

Neuronal Networks

Formulae Collection

Marvin Ritter
marvin.ritter@gmail.com

23rd September 2015

About this document

This is my personal collection of formulae in the field of neural networks. Where I felt the need for it I also added explanations and derivations. I started it in preparation to my exam on the neural nets at Karlsruhe Institute of Technology (KIT). Even so the content is similar to the course topics it is neither limited to it nor guaranteed to cover it completely. *These are not official lecture notes and may contain errors and lag completeness.*

Corrections, supplements (or wishes for it) and links to good sources/papers are very welcome. Just mail me to marvin.ritter@gmail.com or create a ticket/pull request.

Contents

1	Introduction	4
1.1	Notation	4
1.2	Other sources	5
1.3	Software	5
2	Unsupervised Vector Quantization	6
2.1	Loss Functions for Clustering	6
2.2	k-Means (Lloyd's Algorithm)	7
2.3	Fuzzy k-Means	7
2.4	LBG-Algorithm	8
2.5	Self-Organizing Maps	9
3	Supervised Vector Quantization	10
3.1	Distortion Measures	10
3.2	Learning Vector Quantization	11
3.2.1	LVQ1	11
3.2.2	OLVQ	12
3.2.3	LVQ2.1	12
3.2.4	LVQ3	12
4	Perceptron	13
4.1	Gradient Descent Algorithm	13
5	Multi-Layer Perceptrons	14
5.1	Backpropagation	14
5.1.1	Notation	14
5.1.2	Weight Updates in Output Layer	14
5.1.3	Weight Outputs for Hidden Layers	15
5.1.4	Backpropagation using Matrix Notation	16
5.2	Error Functions	17
5.2.1	Mean-Squared Error	17
5.2.2	Cross Entropy Error	17
5.2.3	Classification Figure of Merit	18
5.2.4	L1 and L2 Regularization	18
5.3	Activation Functions	18
5.3.1	Step Function	18
5.3.2	Sigmoid	19
5.3.3	Softmax Function	19
5.3.4	Hyperbolic Tangent Function	21
5.3.5	Linear Function	22
5.3.6	Rectified Linear Unit	22
5.3.7	Softplus	22
5.3.8	Maxout	22
5.4	Learning Rate	23
5.4.1	Exponential Decaying Learning Rate	23

5.4.2	Performance Scheduling	23
5.4.3	Momentum	23
5.4.4	Resilient Propagation (Rprop)	24
5.4.5	AdaGrad	24
5.4.6	Newbob Scheduling	24
6	Boltzmann Machines	25
6.1	Hopfield Nets	25
6.1.1	Update Procedure	25
6.1.2	Energy function	26
6.1.3	Convergence	26
6.1.4	Associative Memory	26
6.1.5	Training	27
6.1.6	Limitations	27
6.2	Boltzmann Machines	27
6.2.1	Energy-based models	27
6.2.2	Energy of Boltzmann Machines	28
6.2.3	Training	28
6.2.4	Simulated Annealing	28
6.2.5	Why to restrict BMs	28
6.3	Restricted Boltzmann Machines	29
6.3.1	Energy	29
6.4	Training	29
6.5	Stacking RBMs	30
7	Autoencoders	31
7.1	Spares Autoencoders	31
7.2	Denoising Autoencoders	31
7.3	Stacking Autoencoders	33
8	Advanced Neural network architectures	34
8.1	Time Delay Neural Networks (TDNN)	34
8.2	Multi-State Time Delay Neural Networks (MS-TDNN)	34
8.3	Convolutional Neural Networks (CNN)	34
8.4	Recurrent Neural networks (RNN)	34
8.4.1	Vanishing Gradient	34
8.4.2	Backpropagation through Time	34
9	Practical Tricks	35
9.1	Data Normalization	35
9.2	Bottleneck-Features	35
9.3	Weight Decay	35
9.4	Dropout	36
	Acronyms	39

Chapter 1

Introduction

1.1 Notation

Throughout the literature (and history) in the big field of neural networks you will find a whole bunch of different notations and meanings of variable names. Especially when the meaning changes inside a document it can be very confusing for beginners. For this document I will try to stick with one notation. I found it to be very flexible and consistent.

Symbol	Meaning(s)
m	Number of features, size of input vector
k, K	Number of target classes/clusters, size of output vector
n, N	Number of training examples/observations, size of training set
η	Learning rate (see 5.4)
θ	Represents all model parameters
X	Training data
μ_j	Centroid of cluster j , also called codebook vector (with μ as codebook)
$d(x, y)$	Dissimilarity measure, distance or distortion measure, if mention otherwise we use the squared Euclidean distance (see 3.1)
E, J	Error function (see 5.2), other names are: loss function, objective function or criterion function By default we will use the mean squared error.
$w_{ij}^{(l)}$	Weight of connection from neuron i in layer $(l - 1)$ to neuron j in layer l , if the layer is omitted all variables in the equation are from the same layer
\mathbf{W}, \mathbf{w}	Weight matrix, weight vector
$b_j^{(l)}$	bias unit for neuron j in layer l For convenience we will use different meanings throughout the lecture notes: <ul style="list-style-type: none"> • Often the bias is not important for the understanding and left out, but keep in mind that neural network implementation nearly always require bias for good results. • b has an associated weight (usually w_{0j}), than the value of b is always 1 and we only train the weight • if there is no weight, we adjust b_j during training
Σ	Summation operator (e. g. $\sum_{i=1}^n x_i$), covariance matrix or the input function of a neuron
$\phi(x)$	Activation function (see 5.3)
$\sigma(x)$	Sigmoid activation function (see 5.3.2)
$z_j^{(l)}$	Total input of the neuron j in layer l , usually $z_j^{(l)} = b_j^{(l)} + \sum_i w_{ij}^{(l)} a_i^{(l-1)}$
$a_j^{(l)}$	Activation/Output of neuron j in layer l , usually computed from the total input $z_j^{(l)}$, $a_j = \phi(z_j)$
$v_j \in \{0, 1\}$	Activation of visible neuron/unit j in Boltzmann machine
$h_j \in \{0, 1\}$	Activation of hidden neuron/unit j in Boltzmann machine
$b_j(c_j)$	Bias of visible (hidden) neuron j in RBM
δ_{ij}	Kronecker delta, which is defined as $\delta_{ij} = \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{if } i \neq j. \end{cases}$

1.2 Other sources

Todo: Link some generally good sources (demos, tutorials, online courses, books).

Also have a look at the references at the end of the document.

1.3 Software

Todo: List with of good libraries for Machines Learning and Deep Neural Networks

Chapter 2

Unsupervised Vector Quantization

Also called *Clustering*.

Notation used in this chapter:

- N observations/training examples, x_1, \dots, x_N
- fixed number of clusters/classes K
- $C(i) = k$ assigns observation x_i to a cluster k (typical each observation is only in one cluster)
- μ_k centroid of cluster k
- N_k number of observations that belong to cluster k ($N_k = \sum_{C(i)=k} 1$)
- $d(x, y)$ is a dissimilarity measure (usually the squared Euclidean distance¹)
- u_{ik} degree of membership of training example i to cluster k (if used in algorithm)
- $m \geq 1$ the ‘fuzzifier’ parameter for Fuzzy k-Means (2.3). Influences sharpness of the clusters, for $m \rightarrow \infty$ u_{ik} will be close to $\frac{1}{K}$ (not sharp). For m close to 1 the u_{ik} s will be close to 0 or 1. Typical values are between 1 and 2.5.

2.1 Loss Functions for Clustering

- Intra-class scatter: $W(C) = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(j)=k} d(x_i, x_j)$
- Inter-class scatter: $B(C) = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(j) \neq k} d(x_i, x_j)$
- Total scatter: $T(C) = W(C) + B(C) = \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N d(x_i, x_j)$
- Minimizing $W(C)$ is equivalent to maximizing $B(C)$

¹Squared Euclidean distance $d(x, y) = \|x - y\|^2 = (x - y)^T (x - y) = \sum_i (x_i - y_i)^2$

2.2 k-Means (Lloyd's Algorithm)

Finding $C^*(i)$ by enumeration is too time-consuming. Instead we will use iterative greedy descent which leads to a local optimum. The error function is defined as

$$E = \sum_{k=1}^K \sum_{C(i)=k} d(x_i, \mu_k) \quad (2.1)$$

Using the squared Euclidean distance, $d(x_i, \mu_k) = \|x_i - \mu_k\|^2$, will effectively assign each sample to the closet center.

A slightly different definition is to minimize our already defined $W(C)$

$$W(C) = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(j)=k} \|x_i - x_j\|^2 = \sum_{k=1}^K N_k \sum_{C(i)=k} \|x_i - \mu_k\|^2 \quad (2.2)$$

Todo: Add derivate for this equation and why N_k is left out later

Also called *Lloyd's Algorithm* we optimize for J in a greedy fashion.

