

Machine Learning 2

4. Neural Networks

Lars Schmidt-Thieme

Information Systems and Machine Learning Lab (ISMLL)
Institute for Computer Science
University of Hildesheim, Germany

Outline

1. Network Topologies
2. Stochastic Gradient Descent (Backpropagation)
3. Regularization

Syllabus

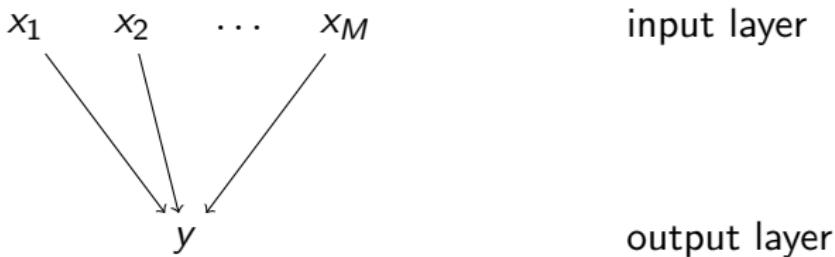
A. Advanced Supervised Learning

Tue. 9.12.	(1)	A.1 Generalized Linear Models
Wed. 10.12.	(2)	A.2 Gaussian Processes
Tue. 16.12.	(3)	A.3 Advanced Support Vector Machines
Wed. 17.12.	(4)	A.4 Neural Networks
Tue. 6.1.	(5)	
Wed. 7.1.	(6)	
Tue. 13.1.	(7)	
Wed. 14.1.	(8)	
Tue. 20.1.	(9)	
Wed. 21.1.	(10)	
Tue. 27.1.	(11)	
Wed. 28.1.	(12)	
Tue. 3.2.	(13)	
Wed. 4.2.	(14)	

Outline

1. Network Topologies
2. Stochastic Gradient Descent (Backpropagation)
3. Regularization

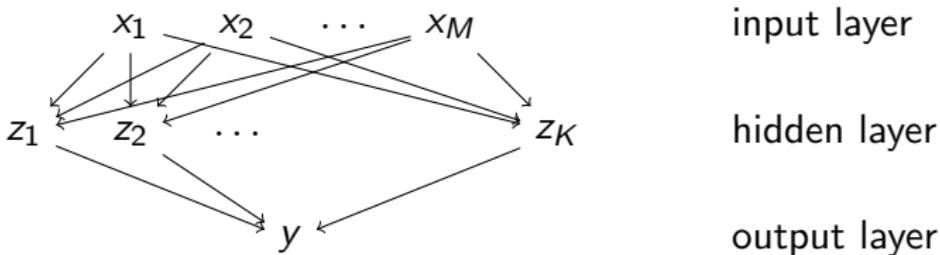
Logistic Regression (0 hidden layers)



logistic regression:

$$\hat{y}(x) := \hat{p}(y = 1 | x) = \text{logistic}(\beta^T x), \quad x \in \mathbb{R}^M$$

Feedforward Neural Network (1 hidden layer)

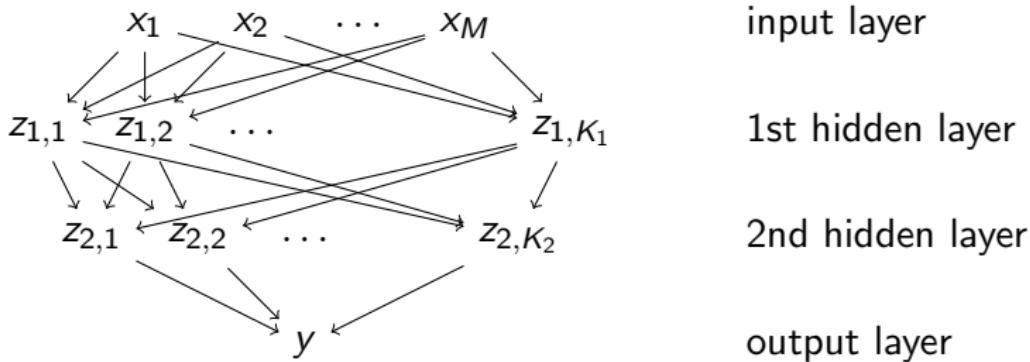


feedforward neural network (1 hidden layer):

$$\hat{y}(x) := \hat{p}(y = 1 \mid x) = \text{logistic}(\beta_2^T z(x))$$

$$z_k(x) := \text{logistic}(\beta_{1,k}^T x), \quad k = 1, \dots, K, x \in \mathbb{R}^M$$

