{(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}})} - y_{(\underline{\mathbf{x}}^{(\alpha)} - \underline{\mathbf{t}}; \underline{\mathbf{w}})} \right)^2$$

Monotony:

$$E_{[\mathbf{w}]}^R = \frac{1}{n_p} \sum_{\mathbf{x}^{(\alpha)} > \mathbf{x}^{(\beta)}} \begin{cases} \left( y_{(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}})} - y_{(\underline{\mathbf{x}}^{(\beta)}; \underline{\mathbf{w}})} \right)^2 & , \text{ if } y_{(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}})} < y_{(\underline{\mathbf{x}}^{(\beta)}; \underline{\mathbf{w}})} \\ 0 & , \text{ else} \end{cases}$$

# Choice of regularization parameter

## Testset-Method

- 1 perform model selection for different values of  $\lambda$   
(on training data)

train	test	validation
-------	------	------------

# Choice of regularization parameter

## Testset-Method

- ① perform model selection for different values of  $\lambda$   
(on training data)
- ② select value of  $\lambda$ , which provides best prediction results  
(on test data)

train	test	validation
-------	------	------------

# Choice of regularization parameter

## Testset-Method

- ① perform model selection for different values of  $\lambda$   
(on training data)
- ② select value of  $\lambda$ , which provides best prediction results  
(on test data)
- ③ estimate generalization performance of selected  $\lambda$   
(on validation data)

train	test	validation
-------	------	------------



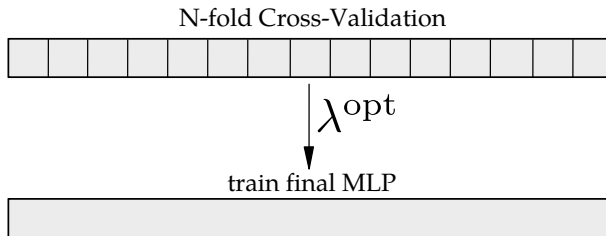
# Choice of regularization parameter

## $n$ -fold cross-validation

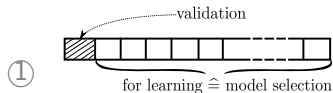
```
for  $\lambda = \lambda_1$  TO  $\lambda_n$  DO do  
  | perform  $n$ -fold cross-validation on data with regularization  $\lambda$   
end
```

pick optimal  $\lambda^{\text{opt}}$  with minimum  $\hat{E}^G$

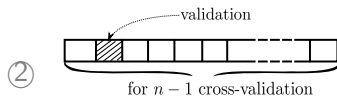
**final model:** train network on all data with  $\lambda^{\text{opt}}$



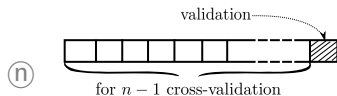
# Validation



- do  $n - 1$  cross-validation for all values of  $\lambda$
- train with best  $\lambda$
- validate with learned model



- do  $n - 1$  cross-validation for all values of  $\lambda$
- train with best  $\lambda$
- validate with learned model



- do  $n - 1$  cross-validation for all values of  $\lambda$
- train with best  $\lambda$
- validate with learned model

■  $\hat{E}^G \hat{=}$  average over all  $n$  validation errors

- Never use test data for model selection (including hyper-parameter search).
- Always embed the whole selection procedure (including hyper-parameter search) within an  $n$ -fold cross-validation run.

## 1.4.7 Classification Problems (Multi-Class)

# Overview inductive learning

data representation



model class



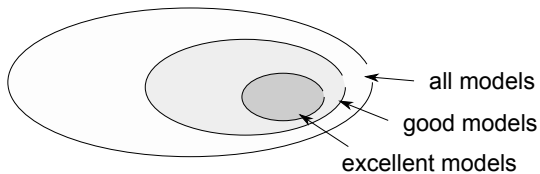
performance measure



optimization



validation



# Data representation

## Prediction of class labels

observations:  $\left\{ \left( \underline{\mathbf{x}}^{(\alpha)}, y_T^{(\alpha)} \right) \right\}, \quad \alpha \in \{1, \dots, p\}$

$c$  classes  $C_k, \quad k \in \{1, \dots, c\} \quad \Rightarrow \quad y_T^{(\alpha)} \in \{C_1, \dots, C_c\}$

# Data representation

## Prediction of class labels

observations:  $\left\{ \left( \underline{\mathbf{x}}^{(\alpha)}, y_T^{(\alpha)} \right) \right\}, \quad \alpha \in \{1, \dots, p\}$

