

Neuronal Networks

Formulae Collection

Marvin Ritter
marvin.ritter@gmail.com

20th 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 |
| 2 | Unsupervised Vector Quantization | 6 |
| 2.1 | Loss Functions for Clustering | 6 |
| 2.2 | k-Means | 7 |
| 2.3 | Fuzzy k-Means | 7 |
| 2.4 | LBG-Algorithm | 8 |
| 2.5 | Self-Organizing Maps | 8 |
| 3 | Supervised Vector Quantization | 9 |
| 3.1 | Distortion Measures | 9 |
| 3.2 | Learning Vector Quantization | 10 |
| 3.2.1 | LVQ1 | 10 |
| 3.2.2 | OLVQ | 11 |
| 3.2.3 | LVQ2.1 | 11 |
| 3.2.4 | LVQ3 | 11 |
| 4 | Multi-Layer Perceptrons | 12 |
| 4.1 | Backpropagation | 12 |
| 4.1.1 | Notation | 12 |
| 4.1.2 | Weight Updates in Output Layer | 12 |
| 4.1.3 | Weight Outputs for Hidden Layers | 13 |
| 4.1.4 | Backpropagation using Matrix Notation | 14 |
| 4.2 | Error Functions | 15 |
| 4.2.1 | Mean-Squared Error | 15 |
| 4.2.2 | Cross Entropy Error | 15 |
| 4.3 | Activation Functions | 16 |
| 4.3.1 | Step Function | 16 |
| 4.3.2 | Sigmoid | 16 |
| 4.3.3 | Softmax Function | 17 |
| 4.3.4 | Hyperbolic Tangent Function | 18 |
| 4.3.5 | Linear Function | 18 |
| 4.3.6 | Rectified Linear Unit | 18 |
| 4.3.7 | Softplus | 20 |
| 4.3.8 | Maxout | 20 |
| 4.4 | Learning Rate | 20 |
| 4.4.1 | Exponential Decaying Learning Rate | 21 |
| 4.4.2 | Performance Scheduling | 21 |
| 4.4.3 | Momentum or Rprop | 21 |
| 4.4.4 | AdaGrad | 21 |
| 4.4.5 | Newbob Scheduling | 22 |
| 5 | Boltzmann Machines | 23 |
| 5.1 | Hopfield Nets | 23 |

| | | |
|----------|-------------------------------|-----------|
| 5.1.1 | Update Procedure | 23 |
| 5.1.2 | Energy function | 24 |
| 5.1.3 | Convergence | 24 |
| 5.1.4 | Associative Memory | 24 |
| 5.1.5 | Training | 25 |
| 5.1.6 | Limitations | 25 |
| 5.2 | Boltzmann Machines | 25 |
| 5.2.1 | Energy | 26 |
| 5.2.2 | Simulated Annealing | 26 |
| 5.2.3 | Why to restrict BMs | 26 |
| 5.3 | Restricted Boltzmann Machines | 26 |
| 5.4 | Stacking RBMs | 27 |
| 6 | Autoencoders | 28 |
| 6.1 | Spares Autoencoders | 28 |
| 6.2 | Denoising Autoencoders | 28 |
| 6.3 | Stacking Autoencoders | 28 |
| 7 | Practical Tricks | 29 |
| 7.1 | Data Normalization | 29 |
| 7.2 | Bottleneck-Features | 29 |
| 7.3 | Weight Decay | 29 |
| | Acronyms | 32 |

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 4.4) |
| μ_j | Centroid of cluster j |
| $d(x, y)$ | Dissimilarity measure, distance or distortion measure, if mention otherwise we use the squared Euclidean distance (see 3.1) |
| E, J | Error function (see 4.2), 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 4.3) |
| $\sigma(x)$ | Sigmoid activation function (see 4.3.2) |
| $z_j^{(l)}$ | 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 $z_j^{(l)}$, $a_j = \phi(z_j)$ |
| δ_{ij} | Kronecker delta, which is defined as $\delta_{ij} = \begin{cases} 1 & \text{if } i = j, \\ 0 & \text{if } i \neq j. \end{cases}$ |

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

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

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

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 $W(C)$

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

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

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

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

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

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

3. **Repeat step 1 and 2 until stopping criteria fulfilled (e. g. centers stop moving/ C(i) unchanged).**

During the course of the k-means algorithm, the loss function monotonically decreases (in both steps) into a local minimum.

