

Machine Intelligence 1

1.3 Multi-layer Perceptrons

Prof. Dr. Klaus Obermayer

Fachgebiet Neuronale Informationsverarbeitung (NI)

WS 2017/2018

1.3.1 Classes of Neural Networks

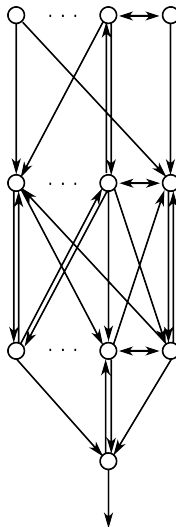
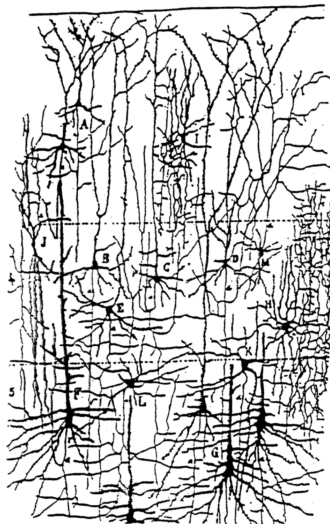
Graphs representing neural networks



Design principles inspired by biology

- architecture:
 - simple elements
 - massively parallel systems
 - low precision and robustness
 - distributed representation of information
 - no separation between data and program
- inductive learning:
 - data driven, adaptive systems
 - learning and self-organization vs. deduction and programming
 - often seen as a plus: biologically inspired learning rules

Graphs representing neural networks



neural network



directed graph

(connectionist) neuron



node of the graph

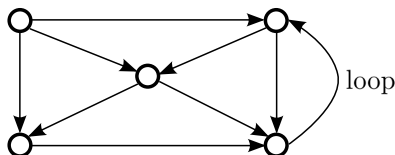
neural connection



weighted edge

Recurrent Neural Networks (RNNs)

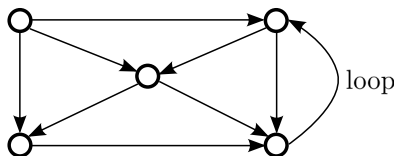
Recurrent networks $\hat{=}$ directed graphs containing cycles



- dynamical systems
- spatio-temporal pattern analysis
- sequence processing
- associative memory and pattern completion

Recurrent Neural Networks (RNNs)

Recurrent networks \cong directed graphs containing cycles



- dynamical systems
- spatio-temporal pattern analysis
- sequence processing
- associative memory and pattern completion

Examples:

Hopfield networks (Hopfield, 1982)

Boltzmann machines (Ackley, Hinton & Sejnowski, 1985)

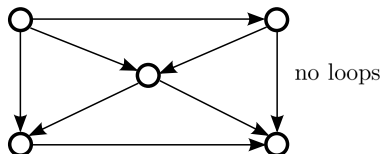
Infinite Impulse Response (IIR) networks (Crochiere & Oppenheim, 1975)

Long Short-Term Memory networks (LSTM, Hochreiter & Schmidhuber, 1997)

Deep Recurrent Neural Networks (Pascanu, Gulcehre, Cyho & Bengio, 2014)

Feedforward Neural Networks (FFN)

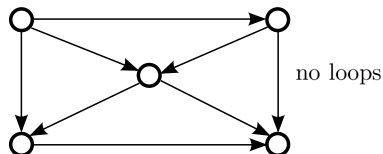
Feedforward Networks \equiv directed acyclic graphs (DAGs)



- association between variables
- prediction of attributes

Feedforward Neural Networks (FFN)

Feedforward Networks \cong directed acyclic graphs (DAGs)



- association between variables
- prediction of attributes

Examples:

Multilayer-**P**erceptron (MLP, Rumelhart, Hinton & Williams, 1986)

Radial **B**asis **F**unction network (RBF, Broomhead & Lowe, 1988)

Support **V**ector **M**achines (SVM, Cortes & Vapnik, 1995)

Deep **B**elief **N**etworks (DBN, Hinton, Osidero & Teh, 2006; Mnih et al., 2016)

Prediction of attributes

real-valued attributes: regression problems

$$\begin{array}{ll} f : \mathbb{R}^N \rightarrow \mathbb{R} & \text{(or subsets) one target value} \\ f : \mathbb{R}^N \rightarrow \mathbb{R}^M & \text{(or subsets) multivariate attributes} \end{array}$$

Prediction of attributes

real-valued attributes: regression problems

$$\begin{array}{ll} f : \mathbb{R}^N \rightarrow \mathbb{R} & \text{(or subsets) one target value} \\ f : \mathbb{R}^N \rightarrow \mathbb{R}^M & \text{(or subsets) multivariate attributes} \end{array}$$

ordinal attributes: classification problems

$$\begin{array}{ll} f : \mathbb{R}^N \rightarrow \mathcal{S} & \text{where } \mathcal{S} \text{ is a set of attributes } \{a_1, a_2, \dots, a_M\} \\ f : \mathbb{R}^N \rightarrow \{-1, +1\} & \text{special case: two class problems} \end{array}$$