$c$  classes  $C_k, \quad k \in \{1, \dots, c\} \quad \Rightarrow \quad y_T^{(\alpha)} \in \{C_1, \dots, C_c\}$

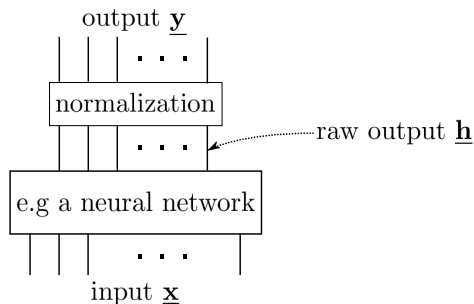
## Prediction of class probabilities

### 1-out-of-c-code

$$y_{Tk}^{(\alpha)} = \begin{cases} 0, & \underline{\mathbf{x}}^{(\alpha)} \notin C_k \\ 1, & \underline{\mathbf{x}}^{(\alpha)} \in C_k \end{cases} \quad \Rightarrow \text{binary vector } \underline{\mathbf{y}}_T^{(\alpha)}, \text{ one non-zero element}$$

Limiting case of probabilities: true labels are known.

# Model class



probabilistic interpretation  
of network output:

$$y_k(\underline{x}; \underline{w}) := P(C_k | \underline{x}; \underline{w})$$

$$0 \leq y_k(\underline{x}; \underline{w}) \leq 1$$

$$\sum_{k=1}^c y_k(\underline{x}; \underline{w}) = 1$$

## Softmax normalization

$$y_k(\underline{x}; \underline{w}) = \frac{\exp(h_k(\underline{x}; \underline{w}))}{\sum_l \exp(h_l(\underline{x}; \underline{w}))}$$



## Performance measure

True probability:  $P_{(C_k|\underline{\mathbf{x}})}$   $\leftrightarrow$  Model prediction:  $P_{(C_k|\underline{\mathbf{x}};\underline{\mathbf{w}})}$

### Kullback-Leibler-Divergence $D_{KL}$

$$D_{KL} = \sum_{k=1}^c \int d\underline{\mathbf{x}} P_{(\underline{\mathbf{x}})} P_{(C_k|\underline{\mathbf{x}})} \ln \left( \frac{P_{(C_k|\underline{\mathbf{x}})}}{P_{(C_k|\underline{\mathbf{x}};\underline{\mathbf{w}})}} \cdot \frac{P_{(\underline{\mathbf{x}})}}{P_{(\underline{\mathbf{x}})}} \right)$$

- distance measure between probability distributions and densities
- non-negative:  $D_{KL} = 0$  iff  $P_{(C|\underline{\mathbf{x}})} \equiv P_{(C|\underline{\mathbf{x}};\underline{\mathbf{w}})}$  (distributions are equal)
- asymmetric:  $D_{KL}(p||q)$  does not generally equal  $D_{KL}(q||p)$

## Performance measure

True probability:  $P_{(C_k|\underline{\mathbf{x}})}$   $\leftrightarrow$  Model prediction:  $P_{(C_k|\underline{\mathbf{x}};\underline{\mathbf{w}})}$

### Kullback-Leibler-Divergence $D_{KL}$

$$D_{KL} = \sum_{k=1}^c \int d\underline{\mathbf{x}} P(\underline{\mathbf{x}}) P_{(C_k|\underline{\mathbf{x}})} \ln \left( \frac{P_{(C_k|\underline{\mathbf{x}})}}{P_{(C_k|\underline{\mathbf{x}};\underline{\mathbf{w}})}} \cdot \frac{P(\underline{\mathbf{x}})}{P(\underline{\mathbf{x}})} \right)$$

- distance measure between probability distributions and densities
- non-negative:  $D_{KL} = 0$  iff  $P_{(C|\underline{\mathbf{x}})} \equiv P_{(C|\underline{\mathbf{x}};\underline{\mathbf{w}})}$  (distributions are equal)
- asymmetric:  $D_{KL}(p||q)$  does not generally equal  $D_{KL}(q||p)$

$$D_{KL} = - \sum_{k=1}^c \int d\underline{\mathbf{x}} P(\underline{\mathbf{x}}) P_{(C_k|\underline{\mathbf{x}})} \ln P_{(C_k|\underline{\mathbf{x}};\underline{\mathbf{w}})} + \underbrace{\sum_{k=1}^c \int d\underline{\mathbf{x}} P(\underline{\mathbf{x}}) P_{(C_k|\underline{\mathbf{x}})} \ln P_{(C_k|\underline{\mathbf{x}})}}_{\text{independent of model parameters}}$$