2.3 Fuzzy k-Means

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

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

TODO: Do we need $\sum_{k=1}^K u_{ik} = 1$?

Our new loss function is:

$$J_m = \sum_{i=1}^N \sum_{k=1}^K u_{ik} d(x_i, \mu_k) \quad (2.3)$$

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

Original Paper: [\[Linde et al., 1980\]](#)

Demo: [external website](#)

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

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

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

$$\mu_1(0) = \frac{1}{N} \sum_{i=1}^N x_i \quad (2.4)$$

$$t = 0 \quad (2.5)$$

2. **Double codebook:**

$$\mu_j(t+1) = (1 + \epsilon) \mu_j(t) \quad \text{for } j = 1, \dots, 2^t \quad (2.6)$$

$$\mu_{2*j+1}(t+1) = (1 - \epsilon) \mu_j(t) \quad \text{for } j = 1, \dots, 2^t \quad (2.7)$$

$$t = t + 1 \quad (2.8)$$

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 $\mu_j(\log_2(K))$.

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

Multi-Layer Perceptrons

4.1 Backpropagation

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

4.1.1 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 $o_j := a_j^{(L)}$.

4.1.2 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:

$$\begin{aligned}
\frac{\partial E}{\partial w_{ij}^{(L)}} &= \frac{\partial E}{\partial z_j^{(L)}} \frac{\partial z_j^{(L)}}{\partial w_{ij}^{(L)}} = \frac{\partial E}{\partial a_j^{(L)}} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \frac{\partial z_j^{(L)}}{\partial w_{ij}^{(L)}} && \text{(chain rule)} \\
\frac{\partial z_j^{(L)}}{\partial w_{ij}^{(L)}} &= \frac{\partial \sum_i w_{ij}^{(L)} * x_i^{(L-1)}}{\partial w_{ij}^{(L)}} = a_i^{(L-1)} && (4.1) \\
\frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} &= \frac{\partial \phi(z_j^{(L)})}{\partial z_j^{(L)}} = (\phi(z_j^{(L)})(1 - \phi(z_j^{(L)}))) = o_j(1 - o_j) && \text{(sigmoid derivative, see 4.3.2)} \\
\frac{\partial E}{\partial a_j^{(L)}} &= \frac{\partial E}{\partial o_j} = \frac{\partial \frac{1}{2} \sum_{i=1}^k (t_i - o_i)^2}{\partial o_j} = (o_j - t_j) && \text{(MSE derivative, see 4.2.1)}
\end{aligned}$$

putting everything together

$$\frac{\partial E}{\partial w_{ij}^{(L)}} = (o_j - t_j) o_j (1 - o_j) a_i^{(L-1)} \quad (4.2)$$

We can now define the error for neuron j as

$$\begin{aligned}
\delta_j^{(L)} &:= \frac{\partial E}{\partial z_j^{(L)}} \\
&= \frac{\partial E}{\partial a_j^{(L)}} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \\
&= (o_j - t_j) o_j (1 - o_j)
\end{aligned} \quad (4.3)$$

, and our weight update

$$\Delta w_{ij}^{(L)} := -\eta \frac{\partial E}{\partial w_{ij}^{(L)}} = -\eta \delta_j^{(L)} a_i^{(L-1)} \quad (4.4)$$

$$w_{ij}^{(L)} \leftarrow w_{ij}^{(L)} + \Delta w_{ij}^{(L)} \quad (4.5)$$

4.1.3 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.

$$\delta_j^{(l)} := \frac{\partial E}{\partial z_j^{(l)}} = \frac{\partial E}{\partial a_j^{(l)}} \frac{\partial a_j^{(l)}}{\partial z_j^{(l)}} \quad (4.6)$$

$$\frac{\partial E}{\partial a_j^{(l)}} = \sum_k \frac{\partial E}{\partial z_k^{(l+1)}} \frac{\partial z_k^{(l+1)}}{\partial a_j^{(l)}} \quad (4.7)$$

$$= \sum_k \frac{\partial E}{\partial z_k^{(l+1)}} \frac{\partial \sum_i w_{ik}^{(l+1)} * a_i^{(l)}}{\partial a_j^{(l)}} \quad (4.8)$$

$$= \sum_k \frac{\partial E}{\partial z_k^{(l+1)}} w_{jk}^{(l+1)} \quad (4.9)$$

$$= \sum_k \delta_k^{(l+1)} w_{jk}^{(l+1)} \quad (4.10)$$

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

$$\delta_j^{(l)} = a_j^{(l)}(1 - a_j^{(l)}) \sum_k \delta_k^{(l+1)} w_{jk}^{(l+1)} \quad (4.11)$$