1. **Initialize:** There are many ways to initialize the centroids:

- distribute uniformly
- randomly (different distribution are possible)
- start with one random centroid and iteratively place to next one by maximizing the distance to all others

2. **Classify:** Assign each observation i to the nearest centroid:

$$C(i) = \arg \min_{1 \leq k \leq K} \|x_i - \mu_k\|^2$$

3. **Recenter:** For each class k , compute a new centroid as the mean of the updated class assignments:

$$\mu_k = \frac{\sum_{C(i)=k} x_i}{N_k}$$

4. **Repeat step 2 and 3 until stopping criteria fulfilled** (e.g. centers stop moving, or equally until $C(i)$ stays unchanged).

During the course of the k-means algorithm, the error function monotonically decreases into a local minimum. The proof is simply done by proofing that both steps will independently decrease E if possible.

2.3 Fuzzy k-Means

Literature: [\[Introduction, 2000\]](#)

In English literature it is sometimes also called *Fuzzy c-Means* (FCM).

Observation does not belong strictly to one cluster, but has each x_i a member ship degree u_{ik} for each cluster k .

$$0 \leq u_{ik} \leq 1 \quad \forall i, k \quad (2.3)$$

$$\sum_k u_{ik} = 1 \quad \forall i \quad (2.4)$$

$$(2.5)$$

Our new error function is

$$E(m) = \sum_{i=1}^N \sum_{k=1}^K u_{ik}^m d(x_i, \mu_k) \quad (2.6)$$

with fuzzifier $m \geq 1$ that controls fuzziness of the clusters. With $m = 1$ the membership degrees will converge to 0 or 1, which implies a crisp partitioning. In the absence of experimentation or domain knowledge, m is commonly set to 2.

1. **Compute degree of membership:**

$$u_{ik} = \left[\sum_{j=1}^K \left(\frac{d(x_i, \mu_k)}{d(x_i, \mu_j)} \right)^{\frac{2}{m-1}} \right]^{-1}$$

2. **Recenter:**

$$\mu_k = \frac{\sum_{i=1}^N u_{ik}^m x_i}{\sum_{i=1}^N u_{ik}^m}$$

3. **Repeat step 1 and 2 until stopping criteria fulfilled.**

2.4 LBG-Algorithm

Literature: [\[Linde et al., 1980\]](#), [Online Demo](#)

The Linde, Buzo, Grey (LBG) Algorithm is an alternative method to design a K -vector codebook (codebook vector = centroid), where $K = 2^x$.

The algorithm takes a parameter a small $\epsilon > 0$ (e. g. $\epsilon = 0.001$).

1. **Initialize** 1-vector starting codebook:

$$\mu_1(0) = \frac{1}{N} \sum_{i=1}^N x_i \quad (2.7)$$

$$t = 0$$

2. **Double codebook:**

$$\begin{aligned} \mu_j(t+1) &= (1 + \epsilon) \mu_j(t) \quad \text{for } j = 1, \dots, 2^t \\ \mu_{2*j+1}(t+1) &= (1 - \epsilon) \mu_j(t) \quad \text{for } j = 1, \dots, 2^t \\ t &= t + 1 \end{aligned} \quad (2.8)$$

3. **Run k-Means** Use k-Means to optimize to improve the $\mu_j(t)$.
4. **Repeat step 2 and 3** until desired number of codebook vectors is reached (for $t = 1, \dots, \log_2(K) - 1$)

Our final codebook will be $\mu_j(\log_2(K))$.

This simple form of the *LBG-Algorithm* doesn't have any advantages over k-means, but if splitting of the centroids is done smart, the final centroids are more likely to be a good representation of the data:

- k-Means initialization is random and not data dependent. Small groups of samples might get equal amount or even more centroids than big groups.
- ϵ can be set to split the centroids according to the direction with the highest variance for each dimension.
- Instead of splitting all centroids only the with biggest variance or biggest number of assigned samples can be splitted.

2.5 Self-Organizing Maps

Sometimes also called self-organizing feature maps.

Chapter 3

Supervised Vector Quantization

3.1 Distortion Measures

Literature: [Linde et al., 1980]

Distortion is a nonnegative measure for an input vector x and a reproduction vector y . Distortion measures may not be true metrics, e. g. be unsymmetric or not fulfill the triangle inequality¹.

Most common is the squared-error distortion, better known as squared Euclidean distance:

$$d(x, y) = \|x - y\|_2^2 = (x - y)^T (x - y) = \sum_{i=1}^N (x_i - y_i)^2 \quad (3.1)$$

Other common distortion measures are the l_ν , or Holder norm,

$$d(x, y) = \left(\sum_{i=1}^N |x_i - y_i|^\nu \right)^{\frac{1}{\nu}} = \|x - y\|_\nu \quad (3.2)$$

and its ν^{th} power, the ν^{th} -law distortion:

$$d(x, y) = \sum_{i=1}^N |x_i - y_i|^\nu \quad (3.3)$$

The holder Norm (3.2) is a distance and fulfills the triangle inequality, the ν^{th} -law distortion not.

All three and many others, as the weighted-squares distortion and the quadratic distortion, only depend on the difference $(x - y)$ and can be described with $d(x, y) = L(x - y)$. A distortion not having this form is the one by Itakura, Saito and Chaffee,

$$d(x, y) = (x - y)^T R(x) (x - y) \quad (3.4)$$

, where $R(x)$ is the autocorrelation matrix.

¹triangle inequality: $d(x, y) \leq d(x, z) + d(z, y)$, for all z

3.2 Learning Vector Quantization

Literature: [Kohonen, 2001, Chapter 6] and [Kohonen, 1990]

Learning vector quantization (LVQ) are used for statistical classification. Supervised learning is used to train multiple codebook vector per class. The class of new observations is then determined by the class of the closest codebook vector.

Kohonen published several versions of LVQ that only differ slightly in the conditions and formulae for updates during the training.

Notation:

- discrete time domain with $t = 0, 1, 2, \dots$
- $x(t)$ is an input sample
- m_i our vector vectors and $m_i(t)$ their sequential values
- classes S_1, \dots, S_K
- each codebook vectors belongs to exactly one class, but each class has multiple codebook vectors
- $d(x, y)$ is a distance measure (usually Euclidean distance is used)
- $c = \arg \min_i d(x, m_i)$ is the index of the closest codebook vector
- δ_{ij} is the Kronecker delta ($\delta_{ij} = 1$ for $i = j$, $\delta_{ij} = 0$ for $i \neq j$)
- $\eta(t)$ learning rate at time t
- w window width (see LVQ2)

For convenience we also define $s(t)$ as

$$s(t) = \begin{cases} +1 & \text{if } x \text{ and } m_c \text{ belong to the same class} \\ -1 & \text{if } x \text{ and } m_c \text{ belong to different classes} \end{cases} \quad (3.5)$$

3.2.1 LVQ1

Update closest codebook vector m_c depending on whether it has the same class as x or not.

If x and $m_c(t)$ belong to the same class set

$$m_c(t+1) = m_c(t) + \eta(t)(x - m_c(t)) \quad (3.6)$$

If x and $m_c(t)$ belong to different classes

$$m_c(t+1) = m_c(t) - \eta(t)(x - m_c(t)) \quad (3.7)$$

Leave the other codebook vectors unchanged

$$m_i(t+1) = m_i(t) \text{ for } i \neq c \quad (3.8)$$

Or in a compressed form:

$$m_i(t+1) = m_i(t) + \eta(t)s(t)\delta_{ci}(x(t) - m_i(t)) \quad (3.9)$$

3.2.2 OLVQ

Determine the optimal learning rate $\eta_i(t)$ for each codebook vector m_i for fastest convergence.

$$\eta_c(t) = \frac{\eta_c(t-1)}{1 + s(t)\eta_c(t-1)} \quad (3.10)$$

Modify update rule from LVQ1.

$$m_i(t+1) = [1 - \eta_i(t)s(t)\delta_{ci}]m_i(t) + \eta_i(t)s(t)\delta_{ci}x(t) \quad (3.11)$$

This modification will not work for LVQ2, but could be applied to LVQ3.

3.2.3 LVQ2.1

m_i and m_j are the two closest codebook vectors to x . One of them must belong to the correct class x and the other to a different class. Moreover, x must be inside a ‘window’. The window is defined around the mid-plane of m_i and m_j and depends on the window width w .

$$\min\left\{\frac{d(x, m_i)}{d(x, m_j)}, \frac{d(x, m_j)}{d(x, m_i)}\right\} > \frac{1-w}{1+w} \quad (3.12)$$

Good values for w are in $[0.2, 0.3]$.