Prediction of attributes

real-valued attributes: regression problems

$$\begin{array}{ll} f : \mathbb{R}^N \rightarrow \mathbb{R} & \text{(or subsets) one target value} \\ f : \mathbb{R}^N \rightarrow \mathbb{R}^M & \text{(or subsets) multivariate attributes} \end{array}$$

ordinal attributes: classification problems

$$\begin{array}{ll} f : \mathbb{R}^N \rightarrow \mathcal{S} & \text{where } \mathcal{S} \text{ is a set of attributes } \{a_1, a_2, \dots, a_M\} \\ f : \mathbb{R}^N \rightarrow \{-1, +1\} & \text{special case: two class problems} \end{array}$$

predicting probabilities

$$\begin{array}{ll} f_k : \mathbb{R}^N \rightarrow [0, 1] \subset \mathbb{R} & \text{probability of attribute } a_k \text{ from } \mathcal{S} \\ \text{s.t. } \sum_{k=1}^M f_k(\cdot) = 1 & \text{constrains all } f_k(\cdot) \text{ to be probabilities} \end{array}$$

Prediction of attributes

real-valued attributes: regression problems

$$\begin{array}{ll} f : \mathbb{R}^N \rightarrow \mathbb{R} & \text{(or subsets) one target value} \\ f : \mathbb{R}^N \rightarrow \mathbb{R}^M & \text{(or subsets) multivariate attributes} \end{array}$$

ordinal attributes: classification problems

$$\begin{array}{ll} f : \mathbb{R}^N \rightarrow \mathcal{S} & \text{where } \mathcal{S} \text{ is a set of attributes } \{a_1, a_2, \dots, a_M\} \\ f : \mathbb{R}^N \rightarrow \{-1, +1\} & \text{special case: two class problems} \end{array}$$

predicting probabilities

$$\begin{array}{ll} f_k : \mathbb{R}^N \rightarrow [0, 1] \subset \mathbb{R} & \text{probability of attribute } a_k \text{ from } \mathcal{S} \\ \text{s.t. } \sum_{k=1}^M f_k(\cdot) = 1 & \text{constrains all } f_k(\cdot) \text{ to be probabilities} \end{array}$$

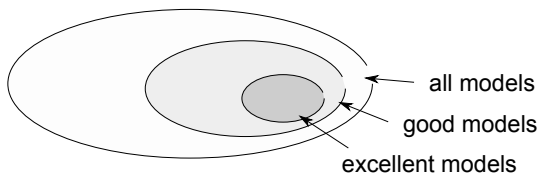
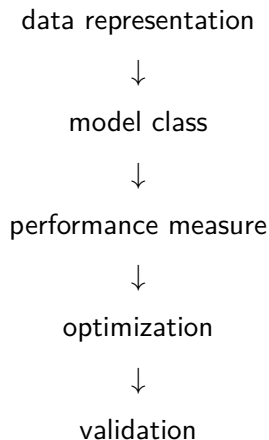
structured input / structured output

$$f : \mathcal{X} \rightarrow \mathcal{Y} \quad \text{where } \mathcal{X} \text{ and } \mathcal{Y} \text{ are sets of structures, e.g. graphs or sentences}$$

see Baklir, Hofmann, Schölkopf, Smola, Taskar and Vishwanathan (Eds.): Predicting structured data, 2007

1.3.2 The Multi-layer-Perceptron for Regression

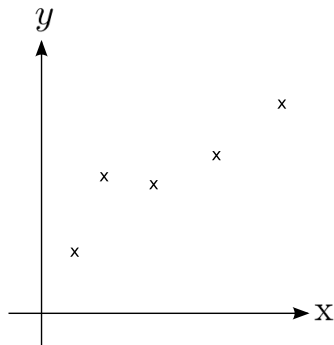
Overview inductive learning



data representation

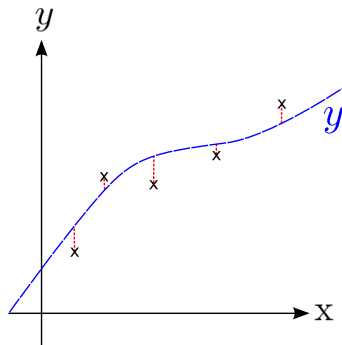
Our example: regression with one real-valued attribute

- input feature vectors: $\underline{\mathbf{x}} \in \mathbb{R}^N$
- output attributes: $y_T \in \mathbb{R}$
- training set: $\{\underline{\mathbf{x}}^{(\alpha)}, y_T^{(\alpha)}\}_{\alpha=1}^p$