Feedforward Neural Network (2 hidden layers)



feedforward neural network (2 hidden layers):

$$\hat{y}(x) := \hat{p}(y = 1 | x) = \text{logistic}(\beta_3^T z_2(x))$$

$$z_{2,k}(x) := \text{logistic}(\beta_{2,k}^T z_1), \quad k = 1, \dots, K_2$$

$$z_{1,k}(x) := \text{logistic}(\beta_{1,k}^T x), \quad k = 1, \dots, K_1, x \in \mathbb{R}^M$$

Different Targets y

Binary classification:

$$\hat{y}(x) := \hat{p}(y = 1 \mid x) = \text{logistic}(\beta_{L+1}^T z_L(x))$$

Regression:

$$\hat{y}(x) := \beta_{L+1}^T z_L(x)$$

Regression with multiple outputs:

$$\hat{y}(x) := \beta_{L+1} z_L(x), \quad \beta \in \mathbb{R}^{K_{\text{out}} \times K_L}$$

Multi-class classification:

$$\hat{y}(x) := \hat{p}(y \mid x) = \text{softmax}(\beta_{L+1} z_L(x))$$

Notes:

- ▶ L hidden layers
- ▶ at hidden nodes always are logistic/sigmoid functions
(activation function, transfer function).

Network Topologies

- ▶ **feedforward neural network** (aka **multiplayer perceptron**, MLP)

- ▶ often just a single hidden layer is used
 - ▶ NN with single hidden layer is already a **universal approximator**
- ▶ **skip arcs** can be used to connect layers skipping a hidden layer
- ▶ sometimes layers are not connected completely, but have **sparse connections**.
- ▶ nodes & connections always form a DAG

- ▶ **recurrent neural network**

- ▶ neural networks with backward connections / not a DAG.
- ▶ used in language modeling
- ▶ no simple probabilistic interpretation

- ▶ **Hopfield networks / associative memory:**

- ▶ symmetric connections between hidden units
- ▶ probabilistic counterpart: Boltzmann machine.

Outline

1. Network Topologies
2. Stochastic Gradient Descent (Backpropagation)
3. Regularization

SGD / Loss

feedforward neural network, L hidden layers with K_1, \dots, K_L nodes each:

$$z_{\ell,k}(x) := s(\beta_{\ell,k}^T z_{\ell-1}), \quad \beta_{\ell,k} \in \mathbb{R}^{K_{\ell-1}}, k = 1, \dots, K_\ell$$

$$z_{0,k} := x_k, \quad k = 1, \dots, K_0 := M$$

$$\hat{y}_k(x) := z_{L+1,k}(x), \quad k = 1, \dots, K_{L+1} := \dim \mathcal{Y}, \text{ usually } = 1$$

$$f(\beta) := \frac{1}{N} \sum_{n=1}^N \ell(y_n, \hat{y}(x)) + \frac{\lambda}{2} \|\beta\|^2$$

gradient for single sample x_n, y_n :

$$\frac{\partial f}{\partial \beta_{\ell,k}} = \ell'(y_n, \hat{y}(x_n)) \sum_{k'=1}^{K_{L+1}} \frac{\partial \hat{y}_{k'}(x_n)}{\partial \beta_{\ell,k}} + \lambda \beta_{\ell,k} = \ell'(y_n, \hat{y}_n) \sum_{k'=1}^{K_{L+1}} \frac{\partial z_{L+1,k'}}{\partial \beta_{\ell,k}} + \lambda \beta_{\ell,k}$$