Assume that m_j belongs to the same class as x ($\Rightarrow x$ and m_i belong to different class), then m_i and m_j are updated simultaneously and similar to LVQ1.

$$\begin{aligned} m_i(t+1) &= m_i(t) - \eta(t)[x(t) - m_i(t)] \\ m_j(t+1) &= m_j(t) + \eta(t)[x(t) - m_j(t)] \end{aligned} \quad (3.13)$$

In the original LVQ2 m_i had to be the closest codebook vector and belong to a different class.

3.2.4 LVQ3

LVQ3 has the same update rule as LVQ2.1, but in case that x, m_i and m_j all belong to the same class we will perform the following update:

$$m_k(t+1) = m_k(t) + \epsilon\eta(t)[x(t) - m_k(t)] \quad \text{for } k \in i, j \quad (3.14)$$

Applicable values for ϵ are between 0.1 and 0.5 and seem to depend on the window size.

Chapter 4

Perceptron

4.1 Gradient Descent Algorithm

This is the base for most learning algorithms and follows a very simple approach. Derive to error function E with respect to the set of parameters θ and use the gradient and the learning rate η to update θ .

```
while not converged do  
   $\theta_j \leftarrow \theta_j - \eta \frac{\partial E}{\partial \theta_j} \quad (j = 1, \dots, |\theta|)$   
end
```

Algorithm 1: Gradient descent algorithm

Chapter 5

Multi-Layer Perceptrons

5.1 Backpropagation

Literature: [Mitchell, 1997, Chapter 4.5], [Duda et al., 2000, Chapter 6.3], [Patterson, 1997] and the original paper [Rumelhart et al., 1986]

5.1.1 Notation

- m inputs/features, $x \in \mathbb{R}^m$
- k target outputs, $t \in \mathbb{R}^k$
- n training examples of form $(x, t) \in \mathbb{R}^m \times \mathbb{R}^k$
- L layers $(1, \dots, L)$
- E : error function (e.g. $E_{\text{MSE}} = \frac{1}{2} \sum_{i=1}^k (t_i - o_i^{(L)})^2$)
- $\phi(y)$: activation function (e.g. $\phi(y) = \sigma(y) = \frac{1}{1+e^{-y}}$)
- $\delta_j^{(l)}$: error of neuron j in layer l

For convenience we define $a_j^{(0)} := x_j$ and $o_j := a_j^{(L)}$.

5.1.2 Weight Updates in Output Layer

Backpropagation is generalization of the delta-rule. We minimize our error function E by ‘going down’ along the gradient. For a single training example (x, t) , with x as input and t and desired output, E is calculated from o and t , where o is the result of a forward propagation using x and our weights $w_{ij}^{(l)}$. As the training example is fixed we can only adjust $w_{ij}^{(l)}$ to minimize E .

We start with the gradient in our output layer:

$$\begin{aligned}
\frac{\partial E}{\partial w_{ij}^{(L)}} &= \frac{\partial E}{\partial z_j^{(L)}} \frac{\partial z_j^{(L)}}{\partial w_{ij}^{(L)}} = \frac{\partial E}{\partial a_j^{(L)}} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \frac{\partial z_j^{(L)}}{\partial w_{ij}^{(L)}} && \text{(chain rule)} \\
\frac{\partial z_j^{(L)}}{\partial w_{ij}^{(L)}} &= \frac{\partial \sum_i w_{ij}^{(L)} * x_i^{(L-1)}}{\partial w_{ij}^{(L)}} = a_i^{(L-1)} && (5.1) \\
\frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} &= \frac{\partial \phi(z_j^{(L)})}{\partial z_j^{(L)}} = (\phi(z_j^{(L)})(1 - \phi(z_j^{(L)}))) = o_j(1 - o_j) && \text{(sigmoid derivative, see 5.3.2)} \\
\frac{\partial E}{\partial a_j^{(L)}} &= \frac{\partial E}{\partial o_j} = \frac{\partial \frac{1}{2} \sum_{i=1}^k (t_i - o_i)^2}{\partial o_j} = (o_j - t_j) && \text{(MSE derivative, see 5.2.1)}
\end{aligned}$$

putting everything together

$$\frac{\partial E}{\partial w_{ij}^{(L)}} = (o_j - t_j) o_j (1 - o_j) a_i^{(L-1)} \quad (5.2)$$

We can now define the error for neuron j as

$$\begin{aligned}
\delta_j^{(L)} &:= \frac{\partial E}{\partial z_j^{(L)}} \\
&= \frac{\partial E}{\partial a_j^{(L)}} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \\
&= (o_j - t_j) o_j (1 - o_j)
\end{aligned} \quad (5.3)$$

, and our weight update

$$\Delta w_{ij}^{(L)} := -\eta \frac{\partial E}{\partial w_{ij}^{(L)}} = -\eta \delta_j^{(L)} a_i^{(L-1)} \quad (5.4)$$

$$w_{ij}^{(L)} \leftarrow w_{ij}^{(L)} + \Delta w_{ij}^{(L)} \quad (5.5)$$

5.1.3 Weight Outputs for Hidden Layers

For hidden layers we need to propagate the error back to the neuron j in layer l . Intuitively what are doing, is just inserting the error into the output layer and propagate backwards to the

input layer using.

$$\delta_j^{(l)} := \frac{\partial E}{\partial z_j^{(l)}} = \frac{\partial E}{\partial a_j^{(l)}} \frac{\partial a_j^{(l)}}{\partial z_j^{(l)}} \quad (5.6)$$

$$\frac{\partial E}{\partial a_j^{(l)}} = \sum_k \frac{\partial E}{\partial z_k^{(l+1)}} \frac{\partial z_k^{(l+1)}}{\partial a_j^{(l)}} \quad (5.7)$$

$$= \sum_k \frac{\partial E}{\partial z_k^{(l+1)}} \frac{\partial \sum_i w_{ik}^{(l+1)} * a_i^{(l)}}{\partial a_j^{(l)}} \quad (5.8)$$

$$= \sum_k \frac{\partial E}{\partial z_k^{(l+1)}} w_{jk}^{(l+1)} \quad (5.9)$$

$$= \sum_k \delta_k^{(l+1)} w_{jk}^{(l+1)} \quad (5.10)$$

$$\frac{\partial a_j^{(l)}}{\partial z_j^{(l)}} = a_j^{(l)}(1 - a_j^{(l)}) \quad (\text{sigmoid derivative})$$

$$\delta_j^{(l)} = a_j^{(l)}(1 - a_j^{(l)}) \sum_k \delta_k^{(l+1)} w_{jk}^{(l+1)} \quad (5.11)$$

This leads to our weight updates for hidden neurons.

$$\Delta w_{ij}^{(l)} = -\eta \frac{\partial E}{\partial w_{ij}^{(l)}} \quad (5.12)$$

$$= -\eta \frac{\partial E}{\partial z_j^{(l)}} \frac{\partial z_j^{(l)}}{\partial w_{ij}^{(l)}} \quad (5.13)$$

$$= -\eta \delta_j^{(l)} a_i^{(l-1)} \quad (5.14)$$

$$w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} + \Delta w_{ij}^{(l)} \quad (5.15)$$

5.1.4 Backpropagation using Matrix Notation

The algorithm can also be formulated using matrices.

- $X \in \mathbb{R}^{n \times m}$ input of training data
- $T \in \mathbb{R}^{n \times k}$ target output for training data
- $W^{(1)}, \dots, W^{(L)}$ weight matrices
- $\delta^{(1)}, \dots, \delta^{(L)}$ error matrices

TODO, the following might me incomplete and wrong

$$A^{(0)} = X \quad (5.16)$$

$$Z^{(l)} = A^{(l-1)} * W^{(l)} \quad (5.17)$$

$$A^{(l)} = \phi(Z^{(l)}) \quad (5.18)$$

$$O = A^{(L)} \quad (5.19)$$

$$E = \frac{1}{2} \sum (T - O) \circ (T - O) \quad (5.20)$$

$$\delta^{(L)} = (O - T) \circ O \circ (1 - O) \quad (5.21)$$

$$\delta^{(l)} = A^{(l)}(1 - A^{(l)})W^{(l+1)}\delta^{(l+1)}\Delta W^{(l)} = -\eta \delta^{(l)} A^{(l-1)} \quad (5.22)$$

$$\hat{W}^{(l)} = W^{(l)} + \Delta W^{(l)} \quad (5.23)$$

5.2 Error Functions

Sometimes also called objective functions or loss-functions.