# Cross entropy

$$E^G \quad \equiv \quad - \sum_{k=1}^c \int d\underline{\mathbf{x}} \underbrace{P_{(\underline{\mathbf{x}})} P_{(C_k|\underline{\mathbf{x}})}}_{\text{unknown!}} \ln P_{(C_k|\underline{\mathbf{x}};\underline{\mathbf{w}})}$$

# Cross entropy

$$E^G \equiv - \sum_{k=1}^c \int d\underline{\mathbf{x}} \underbrace{P_{(\underline{\mathbf{x}})} P_{(C_k|\underline{\mathbf{x}})}}_{\text{unknown!}} \ln P_{(C_k|\underline{\mathbf{x}};\underline{\mathbf{w}})}$$

## ■ multinomial distribution

$$P_{(\underline{\mathbf{y}}_T|\underline{\mathbf{x}},\underline{\mathbf{w}})} = \prod_{k=1}^c (y_{k(\underline{\mathbf{x}};\underline{\mathbf{w}})})^{y_{T^k}},$$

$$y_{k(\underline{\mathbf{x}};\underline{\mathbf{w}})} \geq 0, \quad \forall k, \quad \sum_{k=1}^c y_{k(\underline{\mathbf{x}};\underline{\mathbf{w}})} = 1, \quad \forall \underline{\mathbf{x}}, \forall \underline{\mathbf{w}}$$

# Cross entropy

$$E^G \equiv - \sum_{k=1}^c \int d\underline{\mathbf{x}} \underbrace{P(\underline{\mathbf{x}}) P(C_k|\underline{\mathbf{x}})}_{\text{unknown!}} \ln P(C_k|\underline{\mathbf{x}}; \underline{\mathbf{w}})$$

## ■ multinomial distribution

$$P(\underline{\mathbf{y}}_T | \underline{\mathbf{x}}, \underline{\mathbf{w}}) = \prod_{k=1}^c (y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})})^{y_{T^k}},$$

$$y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})} \geq 0, \quad \forall k, \quad \sum_{k=1}^c y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})} = 1, \quad \forall \underline{\mathbf{x}}, \forall \underline{\mathbf{w}}$$

## ■ mathematical expectation $\rightarrow$ empirical average over training set:

$$E^T = \frac{1}{p} \sum_{\alpha=1}^p e^{(\alpha)} = -\frac{1}{p} \sum_{\alpha=1}^p \sum_{k=1}^c y_{T^k}^{(\alpha)} \ln (y_{k(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}})})$$

# Optimization via gradient descent (on-line)

$$\frac{\partial e^{(\alpha)}}{\partial \underline{\mathbf{w}}} = -\frac{\partial}{\partial \underline{\mathbf{w}}} \sum_{k=1}^c y_{T^k}^{(\alpha)} \ln(y_{k(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}})}) \quad \text{with} \quad y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})} = \frac{\exp h_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\sum_l \exp h_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}$$

= see blackboard...

# Optimization via gradient descent (on-line)

$$\begin{aligned}
 \frac{\partial e^{(\alpha)}}{\partial \underline{\mathbf{w}}} &= -\frac{\partial}{\partial \underline{\mathbf{w}}} \sum_{k=1}^c y_{Tk}^{(\alpha)} \ln(y_{k(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}})}) && \text{with} \quad y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})} = \frac{\exp h_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\sum_l \exp h_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})}} \\
 &= \sum_{k=1}^c \left( y_{k(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}})} - y_{Tk}^{(\alpha)} \right) \underbrace{\frac{\partial h_{k(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}}}_{\leadsto \text{backprop}}
 \end{aligned}$$

# Validation

- test-set method or n-fold cross-validation
- BUT: using the cross-entropy based cost

$$\hat{E}^G = \frac{1}{q} \sum_{\beta=1}^q e^{(\beta)} = -\frac{1}{q} \sum_{\beta=1}^q \sum_{k=1}^c y_{Tk}^{(\beta)} \ln (y_{k(\underline{\mathbf{x}}^{(\beta)}; \underline{\mathbf{w}})})$$