Our example: regression with one real-valued attribute

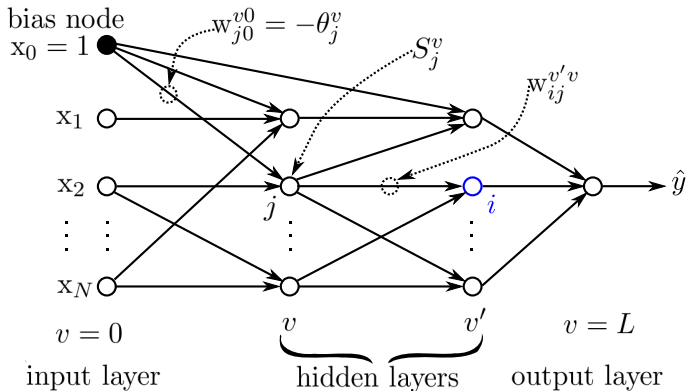
- input feature vectors: $\underline{\mathbf{x}} \in \mathbb{R}^N$
- output attributes: $y_T \in \mathbb{R}$
- training set: $\{\underline{\mathbf{x}}^{(\alpha)}, y_T^{(\alpha)}\}_{\alpha=1}^p$
- label function: $y_T = y(\underline{\mathbf{x}}) + \text{noise}$



model class

Multi-layer perceptrons (MLPs)

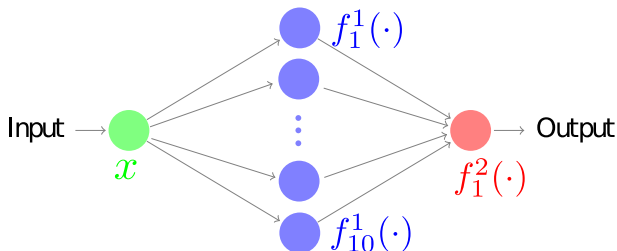
- layered FFN $\hat{y}(\cdot; \underline{\mathbf{w}})$ models label function $y(\cdot)$



- $S_i^{v'} = f_i^{v'}(h_i^{v'}) \stackrel{\text{e.g.}}{=} \tanh(h_i^{v'}), \quad h_i^{v'} = \sum_{j=0}^N w_{ij}^{v'v} S_j^v$

Model class: one example

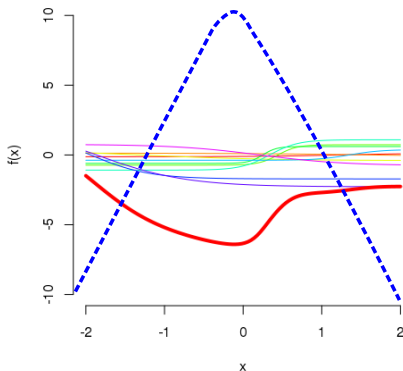
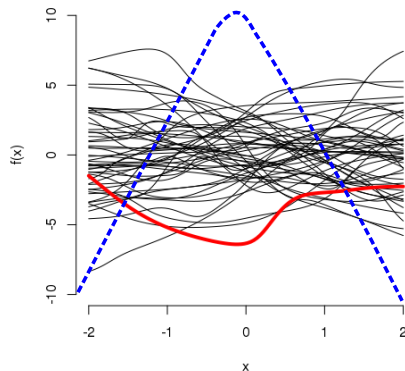
1 input, 1 output, 1 hidden layer with 10 units



$$\hat{y}(x) = f_1^2 \left(\sum_{i=1}^{10} w_{1i}^{21} f_i^1(w_{i1}^{10} x - w_{i0}^{10}) \right)$$

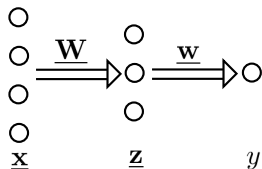
$$f_i^1(a) = \tanh(a) \text{ \& } f_1^2(a) = a$$

Model class: one example

Input-Output of 1 NN**50 randomly drawn NNs**

(parameters drawn from Gaussian distribution)

Linear transfer functions



$$\begin{aligned}
 y &= \underline{w}^\top \underline{z} = \underline{w}^\top \underline{W} \underline{x} \\
 &= (\underline{W}^\top \underline{w})^\top \underline{x} = \underline{\hat{w}}^\top \underline{x} \\
 &\hat{=} \text{connectionist neuron}
 \end{aligned}$$

MLPs are universal approximators

Funahashi (1989)

Let $y_{(\underline{x})}^*$ be a continuous, real valued function over a compact interval K and

$$\hat{y}_{(\underline{x})} = \sum_{i=1}^M w_i^{21} f\left(\sum_{j=1}^N w_{ij}^{10} x_j - \theta_i\right)$$

be a three-layered MLP with a non-constant, bounded, monotonously increasing and continuous function $f : \mathbb{R} \rightarrow \mathbb{R}$.