5.2.1 Mean-Squared Error

As the name suggests the mean-square error (MSE) is defined as the mean (over all training examples) of squared difference between the correct value t_i and the correct value o_i . This big errors are punished harder than small differences.

$$E_{\text{MSE}} = \frac{1}{n} \sum_{i=1}^n (\mathbf{t}_i - \mathbf{o}_i)^2 \quad (5.24)$$

This would derive to

$$\begin{aligned} \frac{\partial E_{\text{MSE}}}{\partial \mathbf{o}_j} &= \frac{1}{n} \sum_{i=1}^n \frac{\partial (\mathbf{t}_i - \mathbf{o}_i)^2}{\partial \mathbf{o}_j} \\ &= \frac{1}{n} \frac{\partial (\mathbf{t}_j - \mathbf{o}_j)^2}{\partial \mathbf{o}_j} \\ &= \frac{1}{n} 2(\mathbf{t}_j - \mathbf{o}_j)(-1) \\ &= \frac{2}{n} (\mathbf{o}_j - \mathbf{t}_j) \end{aligned} \quad (5.25)$$

which is technically fine, but as the fraction is only a constant factor, we will often use a slightly different definition:

$$E_{\text{MSE}} = \frac{1}{2} \sum_{i=1}^k (t_i - o_i)^2 \quad (5.26)$$

$$\frac{\partial E_{\text{MSE}}}{\partial o_j} = -(t_j - o_j) = o_j - t_j \quad (5.27)$$

5.2.2 Cross Entropy Error

Cross-entropy (CE) works great well on classifications tasks. t_i is either 0 or 1 and o_i is the class probability computed by the network $\Rightarrow o_i \in (0, 1]$ as $\log(0)$ is not defined.

$$E_{\text{CE}} = - \sum_{i=1}^k (t_i \log(o_i) + (1 - t_i) \log(1 - o_i)) \quad (5.28)$$

This will derive to:

$$\begin{aligned} \frac{\partial E_{\text{CE}}}{\partial o_j} &= - \frac{\partial}{\partial o_j} \sum_{i=1}^k (t_i \log(o_i) + (1 - t_i) \log(1 - o_i)) \\ &= - \frac{\partial}{\partial o_j} t_j \log(o_j) - \frac{\partial}{\partial o_j} (1 - t_j) \log(1 - o_j) \\ &= - \frac{t_j}{o_j} + \frac{1 - t_j}{1 - o_j} \\ &= \frac{1 - t_j}{1 - o_j} - \frac{t_j}{o_j} \end{aligned} \quad (5.29)$$

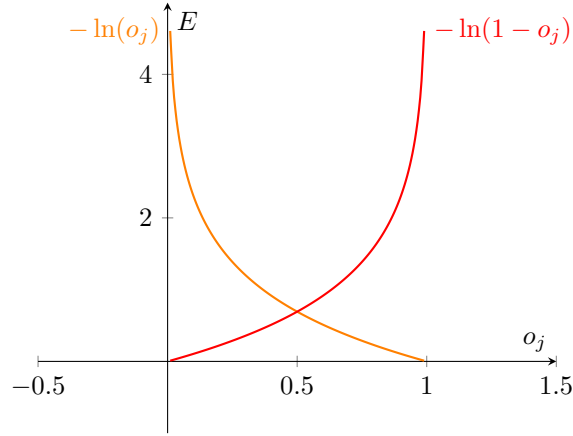


Figure 5.1: Error produced for $t_j = 1$ (orange) and for $t_j = 0$ (red). o_j must be in range $(0, 1]$, otherwise the logarithm can be complex or undefined or the error gets rediculously large.

5.2.3 Classification Figure of Merit

Literature: [Hampshire and Waibel, 1990]

The MSE and the CE functions focus on the difference of the network output \mathbf{o} and target activations \mathbf{t} and try to minimize it. Classification figure of merit (CFM) takes a different approach by only using the target activations to identify the correct class c and then tries to maximize the difference between o_c and all other $o_i (i \neq j)$.

$$E_{\text{CFM}} = \frac{1}{k-1} \sum_{i=1, i \neq c}^k \frac{\alpha}{1 + e^{-\beta(o_c - o_i) + \gamma}} \quad (5.30)$$

Where c is the correct class, k the number of classes and α , β and γ are just parameters of the sigmoid function (scaling, discontinuity and lateral shift).

5.2.4 L1 and L2 Regularization

Small weights tend to perform better and we can modify the error functions from above to penalize big weight by adding additional error terms.

- L1 norm: $\|\mathbf{w}\|_{L1} = \sum_j |\mathbf{w}_j|$
- L2 norm: $\|\mathbf{w}\|_{L2} = \sum_j \mathbf{w}_j^2$
- new error function $E' = E + \alpha_{L1} \|\mathbf{w}\|_{L1} + \alpha_{L2} \|\mathbf{w}\|_{L2}$
- new hyperparamaters α_{L1} and α_{L2} , optimize on development set

Todo: Explain advantages and disadvantages, see <http://www.chioka.in/differences-between-l1-and-l2-as-loss-function-and-regularization>

5.3 Activation Functions

5.3.1 Step Function

$$\phi(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases} \quad (5.31)$$

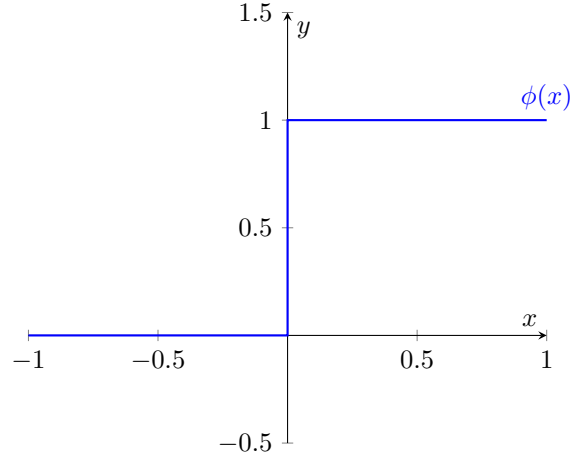


Figure 5.2: Step Function

The derivative is always 0.

5.3.2 Sigmoid

Most common activation function, can saturate and is easy to derivate.

$$\sigma(x) = \frac{1}{1 + e^{-\beta x}} \quad (5.32)$$

And it's derivative:

$$\begin{aligned} \frac{d\sigma(x)}{dx} &= \frac{d}{dx}(1 + e^{-\beta x})^{-1} \\ &= (-1)(1 + e^{-\beta x})^{-2}(-\beta e^{-\beta x}) \\ &= \beta(1 + e^{-\beta x})^{-2}e^{-\beta x} \\ &= \beta \frac{1}{1 + e^{-\beta x}} \frac{e^{-\beta x}}{1 + e^{-\beta x}} \\ &= \beta \sigma(x) \frac{e^{-\beta x} + 1 - 1}{1 + e^{-\beta x}} \\ &= \beta \sigma(x) \left(1 - \frac{1}{1 + e^{-\beta x}}\right) \\ &= \beta \sigma(x)(1 - \sigma(x)) \end{aligned} \quad (5.33)$$

And because there exists algorithm using this, you should have seen the second derivative as well ($\beta = 1$):

$$\begin{aligned} \frac{d\sigma(x)(1 - \sigma(x))}{dx} &= (1 - \sigma(x)) \frac{d\sigma(x)}{dx} + \sigma(x) \frac{d(1 - \sigma(x))}{dx} \\ &= (1 - \sigma(x))\sigma(x)(1 - \sigma(x)) + \sigma(x)(-1)\sigma(x)(1 - \sigma(x)) \\ &= \sigma(x)(1 - \sigma(x))^2 - \sigma(x)^2(1 - \sigma(x)) \end{aligned} \quad (5.34)$$

5.3.3 Softmax Function

In a classification problem we would like \mathbf{a} to be a probability distribution ($\Rightarrow \sum_i a_i = 1$). The softmax function will output a posteriori probability ($p(k|x)$) and the biggest value from the layer

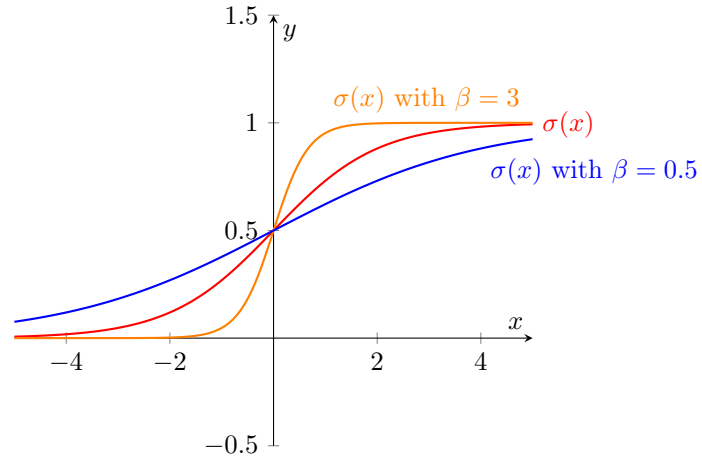


Figure 5.3: Sigmoid function with $\beta = 1$ (red), $\beta = 3$ (orange) and $\beta = 0.5$ (blue).