# Prediction of class labels

## ■ Decision costs

- $\$_{ij}$  : cost for choosing  $C_i$  when  $\underline{x}$  is of class  $C_j$

	patient is sick	patient is not sick
prediction: sick	buy medicine (\$20)	adverse effects (\$100)
prediction: not sick	sick leave (\$1500)	everything is fine (\$0)

- Choose  $C_k$  with minimal predicted costs,

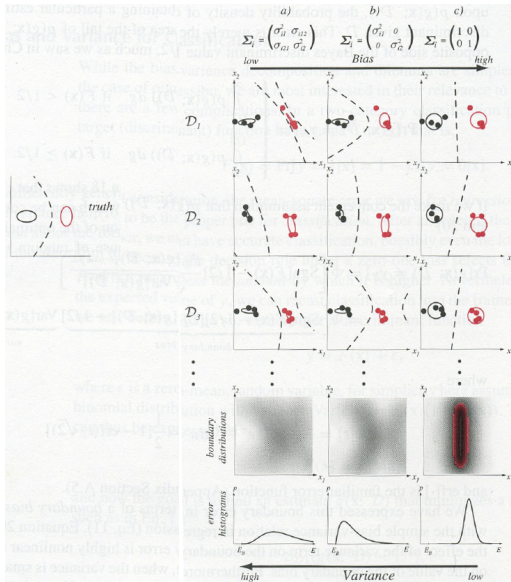
$$\text{i.e. } k = \underset{i}{\operatorname{argmin}} \sum_{j=1}^c \$_{ij} y_j(\underline{x}; \underline{w})$$

# End of Section 1.4

the following slides contain

## OPTIONAL MATERIAL

# Classification example



- bias-variance trade-off also applies to classification ("boundary error")
- BUT – for classification: not additive & variance dominates

# Gradient of softmax normalization

## Network computation of probabilities

$$y_{k(\underline{\mathbf{x}};\underline{\mathbf{w}})} = \frac{\exp(h_{k(\underline{\mathbf{x}};\underline{\mathbf{w}})})}{\sum_l \exp(h_{l(\underline{\mathbf{x}};\underline{\mathbf{w}})})} \quad (\text{softmax function})$$

# Gradient of softmax normalization

## Network computation of probabilities

$$y_{k(\underline{\mathbf{x}};\underline{\mathbf{w}})} = \frac{\exp(h_{k(\underline{\mathbf{x}};\underline{\mathbf{w}})})}{\sum_l \exp(h_{l(\underline{\mathbf{x}};\underline{\mathbf{w}})})} \quad (\text{softmax function})$$

$$\frac{\partial y_{k(\underline{\mathbf{x}};\underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}} = \frac{\exp(h_{k(\underline{\mathbf{x}};\underline{\mathbf{w}})}) \frac{\partial h_{k(\underline{\mathbf{x}};\underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}} \sum_{l=1}^c \exp(h_{l(\underline{\mathbf{x}};\underline{\mathbf{w}})}) - \exp(h_{k(\underline{\mathbf{x}};\underline{\mathbf{w}})}) \frac{\partial}{\partial \underline{\mathbf{w}}} \sum_{l=1}^c \exp(h_{l(\underline{\mathbf{x}};\underline{\mathbf{w}})})}{\left( \sum_{l=1}^c \exp(h_{l(\underline{\mathbf{x}};\underline{\mathbf{w}})}) \right)^2}$$

# Gradient of softmax normalization

## Network computation of probabilities

$$y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})} = \frac{\exp(h_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})})}{\sum_l \exp(h_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})})} \quad (\text{softmax function})$$

$$\begin{aligned} \frac{\partial y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}} &= \frac{\exp(h_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}) \frac{\partial h_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}} \sum_{l=1}^c \exp(h_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})}) - \exp(h_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}) \frac{\partial}{\partial \underline{\mathbf{w}}} \sum_{l=1}^c \exp(h_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})})}{\left( \sum_{l=1}^c \exp(h_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})}) \right)^2} \\ &= y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})} \frac{\partial h_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}} - y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})} \sum_{l=1}^c y_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})} \frac{\partial h_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}} \end{aligned}$$

# Gradient of Cross Entropy

$$\frac{\partial e^{(\alpha)}}{\partial \underline{\mathbf{w}}} = - \frac{\partial}{\partial \underline{\mathbf{w}}} \sum_{k=1}^c y_{T^k}^{(\alpha)} \ln (y_{k(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}})})$$