This leads to our weight updates for hidden neurons.

$$\Delta w_{ij}^{(l)} = -\eta \frac{\partial E}{\partial w_{ij}^{(l)}} \quad (4.12)$$

$$= -\eta \frac{\partial E}{\partial z_j^{(l)}} \frac{\partial z_j^{(l)}}{\partial w_{ij}^{(l)}} \quad (4.13)$$

$$= -\eta \delta_j^{(l)} a_i^{(l-1)} \quad (4.14)$$

$$w_{ij}^{(l)} \leftarrow w_{ij}^{(l)} + \Delta w_{ij}^{(l)} \quad (4.15)$$

4.1.4 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

$$A^{(0)} = X \quad (4.16)$$

$$Z^{(l)} = A^{(l-1)} * W^{(l)} \quad (4.17)$$

$$A^{(l)} = \phi(Z^{(l)}) \quad (4.18)$$

$$O = A^{(L)} \quad (4.19)$$

$$E = \frac{1}{2} \sum (T - O) \circ (T - O) \quad (4.20)$$

$$\delta^{(L)} = (O - T) \circ O \circ (1 - O) \quad (4.21)$$

$$\delta^{(l)} = A^{(l)}(1 - A^{(l)})W^{(l+1)}\delta^{(l+1)}\Delta W^{(l)} = -\eta \delta^{(l)} A^{(l-1)} \quad (4.22)$$

$$\hat{W}^{(l)} = W^{(l)} + \Delta W^{(l)} \quad (4.23)$$

4.2 Error Functions

Sometimes also called objective functions or loss-functions.

4.2.1 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.

$$E_{\text{MSE}} = \frac{1}{n} \sum_{i=1}^n (t_i - o_i)^2 \quad (4.24)$$

This would derive to

$$\begin{aligned} \frac{\partial E_{\text{MSE}}}{\partial o_j} &= \frac{1}{n} \sum_{i=1}^n \frac{\partial (t_i - o_i)^2}{\partial o_j} \\ &= \frac{1}{n} \frac{\partial (t_j - o_j)^2}{\partial o_j} \\ &= \frac{1}{n} 2(t_j - o_j)(-1) \\ &= \frac{2}{n} (o_j - t_j) \end{aligned} \quad (4.25)$$

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

$$E_{\text{MSE}} = \frac{1}{2} \sum_{i=1}^k (t_i - o_i)^2 \quad (4.26)$$

$$\frac{\partial E_{\text{MSE}}}{\partial o_j} = -(t_j - o_j) = o_j - t_j \quad (4.27)$$

4.2.2 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.

$$E_{\text{CE}} = - \sum_{i=1}^k (t_i \log(o_i) + (1 - t_i) \log(1 - o_i)) \quad (4.28)$$

This will derive to:

$$\begin{aligned} \frac{\partial E_{\text{CE}}}{\partial o_j} &= - \frac{\partial}{\partial o_j} \sum_{i=1}^k (t_i \log(o_i) + (1 - t_i) \log(1 - o_i)) \\ &= - \frac{\partial}{\partial o_j} t_j \log(o_j) - \frac{\partial}{\partial o_j} (1 - t_j) \log(1 - o_j) \\ &= - \frac{t_j}{o_j} + \frac{1 - t_j}{1 - o_j} \\ &= \frac{1 - t_j}{1 - o_j} - \frac{t_j}{o_j} \end{aligned} \quad (4.29)$$