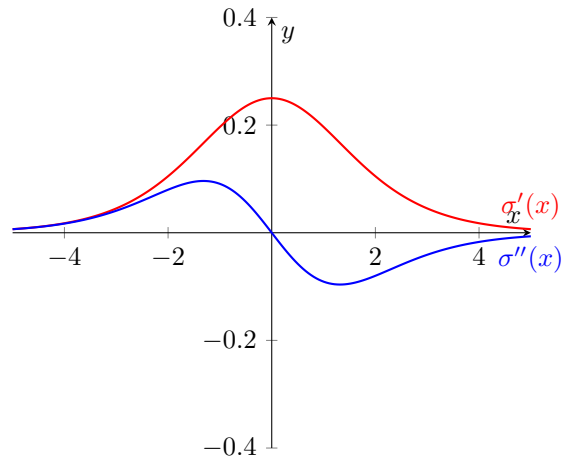


Figure 5.4: First and second derivative of the sigmoid function.

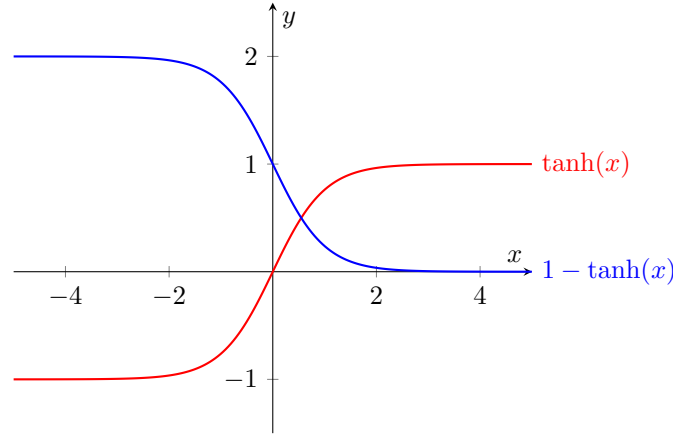


Figure 5.5: Hyperbolic Tangent Function (red) and its derivative (blue).

below is very likely to be close to 1.

$$a_i = \phi(z_i) = \frac{e^{z_i}}{\sum_j e^{z_j}} \quad (5.35)$$

And it's derivative:

$$\begin{aligned} \frac{\partial \phi(z_i)}{\partial z_i} &= \frac{\left(\frac{\partial}{\partial z_i} e^{z_i} \right) (\sum_j e^{z_j}) - e^{z_i} \left(\frac{\partial}{\partial z_i} (\sum_j e^{z_j}) \right)}{(\sum_j e^{y_j})^2} \\ &= \frac{e^{y_i} \sum_j e^{y_j} - e^{y_i} e^{y_i}}{(\sum_j e^{y_j})^2} \\ &= \frac{e^{y_i} \sum_j e^{y_j}}{(\sum_j e^{y_j})^2} - \frac{e^{y_i} e^{y_i}}{(\sum_j e^{y_j})^2} \\ &= \frac{e^{y_i}}{\sum_j e^{y_j}} - \left(\frac{e^{y_i}}{\sum_j e^{y_j}} \right)^2 \\ &= \phi(y_i) - \phi(y_i)^2 \\ &= \phi(y_i)(1 - \phi(y_i)) \end{aligned} \quad (5.36)$$

5.3.4 Hyperbolic Tangent Function

Similar to Sigmoid function, but with outputs in $[-1, 1]$. Also note that if the input has a mean of 0 then so will the output.

Todo: Why is tanh worth the computational overhead over sigmoid?

$$\phi(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \quad (5.37)$$

And it's derivative:

$$\frac{d\phi(x)}{dx} = 1 - \tanh(x) = \frac{(e^x - e^{-x})^2}{(e^x + e^{-x})^2} \quad (5.38)$$

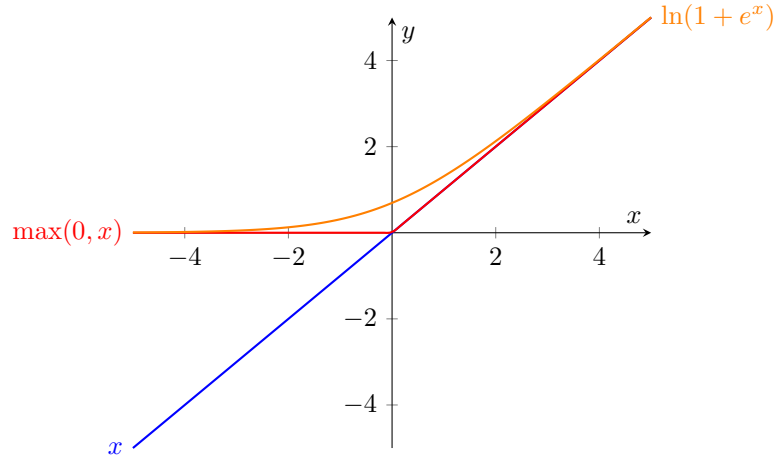


Figure 5.6: Linear Function (blue), Rectified Linear Unit (red), Softplus (orange)

5.3.5 Linear Function

Very simple, but with two big downsides. The gradient is always 1 and multi-layers using it effectively only do a linear combination and therefore can be left out.

$$\phi(x) = x$$

5.3.6 Rectified Linear Unit

The Rectified Linear Unit (ReLU) is more biologically plausible.

$$\phi(x) = \max(0, x) \quad (5.39)$$

And it's derivative

$$\frac{d\phi(x)}{dx} = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases} \quad (5.40)$$

5.3.7 Softplus

Smoothed version of ReLU.

$$\phi(x) = \ln(1 + e^x) \quad (5.41)$$

And it's derivative

$$\frac{d\phi(x)}{dx} = \frac{e^x}{1 + e^x} = \frac{1}{e^{-x} + 1} \quad (5.42)$$

5.3.8 Maxout

Literature: [Goodfellow et al., 2013]

Outputs the maximum of its inputs

$$a_j^{(l)} = \max_i z_i^{(l)}, \quad (j-1)g + 1 \leq i \leq jg \quad (5.43)$$

where $z_i^{(l)} = b_i^{(l)} + \sum_k w_{ki}^{(l)} x_k^{(l-1)}$ is the summed input and g the size of the group on which we take the maximum.

5.4 Learning Rate

Literature: [Duda et al., 2000, Chapter 6.8]

We use the learning rate η in our weight updates:

$$w_{ij} \leftarrow w_{ij} + \eta \delta_j a_i$$

There are several problems with a constant learning rate:

- large learning rate
 - constantly jump past the minimum (\Rightarrow no convergence)
 - neurons may saturate
- small learning rate
 - slow training
 - may get stuck in local minimum

And of course several ways to deal with it them:

- Predetermined piecewise constant learning rate
Use a predetermined sequence of η_i , increment i after every iteration (epoch/batch)
- Exponentially decaying learning rate
- Performance Scheduling
- Weight Dependent Learning Rate Methods (e. g. Momentum, AdaGrad)

5.4.1 Exponential Decaying Learning Rate

We calculate a new learning rate $\eta(t)$ after each iteration (epoch/batch).

$$\eta(t) = \alpha \eta(t-1) = \alpha^t \eta(0) \quad (5.44)$$

With the parameter $\alpha \in (0, 1)$ and we still have to define an initial learning rate $\eta(0)$.

5.4.2 Performance Scheduling

Measure the cross validation error and decrease the learning rate when the error stops improving.

5.4.3 Momentum

Similar to a rock rolling down a hill we let our gradient get some momentum to speed up training while the direction (sign) of the gradient stays the same. We do so by modifying our weight update rule to,

$$w_{ij} \leftarrow w_{ij} + v_{ij}(t) \quad (5.45)$$

$$v_{ij}(t) := \alpha v_{ij}(t-1) - \eta \delta_j a_i^{(l-1)} \quad (5.46)$$

, with α as momentum factor.

A lot of research has gone into the question on how to pick/adjust both η and α , a good start for the discussion is [Yu and Chen, 1997].