SGD / Layers

feedforward neural network, L hidden layers with K_1, \dots, K_L nodes each:

$$z_{\ell,k}(x) := s(\beta_{\ell,k}^T z_{\ell-1}), \quad \beta_{\ell,k} \in \mathbb{R}^{K_{\ell-1}}, k = 1, \dots, K_\ell$$

gradient for single sample x_n, y_n :

$$\frac{\partial z_{\ell',k'}}{\partial \beta_{\ell,k}} = \sum_{\tilde{k}=1}^{K_{\ell-1}} \frac{\partial z_{\ell',k'}}{\partial z_{\ell'-1,\tilde{k}}} \frac{\partial z_{\ell'-1,\tilde{k}}}{\partial \beta_{\ell,k}}$$

$$\frac{\partial z_{\ell',k'}}{\partial z_{\ell,k}} = s'(\beta_{\ell',k'}^T z_{\ell'-1}) \sum_{\tilde{k}=1}^{K_{\ell-1}} \beta_{\ell',k',\tilde{k}} \frac{\partial z_{\ell'-1,\tilde{k}}}{\partial z_{\ell,k}}$$

$$\frac{\partial z_{\ell',k'}}{\partial z_{\ell'-1,k}} = s'(\beta_{\ell',k'}^T z_{\ell'-1}) \beta_{\ell',k',k}$$

$$\frac{\partial z_{\ell',k'}}{\partial \beta_{\ell'-1,k}} = s'(\beta_{\ell',k'}^T z_{\ell'-1}) z_{\ell'-1,k}$$

SGD / Arrange Computations

1. Feedforward:

$$\begin{aligned} z_{0,k} &:= x_k, & k &= 1, \dots, K_0 := M \\ z_{\ell,k} &:= s(\beta_{\ell,k}^T z_{\ell-1}), & \ell &= 1, \dots, L+1, k = 1, \dots, K_\ell \end{aligned}$$

2. Backpropagation:

for $\ell := L+1, \dots, 1, k := 1, \dots, K_\ell$:

for $k' := 1, \dots, K_{\ell+1}$:

$$\frac{\partial z_{L+1,k'}}{\partial z_{\ell,k}} = \sum_{\tilde{k}=1}^{K_{\ell+1}} \frac{\partial z_{L+1,k'}}{\partial z_{\ell+1,\tilde{k}}} \frac{\partial z_{\ell+1,\tilde{k}}}{\partial z_{\ell,k}} = \sum_{\tilde{k}=1}^{K_{\ell+1}} \frac{\partial z_{L+1,k'}}{\partial z_{\ell+1,\tilde{k}}} s'(\beta_{\ell+1,\tilde{k}}^T z_\ell) \beta_{\ell+1,\tilde{k}}, \text{ red}$$

$$\frac{\partial z_{L+1,k'}}{\partial \beta_{\ell,k}} = \frac{\partial z_{L+1,k'}}{\partial z_{\ell,k}} \frac{\partial z_{\ell,k}}{\partial \beta_{\ell,k}} = \frac{\partial z_{L+1,k'}}{\partial z_{\ell,k}} s'(\beta_{\ell,k}^T z_{\ell-1}) z_{\ell-1}$$

$$\frac{\partial f}{\partial \beta_{\ell,k}} = \ell'(y_n, \hat{y}_n) \sum_{k'=1}^{K_{\ell+1}} \frac{\partial z_{L+1,k'}}{\partial \beta_{\ell,k}} + \lambda \beta_{\ell,k}$$