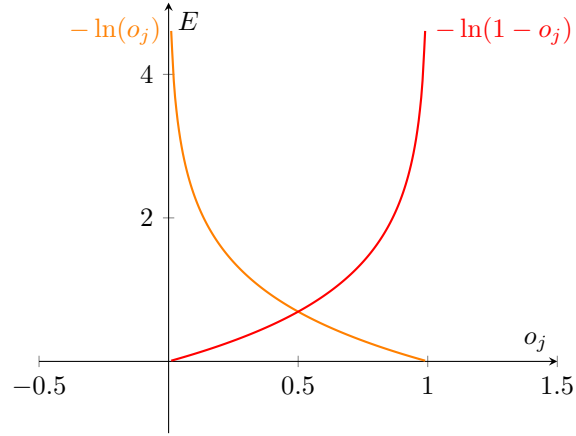


Figure 4.1: Error produced for $t_j = 1$ (orange) and for $t_j = 0$ (red). o_j must be in range $(0, 1]$, otherwise the logarithm can be complex or undefined or the error gets rediciously large.

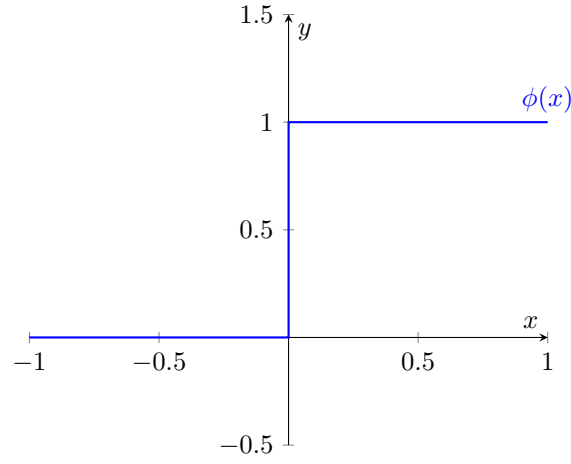


Figure 4.2: Step Function

4.3 Activation Functions

4.3.1 Step Function

$$\phi(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases} \quad (4.30)$$

The derivative is always 0.

4.3.2 Sigmoid

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

$$\sigma(x) = \frac{1}{1 + e^{-\beta x}} \quad (4.31)$$

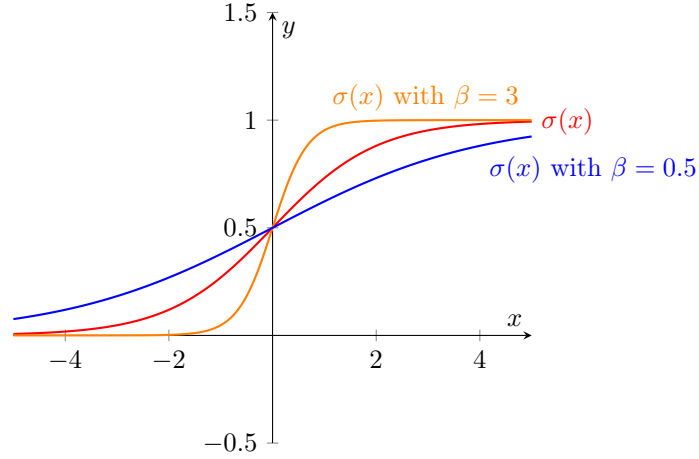


Figure 4.3: Sigmoid function with $\beta = 1$ (red), $\beta = 3$ (orange) and $\beta = 0.5$ (blue).

And it's derivative:

$$\begin{aligned}
 \frac{d\sigma(x)}{dx} &= \frac{d}{dx}(1 + e^{-\beta x})^{-1} \\
 &= (-1)(1 + e^{-\beta x})^{-2}(-\beta e^{-\beta x}) \\
 &= \beta(1 + e^{-\beta x})^{-2}e^{-\beta x} \\
 &= \beta \frac{1}{1 + e^{-\beta x}} \frac{e^{-\beta x}}{1 + e^{-\beta x}} \\
 &= \beta \sigma(x) \frac{e^{-\beta x} + 1 - 1}{1 + e^{-\beta x}} \\
 &= \beta \sigma(x) \left(1 - \frac{1}{1 + e^{-\beta x}}\right) \\
 &= \beta \sigma(x)(1 - \sigma(x))
 \end{aligned} \tag{4.32}$$