5.4.4 Resilient Propagation (Rprop)

Literature: [Riedmiller and Braun, 1994]

Same idea as *Momentum* but just a bit more complicated. Only the sign of the gradient $\frac{\partial E(t)}{\partial w_{ij}}$ is used to adjust the values $v_{ij}(t)$. It also introduces several new hyperparameters (v_0 , v_{\min} , v_{\max} , α^+ , α^-) instead of just one.

$$w_{ij} \leftarrow \begin{cases} w_{ij} - v_{ij}(t) & \text{if } \frac{\partial E(t)}{\partial w_{ij}} > 0 \\ w_{ij} + v_{ij}(t) & \text{if } \frac{\partial E(t)}{\partial w_{ij}} < 0 \end{cases} \quad (5.47)$$

$$v_{ij}(0) = v_0 \quad (5.48)$$

$$v_{ij}(t) = \begin{cases} \alpha^+ v_{ij}(t-1) & \text{if } \frac{\partial E(t-1)}{\partial w_{ij}} \frac{\partial E(t)}{\partial w_{ij}} > 0 \\ \alpha^- v_{ij}(t-1) & \text{if } \frac{\partial E(t-1)}{\partial w_{ij}} \frac{\partial E(t)}{\partial w_{ij}} < 0 \\ v_{ij}(t-1) & \text{else} \end{cases} \quad (5.49)$$

In practice the starting weight update v_0 should be a small and $v_{ij}(t)$ should be limited to $[v_{\min}, v_{\max}]$. Good values for α^- and α^+ are 0.5 and 1.2.

5.4.5 AdaGrad

Literature: [Zeiler, 2012, Duchi et al., 2011, Dyer,]

AdaGrad alters the weight update to adapt based on historical information, so that frequently occurring features in the gradients get small learning rates and infrequent features get higher ones.

$$g_{ij} = \frac{\partial E}{\partial w_{ij}} \quad (5.50)$$

$$\eta_{ij}(t) = \frac{\eta_0}{\sqrt{\sum_{\tau=1}^t g_{ij}(\tau)^2}} g_{ij}(t) \quad (5.51)$$

$$w_{ij}(t) \leftarrow w_{ij}(t-1) - \eta_{ij}(t) g_{ij}(t) \quad (5.52)$$

This eliminates the hyperparameter η , but introduces an initial learning rate η_0 and has some more downsides [Zeiler, 2012]. One way to improve it is the limit the window for the sum to size w .

$$\eta_{ij}(t) = \frac{\eta_0}{\sqrt{\sum_{\tau=1}^w g_{ij}(t-\tau)^2}} g_{ij}(t) \quad (5.53)$$

AdaDelta contains further improvements.

5.4.6 Newbob Scheduling

Combination of exponential decaying learning rate and performance scheduling. Easy to implement and solid results.

Phase 1 Start with constant learning rate $\eta(0)$

Phase 2 Once the cross-validation error stops decreasing switch to exponential decaying learning rate

Terminate End training after cross-validation error stops decreasing again

Chapter 6

Boltzmann Machines

6.1 Hopfield Nets

Literature: [Patterson, 1997, Chapter 5.5], [Storkey, 1997]

In the following we will only consider binary Hopfield nets, in which the neurons are limited to states 0 (or $-/-1$) and 1 (or $+/+1$) (sometimes also denoted as $+$ and $-$ or -1 and $+1$). Hopfield nets are not very efficient, but good to show the principle of neural nets and a good introduction to Boltzmann machines.

Hopfield nets consist of a single layer of neurons, that is fully connected and acts both as input and output layer. Weights are stored in a weight matrix W , where the entry w_{ij} corresponds to the connection weight from neuron i to neuron j . We will not allow self loops ($w_{ii} = 0$) and the connections will be symmetric ($w_{ij} = w_{ji}$). There exist variations without those limitations.

6.1.1 Update Procedure

1. Set state x to our input
2. Update neuron state using

$$x_j = \phi \left(\sum_i w_{ij} x_i \right)$$

where $\phi(z)$ is the step function

$$\phi(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{if } z \leq 0 \end{cases} \quad (6.1)$$

The updates can be performed in different ways:

- synchronous: all neurons get the same input (parallel, we cannot guarantee that the energy function E , see next section, decreases)
 - asynchronous (sequential): one neuron at the time (next neuron gets updated state), order can be fixed or random (random order is closest to biological neural nets)
 - combination are possible (group of neurons in parallel, similar to mini-batch training)
3. Repeat step 2 until x stabilizes.
 4. State x is our output

x is now a state that minimizes the energy function E (see next section). If x is a *retrieval state* it matches exactly one of our training patterns, otherwise it is a *spurious state* and can not be recognized.

If we want to remember more than the state vector (e.g. class of the pattern) we have to store these ourselves during training and use x to look them up.

6.1.2 Energy function

We define the energy in our Hopfield net as:

$$E = -\frac{1}{2} \sum_i \sum_{j \neq i} x_i x_j w_{ij} \quad (6.2)$$

E is also the objective function that we try to minimize during the update procedure. A state x that is a local minimum of E is also called attractor.

6.1.3 Convergence

If only one neuron j is updated at the time, the update will always lead to the same or lower energy.

- $x_j(t+1) = x_j(t) \Rightarrow$ state did not change, energy stays the same
- $x_j(t+1) = 1 - x_j(t) \Rightarrow$ state did change, compute the energy change

$$E_j = -\frac{1}{2} \sum_i x_i x_j w_{ij} \quad (6.3)$$

$$\Delta E = E_j(t+1) - E_j(t) \quad (6.4)$$

$$= -\frac{1}{2} \left[x_j(t+1) \sum_i x_i w_{ij} - x_j(t) \sum_i x_i w_{ij} \right] \quad (6.5)$$

$$= -\frac{1}{2} \Delta x_j \sum_i x_i w_{ij} \quad \text{with} \quad \Delta x_j = x_j(t+1) - x_j(t) \quad (6.6)$$

Change from 0 to 1

$$\Delta x_j = 1, \sum_i x_i w_{ij} > 0 \Rightarrow \Delta E_j \leq 0 \quad (6.7)$$

Change from 1 to 0

$$\Delta x_j = -1, \sum_i x_i w_{ij} \leq 0 \Rightarrow \Delta E_j \leq 0 \quad (6.8)$$

6.1.4 Associative Memory

An associative memory can store information, e.g. a binary vector representing the state of a Hopfield network, and retrieve it from only partial information. The number of patterns that can be learned by a Hopfield net is called the capacity C and depends on the number of neurons m and the used training algorithm.

If we want to link additional information to a state vector (e.g. the class of it) we have to store this ourselves.

6.1.5 Training

A simple non-incremental learning rule is:

$$w_{ij} = \sum_{x \in X} (2x_i - 1)(2x_j - 1) \quad (6.9)$$

Assuming uncorrelated patterns, a net trained with this method has can store up to C patterns (capacity). C is bound by

$$\frac{m}{4 \ln(m)} < C < \frac{m}{2 \ln m} \quad (6.10)$$

, where m is the number of neurons.

See [Storkey, 1997] for more on different learning rules and their effect on the capacity.

6.1.6 Limitations

- Found stable state is not guaranteed the most similar pattern to the input pattern
- Not all memories are remembered with same emphasis (attractors regions have different sizes)
- Spurious states can occur
- Efficiency is not good

6.2 Boltzmann Machines

Literature: [Bengio, 2009]

Boltzmann machines are stochastic recurrent neural networks. They work similar *Hopfield nets* in that they have a binary state vector ($s_j \in 0, 1$), which is sometimes divided in visible states (neurons) v_j and hidden neurons h_j . But while the activation calculation in *Hopfield nets* is deterministically, it is stochastic for Boltzmann machine (BM).

$$z_j = b_j + \sum_i v_i w_{ij} \quad (6.11)$$

The probability that the activation of neuron s_j is set to 1 is calculated using the sigmoid function,

$$p(s_j = 1) = \sigma(v_j) = \frac{1}{1 + e^{-z_j}} \quad (6.12)$$

, otherwise the activation of s_j becomes 0.

Further we define the network architecture as follows

- Network is fully connected
- No self connections, $w_{ii} = 0$
- Undirected/symmetric, $w_{ij} = w_{ji}$

6.2.1 Energy-based models

Todo: Derivations of energy-based models, see [Bengio, 2009]

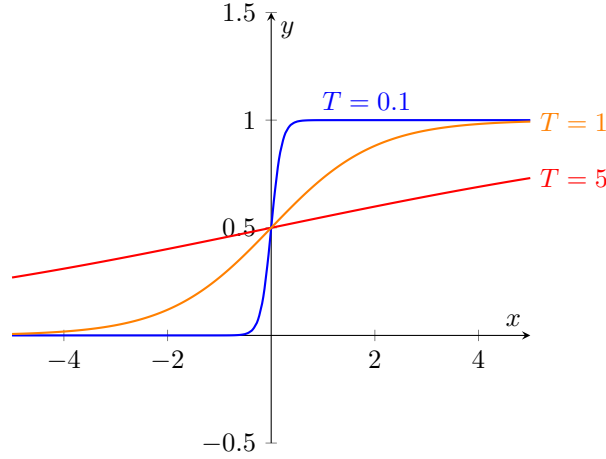


Figure 6.1: Sigmoid function (orange), with high temperature (red) and low temperature (blue).