MLPs are universal approximators

Funahashi (1989)

Let $y^*_{(\underline{x})}$ be a continuous, real valued function over a compact interval K and

$$\hat{y}_{(\underline{x})} = \sum_{i=1}^M w_i^{21} f\left(\sum_{j=1}^N w_{ij}^{10} x_j - \theta_i\right)$$

be a three-layered MLP with a non-constant, bounded, monotonously increasing and continuous function $f : \mathbb{R} \rightarrow \mathbb{R}$.

Then there exists a set of parameters $M, N \in \mathbb{N}$ and $w_i^{21}, w_{ij}^{10}, \theta_i \in \mathbb{R}$ such that for every $\varepsilon > 0$:

$$\max_{\underline{x} \in K} \left| \hat{y}_{(\underline{x})} - y^*_{(\underline{x})} \right| \leq \varepsilon$$

Funahashi (1989) On the approximate realization of continuous mappings by neural networks. *Neur Netw*, 2:183–192

Hornik et al. (1989) Multilayer Feedforward Networks are Universal Approximators. *Neur Netw*, 2:359–366.

1.3.3 Performance Measures and Model Selection

Cost functions

$$\underbrace{\underline{\mathbf{x}} \in \mathbb{R}^N}_{\text{feature vector}} \longrightarrow \underbrace{y \in \mathbb{R}}_{\text{attribute}}$$

y_T : true value of attribute

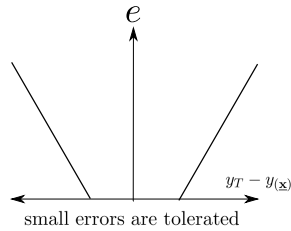
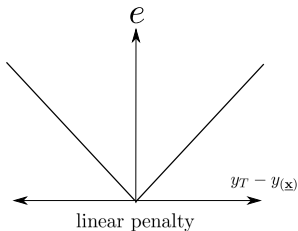
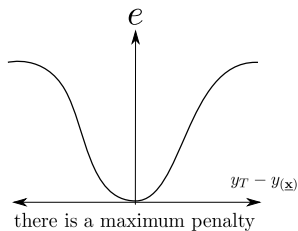
$y(\underline{\mathbf{x}})$: predicted value of attribute (e.g. by MLP)

Cost functions

$$\underbrace{\underline{\mathbf{x}} \in \mathbb{R}^N}_{\text{feature vector}} \longrightarrow \underbrace{y \in \mathbb{R}}_{\text{attribute}}$$

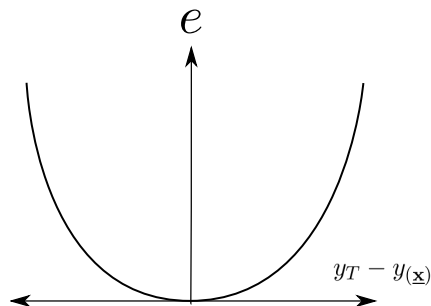
y_T : true value of attribute
 $y(\underline{\mathbf{x}})$: predicted value of attribute (e.g. by MLP)

individual cost $e(y_T, \underline{\mathbf{x}})$



several choices \Rightarrow predictor will depend on error measure!

Lecture example: quadratic error



$$e(y_T, \underline{x}) = \frac{1}{2} \left(y(\underline{x}) - y_T \right)^2$$

- Gaussian noise on attributes
- sensitive against “outliers”

Performance measure

Generalization error

$$E^G := \langle e \rangle_{y_T, \underline{x}} = \iint d\underline{x} dy_T P_{(y_T, \underline{x})} e_{(y_T, \underline{x})}$$

$P_{(y_T, \underline{x})}$: joint Probability Density Function (PDF) of observations

Performance measure

Generalization error

$$E^G := \langle e \rangle_{y_T, \underline{\mathbf{x}}} = \iint d\underline{\mathbf{x}} dy_T P_{(y_T, \underline{\mathbf{x}})} e_{(y_T, \underline{\mathbf{x}})}$$

$P_{(y_T, \underline{\mathbf{x}})}$: joint Probability Density Function (PDF) of observations

“good” predictor: low value of E^G

“bad” predictor: high value of E^G

$$E^G \stackrel{!}{=} \min$$

but: $P(y_T | \underline{\mathbf{x}})$ is not known

Principle of Empirical Risk Minimization (ERM)