# Gradient of Cross Entropy

$$\frac{\partial e^{(\alpha)}}{\partial \underline{\mathbf{w}}} = -\frac{\partial}{\partial \underline{\mathbf{w}}} \sum_{k=1}^c y_{Tk}^{(\alpha)} \ln(y_{k(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}})}) \quad \text{with} \quad y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})} = \frac{\exp h_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\sum_l \exp h_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}$$



# Gradient of Cross Entropy

$$\begin{aligned}\frac{\partial e^{(\alpha)}}{\partial \underline{\mathbf{w}}} &= -\frac{\partial}{\partial \underline{\mathbf{w}}} \sum_{k=1}^c y_{Tk}^{(\alpha)} \ln(y_{k(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}})}) \\ &= -\sum_{k=1}^c \frac{y_{Tk}}{y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}} \cdot \frac{\partial y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}}\end{aligned}$$

with  $y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})} = \frac{\exp h_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\sum_l \exp h_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}$

see derivation on Slide 3

# Gradient of Cross Entropy

$$\begin{aligned}
 \frac{\partial e^{(\alpha)}}{\partial \underline{\mathbf{w}}} &= -\frac{\partial}{\partial \underline{\mathbf{w}}} \sum_{k=1}^c y_{Tk}^{(\alpha)} \ln(y_{k(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}})}) && \text{with } y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})} = \frac{\exp h_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\sum_l \exp h_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})}} \\
 &= -\sum_{k=1}^c \frac{y_{Tk}}{y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}} \cdot \frac{\partial y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}} && \text{see derivation on Slide 3} \\
 &= -\sum_{k=1}^c \frac{y_{Tk}}{y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}} \left( y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})} \frac{\partial h_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}} - y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})} \sum_{l=1}^c y_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})} \frac{\partial h_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}} \right)
 \end{aligned}$$

# Gradient of Cross Entropy

$$\begin{aligned}
 \frac{\partial e^{(\alpha)}}{\partial \underline{\mathbf{w}}} &= -\frac{\partial}{\partial \underline{\mathbf{w}}} \sum_{k=1}^c y_{Tk}^{(\alpha)} \ln(y_{k(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}})}) && \text{with } y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})} = \frac{\exp h_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\sum_l \exp h_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})}} \\
 &= -\sum_{k=1}^c \frac{y_{Tk}}{y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}} \cdot \frac{\partial y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}} && \text{see derivation on Slide 3} \\
 &= -\sum_{k=1}^c \frac{y_{Tk}}{y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}} \left( y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})} \frac{\partial h_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}} - y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})} \sum_{l=1}^c y_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})} \frac{\partial h_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}} \right) \\
 &= -\sum_{k=1}^c y_{Tk} \frac{\partial h_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}} + \underbrace{\left( \sum_{k=1}^c y_{Tk} \right)}_{=1} \left( \sum_{l=1}^c y_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})} \frac{\partial h_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}} \right)
 \end{aligned}$$

# Gradient of Cross Entropy

$$\begin{aligned}
 \frac{\partial e^{(\alpha)}}{\partial \underline{\mathbf{w}}} &= -\frac{\partial}{\partial \underline{\mathbf{w}}} \sum_{k=1}^c y_{Tk}^{(\alpha)} \ln(y_{k(\underline{\mathbf{x}}^{(\alpha)}; \underline{\mathbf{w}})}) && \text{with } y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})} = \frac{\exp h_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\sum_l \exp h_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})}} \\
 &= -\sum_{k=1}^c \frac{y_{Tk}}{y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}} \cdot \frac{\partial y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}} && \text{see derivation on Slide 3} \\
 &= -\sum_{k=1}^c \frac{y_{Tk}}{y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}} \left( y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})} \frac{\partial h_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}} - y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})} \sum_{l=1}^c y_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})} \frac{\partial h_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}} \right) \\
 &= -\sum_{k=1}^c y_{Tk} \frac{\partial h_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}} + \underbrace{\left( \sum_{k=1}^c y_{Tk} \right)}_{=1} \left( \sum_{l=1}^c y_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})} \frac{\partial h_{l(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}} \right) \\
 &= \sum_{k=1}^c \underbrace{\left( y_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})} - y_{Tk} \right)}_{\leadsto \text{backprop}} \frac{\partial h_{k(\underline{\mathbf{x}}; \underline{\mathbf{w}})}}{\partial \underline{\mathbf{w}}}
 \end{aligned}$$