6.2.2 Energy of Boltzmann Machines

$$E(s) = - \sum_i b_i s_i - \sum_{i,j} s_i w_{ij} s_j \quad (6.13)$$

$$E(v, h) = - \sum_i b_i v_i - \sum_i c_i h_i - \sum_{i,j} v_i w_{ij}^{vh} h_j - \sum_{i,j} v_i w_{ij}^{vv} v_j - \sum_{i,j} h_i w_{ij}^{hh} h_j \quad (6.14)$$

6.2.3 Training

Training is done by calculating the gradient of the log-likelihood derived by the parameter set θ .

$$\frac{\partial \log p(v)}{\partial \theta} \quad (6.15)$$

Todo: Training of BMs

6.2.4 Simulated Annealing

Use temperature T to allow more changes in the beginning (e. g. jumps out of bad local minima) by using the sigmoid function with $\beta = \frac{1}{T}$.

$$p(s_j = 1) = \frac{1}{1 + e^{\frac{-z_j}{T}}} \quad (6.16)$$

Figure 6.1 shows curves with different temperature.

6.2.5 Why to restrict BMs

Unrestricted Boltzmann machines are very powerful and can compute every function. But due the complex net structure the training is very slow and computational expensive. RBMs applying constraints to the network architecture, but can still compute any function.

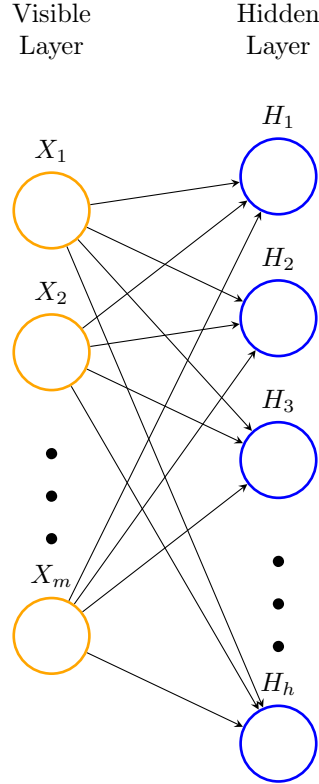


Figure 6.2: Architecture of a RBM

6.3 Restricted Boltzmann Machines

Literature: [Bengio, 2009]

For an efficient training Boltzmann machines can be restricted to have a bipartite graph with one set of visible neurons and one set of hidden neurons. As shown in figure 6.2 there are only visible-hidden and hidden-visible connections (still symmetric). Therefore hidden units h_j only depend on the visible units v_j and vice versa.

$$p(h_j = 1|v) = \sigma(c_j + \sum_i w_{ij}v_i) \quad (6.17)$$

$$p(v_j = 1|h) = \sigma(b_j + \sum_i w_{ij}h_i) \quad (6.18)$$

With b_j as the biases for the visible units and c_j for the hidden units.

6.3.1 Energy

$$E(v, h) = - \sum_{i,j} w_{ij}v_ih_j - \sum_i a_iv_i - \sum_i b_ih_i \quad (6.19)$$

6.4 Training

Literature: [Hinton et al., 2012]

The most common algorithm used is contrastive divergence (CD) used inside a gradient-descent and performing Gibbs Sampling. A single-step contrastive divergence (CD-1) procedure for a single training example can be summarized as follows:

1. Sample hidden units \mathbf{h} from training example \mathbf{v}
2. Sample reconstruction \mathbf{v}' of visible units using \mathbf{h} and then resample \mathbf{h}' from it. (Gibbs sampling step)
3. $w_{ij} \leftarrow w_{ij} - \eta (\mathbf{v}\mathbf{h}^T - \mathbf{v}'\mathbf{h}'^T)$

6.5 Stacking RBMs

Chapter 7

Autoencoders

- Hidden layer can be any size
- Output layer must have size of input layer
- Weight matrix \mathbf{W} and bias \mathbf{b} can be learning using backpropagation
- $\mathbf{W}^{(2)} = (\mathbf{W}^{(1)})^T$
- Non-linear activation function is required, otherwise autoencoder might learn pass through and hidden layer could be similar to principal component analysis (PCA)

7.1 Sparse Autoencoders

Literature: [Ng, 2011]

Favor sparse weight matrices by forcing the average neuron activation $\hat{\rho}_j$ to be close to ρ 0.2.

$$\hat{\rho}_j = \frac{1}{X} \sum_{x \in X} \mathbf{w}_j \mathbf{x} \quad (7.1)$$

This can be done by adding the Kullback–Leibler divergence,

$$KL(\rho || \hat{\rho}_j) = \rho \log \frac{\rho}{\hat{\rho}_j} + (1 - \rho) \log \frac{1 - \rho}{1 - \hat{\rho}_j} \quad (7.2)$$

, to the error function E ,

$$E_{\text{sparse}} = E + \sum_{j \in \text{hidden}} KL(\rho || \hat{\rho}_j) \quad (7.3)$$

Figure 7.2 shows the error introduced by the average neuron activation of a single hidden neuron.

7.2 Denoising Autoencoders

- Create corrupted X' of training data X , e. g. using Gaussian noise or by setting values to zero
- During training use X' as input, but target X as output
- Resulting autoencoder finds good features for noisy input

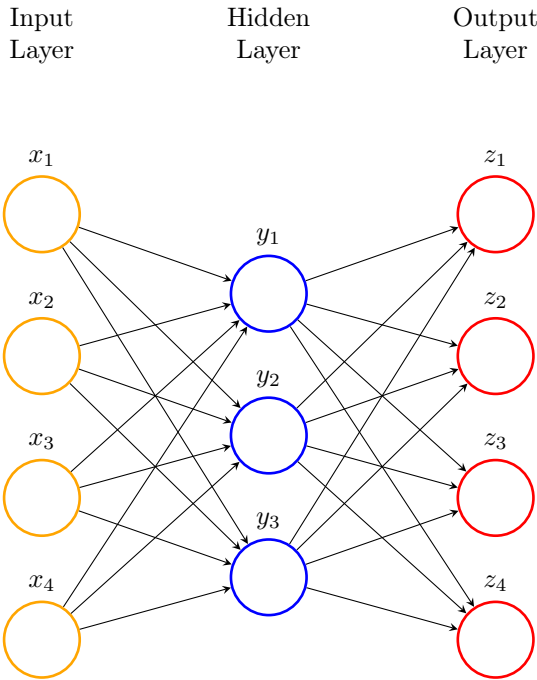


Figure 7.1: Simple Autoencoder

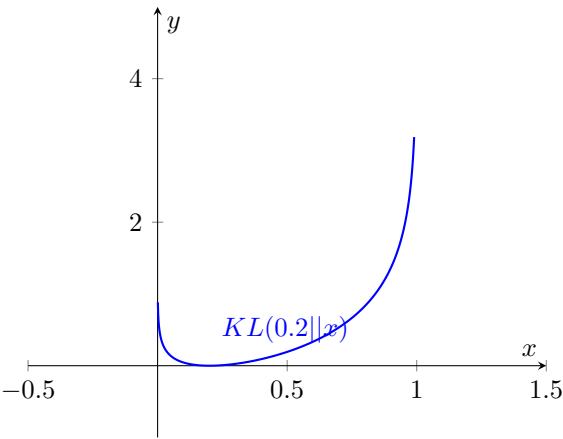


Figure 7.2: Kullback–Leibler divergence $KL(\rho||\hat{\rho}_j)$ with $\rho = 0.2$

7.3 Stacking Autoencoders

- Train first autoencoder AE_1 using the original training data X
- Input X into AE_1 and use output as X'
- Train second autoencoder AE_2 using X'
- Use AE_2 to transform X' to X''
- Continue to create EA_3, EA_4, \dots
- Use weight matrices and biases from autoencoders to stack a deep neural network (DNN)

Chapter 8

Advanced Neural network architectures

8.1 Time Delay Neural Networks (TDNN)

8.2 Multi-State Time Delay Neural Networks (MS-TDNN)

Literature: [[Haffner and Waibel, 1992](#), [Hild and Waibel, 1993](#)]

8.3 Convolutional Neural Networks (CNN)

8.4 Recurrent Neural networks (RNN)

8.4.1 Vanishing Gradient

8.4.2 Backpropagation through Time

Chapter 9

Practical Tricks

9.1 Data Normalization

Numerical needs to be normalized because neural network (NN) function best with inputs between in range $[0, 1]$ or $[-1, +1]$. There are two common normalizations, the min-max normalization (sometimes called feature scaling),

$$\hat{x}_i = \frac{x_i - \min(x)}{\max(x) - \min(x)} \quad (9.1)$$

, which will transform the smallest value to 0 and the biggest to 1 and everything else linearly in between, and gaussian normalization,

$$\hat{x}_i = \frac{x_i - \text{mean}(x)}{\text{std}(x)} = \frac{x_i - \mu}{\sigma} \quad (9.2)$$

, will transform x to have zero mean and a standard deviation of one. Other names for int are standard score or z-scores.