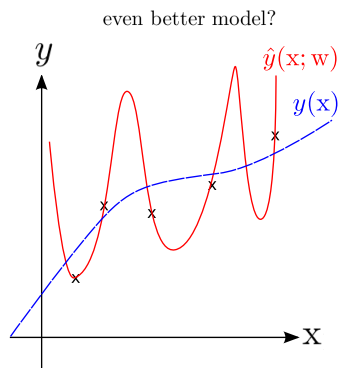
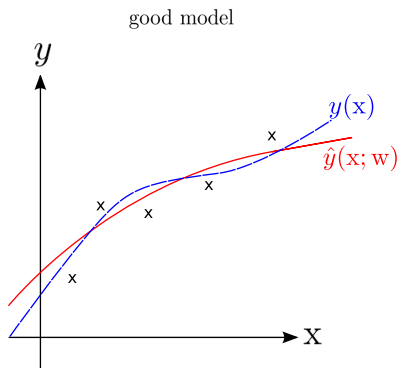
math. expectation (E^G) \rightarrow empirical average (E^T)

$$E^G := \int d\mathbf{x} dy_T P(y_T, \mathbf{x}) e(y_T, \mathbf{x}) \rightarrow E^T := \frac{1}{p} \sum_{\alpha=1}^p e^{(\alpha)}$$

generalization error \rightarrow training error

$$E^G \stackrel{!}{=} \min \rightarrow E^T \stackrel{!}{=} \min$$

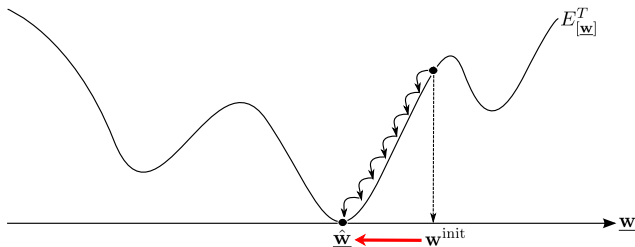
Consequences?



1.3.4 Optimization of Model Parameters: Gradient Descent

Gradient descent

training error E^T for the given training set: $E_{[\underline{\mathbf{w}}]}^T = \frac{1}{p} \sum_{\alpha=1}^p e_{[\underline{\mathbf{w}}]}^{(\alpha)}$



$$w_{ij}^{v'v}(t+1) = w_{ij}^{v'v}(t) - \underbrace{\hat{\eta}}_{\text{learning step}} \underbrace{\frac{\partial E_{[\underline{\mathbf{w}}]}^T}{\partial w_{ij}^{v'v}}}_{\text{gradient vector}}$$

Calculation of the gradient

$$\begin{aligned}
 \frac{\partial E_{\underline{\mathbf{w}}}^T}{\partial w_{ij}^{v'v}} &= \frac{1}{p} \sum_{\alpha=1}^p \underbrace{\frac{\partial e_{\underline{\mathbf{w}}}^{(\alpha)}}{\partial w_{ij}^{v'v}}}_{\text{individual cost}} = \frac{1}{p} \sum_{\alpha=1}^p \underbrace{\frac{\partial e_{\underline{\mathbf{w}}}^{(\alpha)}}{\partial y(\underline{\mathbf{x}}^{(\alpha)}, \underline{\mathbf{w}})}}_{\text{factor depending on cost function}} \cdot \underbrace{\frac{\partial y(\underline{\mathbf{x}}^{(\alpha)}, \underline{\mathbf{w}})}{\partial w_{ij}^{v'v}}}_{\text{factor depending on model class (e.g. MLP)}}
 \end{aligned}$$

Calculation of the error term

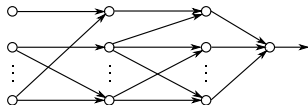
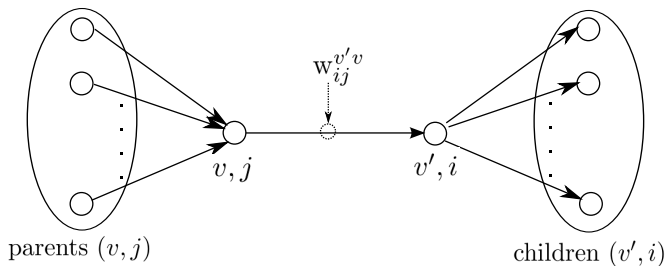
Quadratic Error

$$e(y_T, \underline{\mathbf{x}}) = \frac{1}{2} (y_T - y(\underline{\mathbf{x}}))^2 \quad \Rightarrow \quad \frac{\partial e^{(\alpha)}_{[\underline{\mathbf{w}}]}}{\partial y_{(\underline{\mathbf{x}}^{(\alpha)}, \underline{\mathbf{w}})}} = y_{(\underline{\mathbf{x}}^{(\alpha)}, \underline{\mathbf{w}})} - y_T$$