SGD / Arrange Computations

2. Backpropagation:

for $\ell := L+1, \dots, 1$, $k := 1, \dots, K_\ell$:

for $k' := 1, \dots, K_{L+1}$:

$$\frac{\partial z_{L+1,k'}}{\partial z_{\ell,k}} = \sum_{\tilde{k}=1}^{K_{\ell+1}} \frac{\partial z_{L+1,k'}}{\partial z_{\ell+1,\tilde{k}}} \frac{\partial z_{\ell+1,\tilde{k}}}{\partial z_{\ell,k}} = \sum_{\tilde{k}=1}^{K_{\ell+1}} \frac{\partial z_{L+1,k'}}{\partial z_{\ell+1,\tilde{k}}} s'(\beta_{\ell+1,\tilde{k}}^T z_\ell) \beta_{\ell+1,\tilde{k},\text{red}}$$

$$\frac{\partial z_{L+1,k'}}{\partial \beta_{\ell,k}} = \frac{\partial z_{L+1,k'}}{\partial z_{\ell,k}} \frac{\partial z_{\ell,k}}{\partial \beta_{\ell,k}} = \frac{\partial z_{L+1,k'}}{\partial z_{\ell,k}} s'(\beta_{\ell,k}^T z_{\ell-1}) z_{\ell-1}$$

$$\frac{\partial f}{\partial \beta_{\ell,k}} = \ell'(y_n, \hat{y}_n) \sum_{k'=1}^{K_{\ell+1}} \frac{\partial z_{L+1,k'}}{\partial \beta_{\ell,k}} + \lambda \beta_{\ell,k}$$

$$\beta_{\ell,k}^{(t)} = \beta_{\ell,k}^{(t)} - \eta^{(t)} \frac{\partial f}{\partial \beta_{\ell,k}}$$

Outline

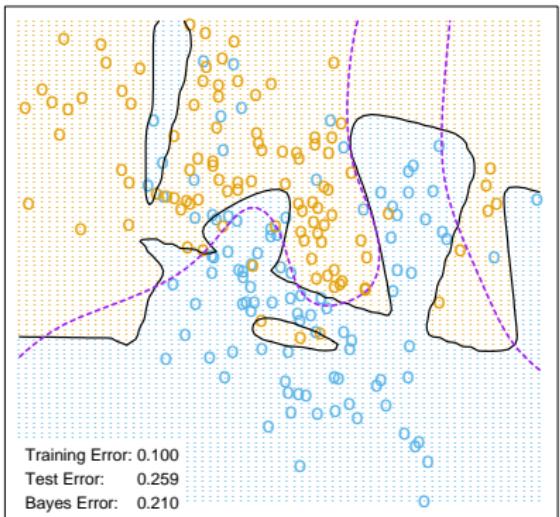
1. Network Topologies
2. Stochastic Gradient Descent (Backpropagation)
3. Regularization

Regularization of Neural Networks

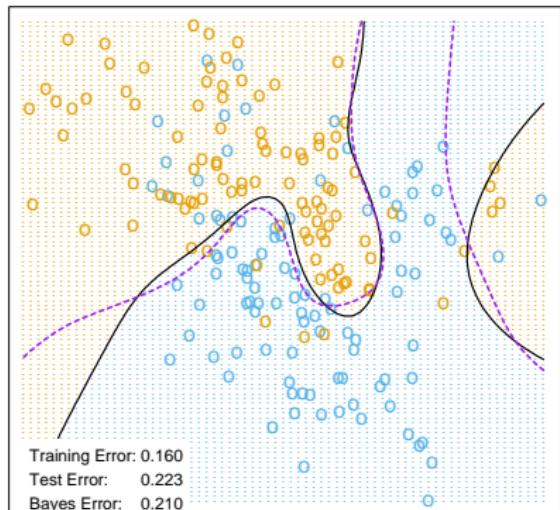
- ▶ L2 regularization
 - ▶ aka **weight decay**
 - ▶ most frequently used method
- ▶ L1 regularization
- ▶ **early stopping**
- ▶ use a random sample of connections
 - ▶ **drop connect**

L2 regularization / Example

Neural Network - 10 Units, No Weight Decay



Neural Network - 10 Units, Weight Decay=0.02



[HTFF05, p. 399]

Further Readings

- ▶ See [Mur12, chapter 16.5] and [HTFF05, chapter 11].

References

 Trevor Hastie, Robert Tibshirani, Jerome Friedman, and James Franklin.

The elements of statistical learning: data mining, inference and prediction, volume 27.
2005.

 Kevin P. Murphy.

Machine learning: a probabilistic perspective.
The MIT Press, 2012.