Those are the two most common methods, but depending on the input there might be more.

9.2 Bottleneck-Features

9.3 Weight Decay

Literature: [Connect et al., 1992] and [Duda et al., 2000, Chapter 6.8]

To help generalizing by simplifying a network and to avoid overfitting one can impose a heuristic that the weights stay small. This may not always lead to improved network performance, but experience showed that it helps in many cases and even can be used together with momentum (see 5.4.3).

It is also every simple, after each weight update we simply “decay” the weights with a parameter $\epsilon \in (0, 1)$.

$$w_{ij} \leftarrow w_{ij}(1 - \epsilon) \quad (9.3)$$

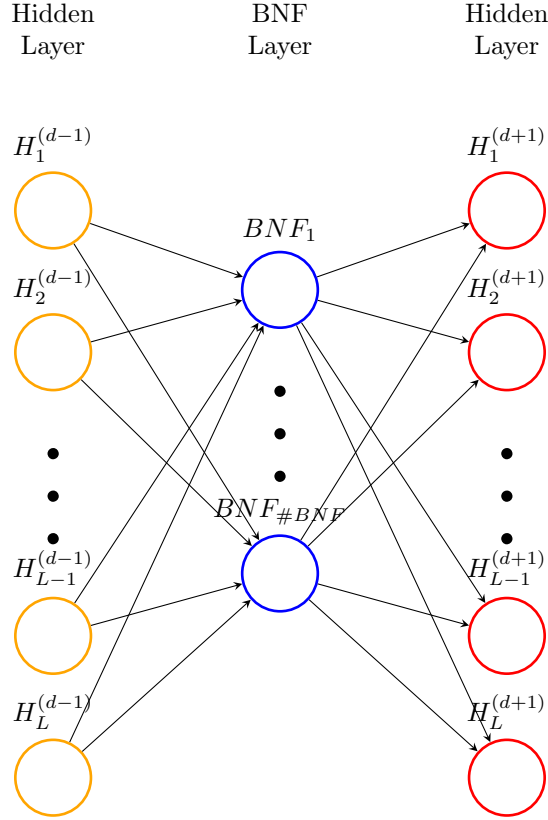


Figure 9.1: Neural network with Bottleneck features

combined with the weight update we get

$$w_{ij} \leftarrow (w_{ij} - \eta \frac{\partial E}{\partial w_{ij}})(1 - \epsilon) \quad (9.4)$$

$$\Delta w_{ij} = -\epsilon w_{ij} - (1 - \epsilon) \eta \frac{\partial E}{\partial w_{ij}} \quad (9.5)$$

$$= -\epsilon w_{ij} - \eta' \frac{\partial E}{\partial w_{ij}} \quad (9.6)$$

$$= -\epsilon w_{ij} - \eta' \frac{\partial E}{\partial w_{ij}} \quad (9.7)$$

Those weights that are need to solve the problem will increase enough with the weight update to compensate the small decrease, but others will get smaller and smaller until they can be eliminated completely. It can be shown that the above is equivalent to gradient descent we a modified error function:

$$E_{\text{weight decay}} = E + \frac{2\epsilon}{\eta} \mathbf{w}^T \mathbf{w} \quad (9.8)$$

9.4 Dropout

Literature: [Hinton et al., 2012, Hinton, 2014, Wang and Manning, 2013]

- Randomly drop out hidden units (set to zero) and input features during training.
- Each time a different model is trained and the final model is the mean of all trained models.
- Prevents feature co-adaptation detectors and overfitting.

Bibliography

- [Bengio, 2009] Bengio, Y. (2009). Learning deep architectures for ai.
- [Connect et al., 1992] Connect, a. K., Krogh, a., and Hertz, J. a. (1992). A simple weight decay can improve generalization. *Advances in Neural Information Processing Systems*, 4:950–957.
- [Duchi et al., 2011] Duchi, J., Hazan, E., and Singer, Y. (2011). Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12:2121–2159.
- [Duda et al., 2000] Duda, R. O., Hart, P. E., and Stork, D. G. (2000). *Pattern Classification*. JOHN WILEY & SONS, INC.
- [Dyer,] Dyer, C. Notes on adagrad. *Ark.Cs.Cmu.Edu*, pages 1–3.
- [Goodfellow et al., 2013] Goodfellow, I. J., Warde-farley, D., and Courville, A. (2013). Maxout networks.
- [Haffner and Waibel, 1992] Haffner, P. and Waibel, A. (1992). Multi-state time delay networks for continuous speech recognition. In *Advances in Neural Information Processing Systems 4*, pages 135–142. Morgan-Kaufmann.
- [Hampshire and Waibel, 1990] Hampshire, J. B. and Waibel, A. H. (1990). Novel objective function for improved phoneme recognition using time-delay neural networks. *IEEE Transactions on Neural Networks*, 1(2):216–228.
- [Hild and Waibel, 1993] Hild, H. and Waibel, A. (1993). Connected letter recognition with a multi-state time delay neural network. In *In 3rd European Conference on Speech, Communication and Technology (EUROSPEECH) 93*.
- [Hinton, 2014] Hinton, G. (2014). Dropout : A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research (JMLR)*, 15:1929–1958.
- [Hinton et al., 2012] Hinton, G. E., Srivastava, N., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. R. (2012). Improving neural networks by preventing co-adaptation of feature detectors. pages 1–18.
- [Introduction, 2000] Introduction, I. (2000). A fuzzy clustering model of data and fuzzy. pages 302–307.
- [Kohonen, 1990] Kohonen, T. (1990). Improved versions of learning vector quantization. *International Joint Conference on Neural Networks*.
- [Kohonen, 2001] Kohonen, T. (2001). *Self-Organizing Maps*, volume 30 of *Springer Series in Information Sciences*. Springer.
- [Linde et al., 1980] Linde, Y., Buzo, A., and Gray, R. (1980). An algorithm for vector quantizer design. *IEEE Transactions on Communications*, 28(1):84–95.
- [Mitchell, 1997] Mitchell, T. M. (1997). *Machine Learning*, volume 4 of *McGraw-Hill Series in Computer Science*. McGraw-Hill.
- [Ng, 2011] Ng, A. (2011). Sparse autoencoder. *CS294A Lecture notes*, pages 1–19.

- [Patterson, 1997] Patterson, D. W. (1997). *Kuenstliche neuronale Netze - Das Lehrbuch*. Prentice Hall, Muenchen [u.a.].
- [Riedmiller and Braun, 1994] Riedmiller, M. and Braun, H. (1994). Rprop - description and implementation details. *Neural Networks*, (January):1–2.
- [Rumelhart et al., 1986] Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1986). Learning internal representations by error propagation. In *Parallel Distributed Processing: Explorations in the Microstructure of Cognition*, volume 1, pages 318–362.
- [Storkey, 1997] Storkey, A. (1997). Increasing the capacity of a hopfield network without sacrificing functionality. In *Artificial Neural Networks-ICANN'97*, pages 451–456. Springer.
- [Wang and Manning, 2013] Wang, S. I. and Manning, C. D. (2013). Fast dropout training. *Proceedings of the 30th International Conference on Machine Learning*, 28:118–126.
- [Yu and Chen, 1997] Yu, X. H. and Chen, G. A. (1997). Efficient backpropagation learning using optimal learning rate and momentum. *Neural Networks*, 10(3):517–527.
- [Zeiler, 2012] Zeiler, M. D. (2012). Adadelta: An adaptive learning rate method. page 6.

Acronyms

CE cross-entropy. 13

DNN deep neural network. 9, 19

KIT Karlsruhe Institute of Technology. 3

LVQ learning vector quantization. 8

MSE mean-square error. 13

NN neural network. 19

Todo list

Todo: Link some generally good sources (demos, tutorials, online courses, books).	5
Todo: List with of good libraries for Machines Learning and Deep Neural Networks . . .	5
Todo: Add derivate for this equation and why N_k is left out later	7
Todo: Explain advantages and disadvantages, see http://www.chioka.in/differences-between-l1-and-l2-as-1	
Todo: Why is tanh worth the computational overhead over sigmoid?	21
Todo: Derivations of energy-based models, see [Bengio, 2009]	27
Todo: Training of BMs	28