1.3.5 The Backpropagation Method

Gradients in neural networks

$$\frac{\partial y(\underline{\mathbf{x}}^{(\alpha)}, \underline{\mathbf{w}})}{\partial w_{ij}^{v'v}}$$



see blackboard for derivation

Summary of the backpropagation method

initialization of weights and thresholds

while stopping criterion not met **do**

gradient $_{ij}^{v'v} := 0, \quad \forall w_{ij}^{v'v}$

for $\alpha \in \{1, \dots, p\}$ **do**

$h_i^0 := x_i^\alpha, \quad \forall i$ // forward propagation

for $v' \in \{1, \dots, L\}$ **do**

$$h_i^{v'} := \sum_{(v', i) \in C(v, j)} w_{ij}^{v'v} f_j^v(h_j^v), \quad \forall i$$

end

$\delta_i^L := f'_i(h_i^L), \quad \forall i$ // backward propagation

for $v' \in \{L-1, \dots, 1\}$ **do**

$$\delta_i^{v'} := f_i^{v'}(h_i^{v'}) \sum_{(\beta, k) \in C(v', i)} \delta_k^\beta w_{ki}^{\beta v'}, \quad \forall i$$

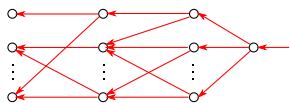
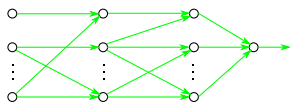
end

gradient $_{ij}^{v'v} := \text{gradient}_{ij}^{v'v} + \frac{\partial e^{(\alpha)}}{\partial y} \delta_i^{v'} f_j^v(h_j^v), \quad \forall w_{ij}^{v'v}$ // average

end

$w_{ij}^{v'v} := w_{ij}^{v'v} - \frac{\eta}{p} \text{gradient}_{ij}^{v'v}, \quad \forall w_{ij}^{v'v}$ // gradient descend step

end



Summary of the backpropagation method

initialization of weights and thresholds

while *stopping criterion not met* **do**

gradient $_{ij}^{v'v} := 0, \quad \forall w_{ij}^{v'v}$

for $\alpha \in \{1, \dots, p\}$ **do**

$h_i^0 := x_i^\alpha, \quad \forall i$ // forward propagation

for $v' \in \{1, \dots, L\}$ **do**

$h_i^{v'} := \sum_{(v', i) \in C(v, j)} w_{ij}^{v'v} f_j^v(h_j^v), \quad \forall i$

end

$\delta_i^L := f'_i(h_i^L), \quad \forall i$ // backward propagation

for $v' \in \{L-1, \dots, 1\}$ **do**

$\delta_i^{v'} := f'^{v'}_i(h_i^{v'}) \sum_{(\beta, k) \in C(v', i)} \delta_k^\beta w_{ki}^{\beta v'}, \quad \forall i$

end

gradient $_{ij}^{v'v} := \text{gradient}_{ij}^{v'v} + \frac{\partial e^{(\alpha)}}{\partial y} \delta_i^{v'} f_j^v(h_j^v), \quad \forall w_{ij}^{v'v}$ // average

end

$w_{ij}^{v'v} := w_{ij}^{v'v} - \frac{\hat{\eta}}{p} \text{gradient}_{ij}^{v'v}, \quad \forall w_{ij}^{v'v}$ // gradient descend step

end

computational and
memory complexity

$\mathcal{O}(n),$

n : number of weights

Initialization and stopping criteria

Initialization: random numbers, such that h_i^v approx. $\mathcal{O}(1)$

→ too **large**?

→ too **small**?

Initialization and stopping criteria

Initialization: random numbers, such that h_i^v approx. $\mathcal{O}(1)$

- too **large**: transfer function saturates and gradients become too small
- too **small**: neurons operate in the linear regime of f

Initialization and stopping criteria

Initialization: random numbers, such that h_i^v approx. $\mathcal{O}(1)$

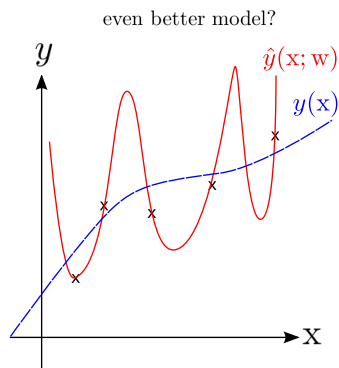
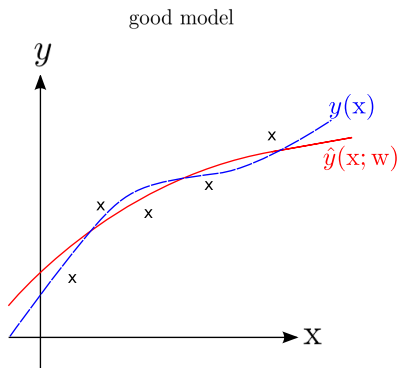
- too **large**: transfer function saturates and gradients become too small
- too **small**: neurons operate in the linear regime of f

Stopping criteria:

- fixed number of iterations
- fixed CPU-time
- E^T falls below a predefined value
- $\frac{\Delta E^T}{E^T}$ falls below a predefined value
- validation criterion fulfilled

1.3.6 Validation of Model Selection

Overfitting



Assessment of prediction quality

Test set method

$$\text{observations} \left\{ \begin{array}{l} \text{training data } \left\{ \left(\underline{\mathbf{x}}^{(\alpha)}, y_T^{(\alpha)} \right) \right\}, \alpha \in \{1, \dots, p\} \\ \rightarrow E^T \text{ selects model parameters} \end{array} \right.$$

$$\hat{E}^T = \frac{1}{p} \sum_{\alpha=1}^q e^{(\alpha)}$$

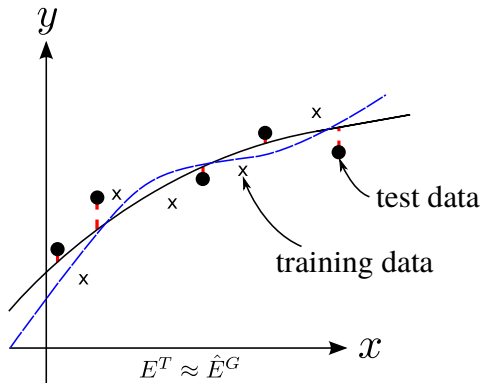
Assessment of prediction quality

Test set method

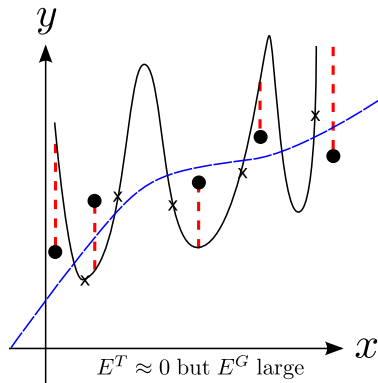
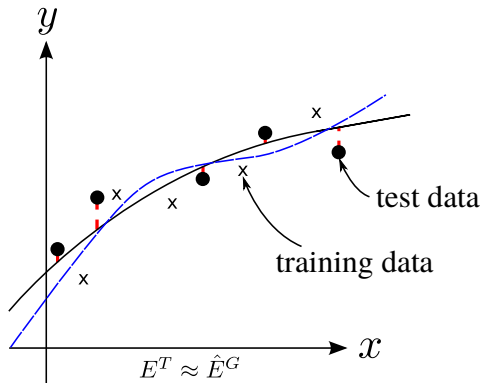
$$\text{observations} \left\{ \begin{array}{l} \text{training data } \left\{ \left(\underline{\mathbf{x}}^{(\alpha)}, y_T^{(\alpha)} \right) \right\}, \alpha \in \{1, \dots, p\} \\ \rightarrow E^T \text{ selects model parameters} \\ \text{test data } \left\{ \left(\underline{\mathbf{x}}^{(\beta)}, y_T^{(\beta)} \right) \right\}, \beta \in \{1, \dots, q\} \\ \rightarrow \hat{E}^G \text{ estimates generalization error} \end{array} \right.$$

$$\hat{E}^G = \frac{1}{q} \sum_{\beta=1}^q e^{(\beta)}$$

Test set method

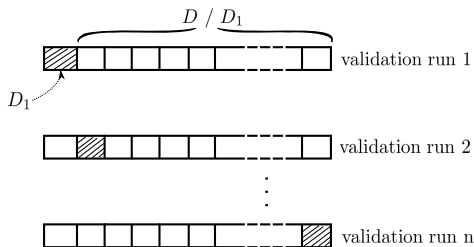


Test set method



Resampling methods: n -fold cross-validation

Observations $D \rightarrow n$ disjunct sets D_j , $\bigcup_{j=1}^n D_j = D$



Training of n networks on the training datasets D / D_j
 Estimation of E^G :

$$\hat{E}^G = \frac{1}{p} \sum_j \sum_{\alpha \in D_j} e^{(\alpha)}$$

Typical choice: $n \in \{5, \dots, 10\}$ $n = p$: leave-one-out cross-validation

Variance of the cross-validation estimator

Remarks

- n-fold cross-validation is only used for estimating E^G
- each of the n folds results in a different solution (set of parameters)
- All data are used for selecting the model parameters

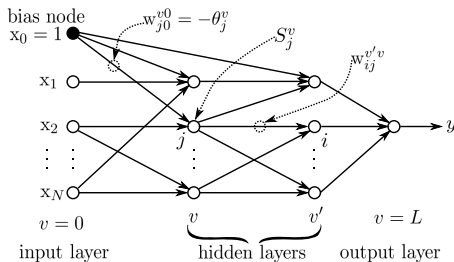
$$\text{var}(\hat{E}^G) = \frac{n-1}{n} \sum_{j=1}^n \left(\underbrace{\hat{E}_j^G}_{\text{test error on } D_j} - \hat{E}^G \right)^2$$