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

$$\begin{aligned}
 \frac{d\sigma(x)(1 - \sigma(x))}{dx} &= (1 - \sigma(x)) \frac{d\sigma(x)}{dx} + \sigma(x) \frac{d(1 - \sigma(x))}{dx} \\
 &= (1 - \sigma(x))\sigma(x)(1 - \sigma(x)) + \sigma(x)(-1)\sigma(x)(1 - \sigma(x)) \\
 &= \sigma(x)(1 - \sigma(x))^2 - \sigma(x)^2(1 - \sigma(x))
 \end{aligned} \tag{4.33}$$

4.3.3 Softmax Function

In a classification problem we would like a_j to be a probability ($\Rightarrow \sum_i a_i = 1$).

$$a_i = \phi(z_i) = \frac{e^{z_i}}{\sum_j e^{z_j}} \tag{4.34}$$

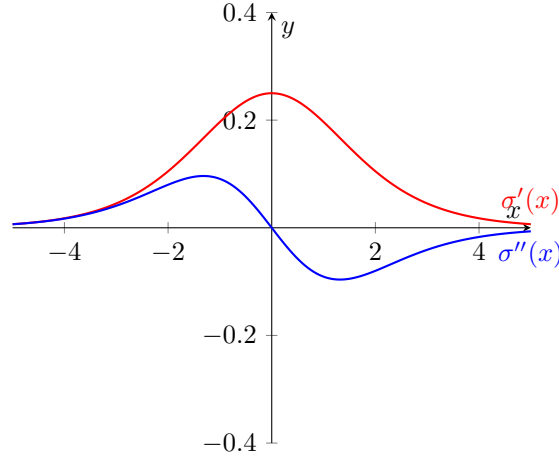


Figure 4.4: First and second derivative of the sigmoid function.

And it's derivative:

$$\begin{aligned}
 \frac{\partial \phi(z_i)}{\partial z_i} &= \frac{\left(\frac{\partial}{\partial z_i} e^{z_i}\right) (\sum_j e^{z_j}) - e^{z_i} \left(\frac{\partial}{\partial z_i} (\sum_j e^{z_j})\right)}{(\sum_j e^{y_j})^2} \\
 &= \frac{e^{y_i} \sum_j e^{y_j} - e^{y_i} e^{y_i}}{(\sum_j e^{y_j})^2} \\
 &= \frac{e^{y_i} \sum_j e^{y_j}}{(\sum_j e^{y_j})^2} - \frac{e^{y_i} e^{y_i}}{(\sum_j e^{y_j})^2} \\
 &= \frac{e^{y_i}}{\sum_j e^{y_j}} - \left(\frac{e^{y_i}}{\sum_j e^{y_j}}\right)^2 \\
 &= \phi(y_i) - \phi(y_i)^2 \\
 &= \phi(y_i)(1 - \phi(y_i))
 \end{aligned} \tag{4.35}$$

4.3.4 Hyperbolic Tangent Function

Similar to Sigmoid function, but if the input has a mean of 0 then so will the output.

$$\phi(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \tag{4.36}$$

And it's derivative:

$$\frac{d\phi(x)}{dx} = 1 - \tanh(x) = \frac{(e^x - e^{-x})^2}{(e^x + e^{-x})^2} \tag{4.37}$$

4.3.5 Linear Function

$$\phi(x) = x \tag{4.38}$$

4.3.6 Rectified Linear Unit

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

$$\phi(x) = \max(0, x) \tag{4.39}$$

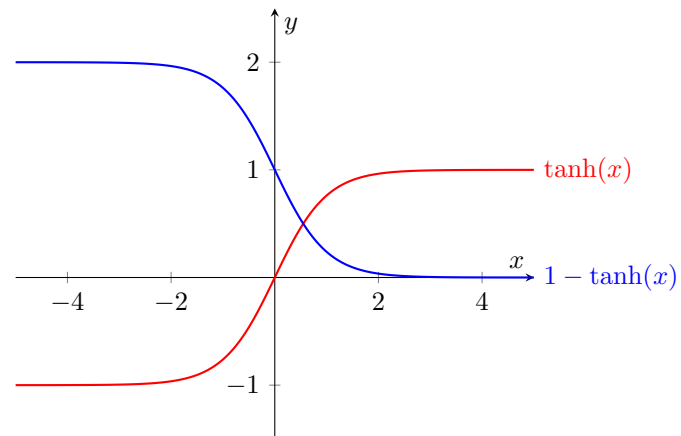


Figure 4.5: Hyperbolic Tangent Function (red) and its derivative (blue).

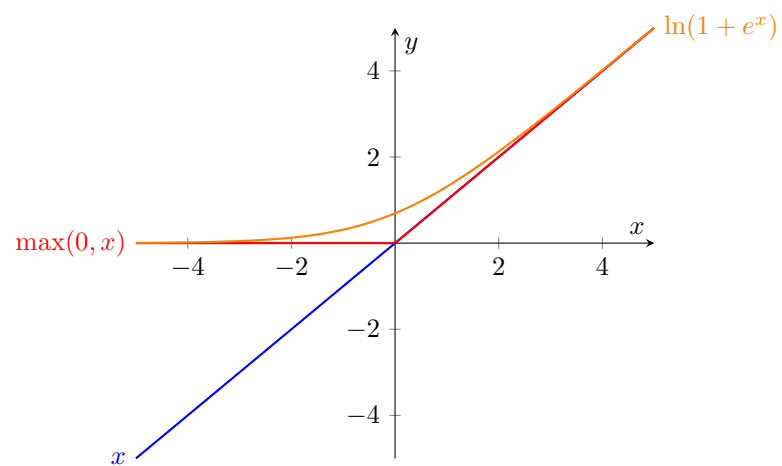


Figure 4.6: Linear Function (blue), Rectified Linear Unit (red), Softplus (orange)

And it's derivative:

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

4.3.7 Softplus

Smoothed version of ReLU.

$$\phi(x) = \ln(1 + e^x) \quad (4.40)$$

And it's derivative:

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

4.3.8 Maxout

Literature: [Goodfellow et al., 2013]

Outputs the maximum of its inputs

$$a_j^{(l)} = \max_i z_i^{(l)}, \quad (j-1)g + 1 \leq i \leq jg \quad (4.41)$$

$$z_i^{(l)} = \sum_k x_k^{(l-1)} w_{ki}^{(l)} + b^{(l)} \quad (4.42)$$

4.4 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)

4.4.1 Exponential Decaying Learning Rate

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

$$\eta(t) = \alpha \eta(t-1) = \alpha^t \eta(0) \quad (4.43)$$

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

4.4.2 Performance Scheduling

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

4.4.3 Momentum or Rprop

Literature: [Riedmiller and Braun, 1994]

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,

$$w_{ij} \leftarrow w_{ij} + v_{ij}(t) \quad (4.44)$$

$$v_{ij}(t) := \alpha v_{ij}(t-1) - \eta \delta_j a_i^{(l-1)} \quad (4.45)$$

, with α as momentum factor.

Todo: Verify that *Rprop* and *Momentum* are really the same or add differences

4.4.4 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.

$$g_{ij} = \frac{\partial E}{\partial w_{ij}} \quad (4.46)$$

$$\eta_{ij}(t) = \frac{\eta_0}{\sqrt{\sum_{\tau=1}^t g_{ij}(\tau)^2}} g_{ij}(t) \quad (4.47)$$

$$w_{ij}(t) \leftarrow w_{ij}(t-1) - \eta_{ij}(t) g_{ij}(t) \quad (4.48)$$

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 .

$$\eta_{ij}(t) = \frac{\eta_0}{\sqrt{\sum_{\tau=1}^w g_{ij}(t-\tau)^2}} g_{ij}(t) \quad (4.49)$$

AdaDelta contains further improvements.

4.4.5 Newbob Scheduling

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

- Start with constant learning rate $\eta(0)$
- Once the cross-validation error stops decreasing switch to exponential decaying learning rate
- End training after cross-validation error stops decreasing again

Chapter 5

Boltzmann Machines

5.1 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.

5.1.1 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

$$\phi(z) = \begin{cases} 1 & \text{if } z > 0 \\ 0 & \text{if } z \leq 0 \end{cases} \quad (5.1)$$

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 that the state vector (e.g. class of the pattern) we have to store these ourself during training and use x to look them up.

5.1.2 Energy function

We define the energy in our Hopfield net as:

$$E = -\frac{1}{2} \sum_i \sum_{j \neq i} x_i x_j w_{ij} \quad (5.2)$$

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.