End of Section 1.3

the following slides contain

OPTIONAL MATERIAL

Multi-layer perceptrons: nomenclature



S_j^v activity of neuron (v, j) (layer, unit)

$S_0^v = x_0 = 1$ activity of **bias** neuron

$S_i^0 = x_i$ **input** to MLP, i^{th} component

$S_i^L = y$ **output** of MLP

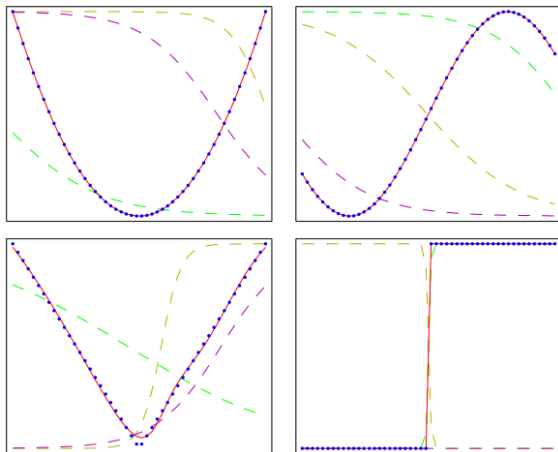
$w_{ij}^{v'v}$ connection **weight** between neurons (v, j) and (v', i)

$w_{j0}^{v0} = \theta_j^v$ connection weight between bias node and neuron (v, j)

$h_i^{v'}$ $= \sum_j w_{ij}^{v'v} S_j^{v-1}$ = total input of neuron (v', i)

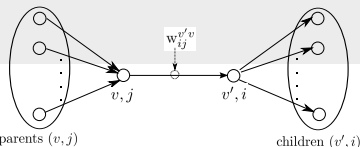
$f_i^{v'}$ transfer function of neuron (v', i)

Illustration: ERM with MLPs (Bishop 2009)



fitted MLPs with 1 hidden layer of 3 neurons

The Credit Assignment Problem



How do the weights $w_{ij}^{v'v}$ of hidden units $S_i^{v'}$ contribute to the error?

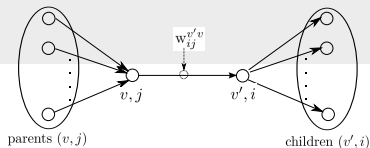
Solution: smart application of the chain rule

$$\frac{\partial y}{\partial w_{ij}^{v'v}} = \underbrace{\frac{\partial y}{\partial h_i^{v'}}}_{:= \delta_i^{v'} \text{ "local error" at neuron } (v', i)} \cdot \underbrace{\frac{\partial h_i^{v'}}{\partial w_{ij}^{v'v}}}_{= S_j^v \text{ activity of neuron } (v, j)}$$

Forward propagation: calculation of activities (parents \rightarrow children)

$$S_j^0 = x_j^{(\alpha)} \rightarrow S_i^{v'} = f\left(\sum_{(v', i) \in C(v, j)} w_{ij}^{v'v} S_j^v\right)$$

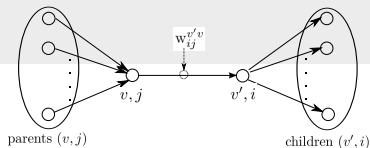
Backpropagation



Backpropagation: calculation of “local errors” (children \rightarrow parents)

$$\delta_i^L = \underbrace{f'(h_i^{v'})}_{=1 \text{ for identity}} \quad \text{and} \quad \delta_i^{v'} = \sum_{(v'', k) \in C(v', i)} \frac{\partial y}{\partial h_k^{v''}} \cdot \frac{\partial h_k^{v''}}{\partial h_i^{v'}}, \quad \forall v' \neq L$$

Backpropagation



Backpropagation: calculation of “local errors” (children \rightarrow parents)

$$\delta_i^L = \underbrace{f'(h_i^{v'})}_{=1 \text{ for identity}} \quad \text{and} \quad \delta_i^{v'} = \sum_{(v'', k) \in C(v', i)} \frac{\partial y}{\partial h_k^{v''}} \cdot \frac{\partial h_k^{v''}}{\partial h_i^{v'}}, \quad \forall v' \neq L$$

rewrite using $\frac{\partial y}{\partial h_k^{v''}} = \delta_k^{v''}$

$$\delta_i^{v'} = \sum_{(v'', k) \in C(v', i)} \delta_k^{v''} \cdot w_{ki}^{v''v'} f'(h_i^{v'}) = f'(h_i^{v'}) \sum_{(v'', k) \in C(v', i)} \delta_k^{v''} w_{ki}^{v''v'}$$

Computational complexity: $O(n)$, n : number of weights & thresholds