5.1.3 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 changed, compute the energy change

$$E_j = -\frac{1}{2} \sum_i x_i x_j w_{ij} \quad (5.3)$$

$$\Delta E = E_j(t+1) - E_j(t) \quad (5.4)$$

$$= -\frac{1}{2} \left[x_j(t+1) \sum_i x_i w_{ij} - x_j(t) \sum_i x_i w_{ij} \right] \quad (5.5)$$

$$= -\frac{1}{2} \Delta x_j \sum_i x_i w_{ij} \quad \text{with} \quad \Delta x_j = x_j(t+1) - x_j(t) \quad (5.6)$$

Change from 0 to 1

$$\Delta x_j = 1, \sum_i x_i w_{ij} > 0 \Rightarrow \Delta E_j \leq 0 \quad (5.7)$$

Change from 1 to 0

$$\Delta x_j = -1, \sum_i x_i w_{ij} \leq 0 \Rightarrow \Delta E_j \leq 0 \quad (5.8)$$

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

5.1.5 Training

A simple non-incremental learning rule is:

$$w_{ij} = \sum_{x \in X} (2x_i - 1)(2x_j - 1) \quad (5.9)$$

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

$$\frac{m}{4 \ln(m)} < C < \frac{m}{2 \ln m} \quad (5.10)$$

, where m is the number of neurons.

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

5.1.6 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

5.2 Boltzmann Machines

Boltzmann machines are stochastic recurrent neural networks.

Boltzmann machines work similar to Hopfield nets. They have a binary state vector ($x_j \in 0, 1$), but the new state a neuron x_j is determined stochastically (based on the total input z_j) on each update of x_j .

$$z_j = b_j + \sum_i x_i w_{ij} \quad (5.11)$$

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

$$p(x_j = 1) = \sigma(z_j) = \frac{1}{1 + e^{-z_j}} \quad (5.12)$$

, otherwise the activation of x_j becomes 0.

In general Boltzmann machines are unrestricted (see 5.3 for restricted Boltzmann machine (RBM)), but have the properties:

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

One can define input neurons that have the fixed input value during evaluation.

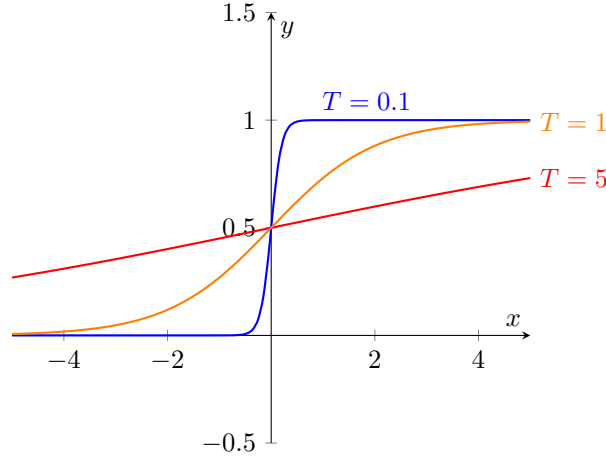


Figure 5.1: Sigmoid function (orange), with high temperature (red) and low temperature (blue).

5.2.1 Energy

The energy of a state x is defined by

$$E(x) = - \sum_i x_i b_i - \frac{1}{2} \sum_{i < j} x_i x_j w_{ij} \quad (5.13)$$

With this we can calculate the probability of a state x

$$p(x) = \frac{e^{-E(x)}}{\sum_y e^{-E(y)}} \quad (5.14)$$

In average updating the neurons will decrease the energy in the network.

5.2.2 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 $\beta = \frac{1}{T}$.

$$p(x_j = 1) = \frac{1}{1 + e^{\frac{-z_j}{T}}} \quad (5.15)$$

Figure 5.1 shows curves with different temperature.

5.2.3 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.

5.3 Restricted Boltzmann Machines

For an efficient training a Boltzmann machines can be restricted to have a bipartite graph with one set of visible neurons and on set of hidden neurons. An shown in figure 5.2 there are only visible-hidden and hidden-visible connection (still symmetric).

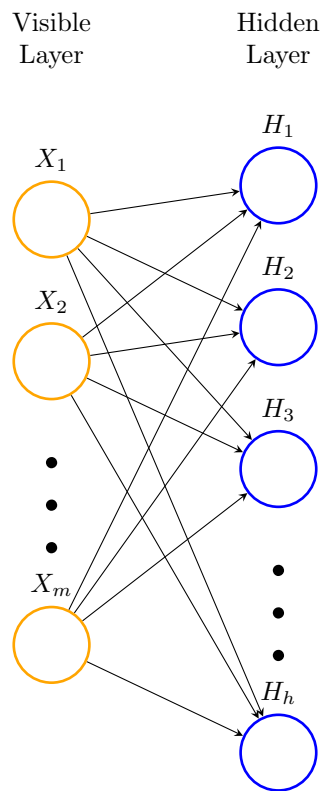


Figure 5.2: Architecture of a RBM

5.4 Stacking RBMs

Chapter 6

Autoencoders

Non-linear activation function required (otherwise autoencoder might learn pass through and hidden layer could be similar to principal component analysis (PCA))

Autoencoders sometimes don't generalize well.

Always tied weights? $W^{(1)} = (W^{(1)})^T$

6.1 Sparse Autoencoders

Literature: [\[Ng, 2011\]](#)

$$KL(\rho || \hat{\rho}_j) = \rho \log \frac{\rho}{\hat{\rho}_j} + (1 - \rho) \log \frac{1 - \rho}{1 - \hat{\rho}_j} \quad (6.1)$$

6.2 Denoising Autoencoders

train with corrupted data X' (gaussian noise)

6.3 Stacking Autoencoders

train first layer with X

train second layer with $X' =$ hidden representation of X

...

Chapter 7

Practical Tricks

7.1 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),

$$\hat{x}_i = \frac{x_i - \min(x)}{\max(x) - \min(x)} \quad (7.1)$$

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

$$\hat{x}_i = \frac{x_i - \text{mean}(x)}{\text{std}(x)} = \frac{x_i - \mu}{\sigma} \quad (7.2)$$

, 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.

7.2 Bottleneck-Features

7.3 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 4.4.3).

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

$$w_{ij} \leftarrow w_{ij}(1 - \epsilon) \quad (7.3)$$

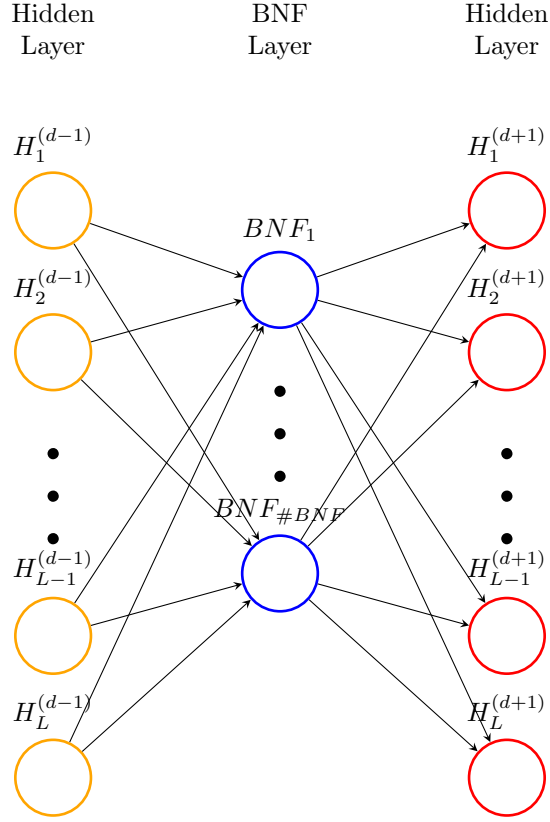


Figure 7.1: Neural network with Bottleneck features

combined with the weight update we get

$$w_{ij} \leftarrow (w_{ij} - \eta \frac{\partial E}{\partial w_{ij}})(1 - \epsilon) \quad (7.4)$$

$$\Delta w_{ij} = -\epsilon w_{ij} - (1 - \epsilon) \eta \frac{\partial E}{\partial w_{ij}} \quad (7.5)$$

$$= -\epsilon w_{ij} - \eta' \frac{\partial E}{\partial w_{ij}} \quad (7.6)$$

$$= -\epsilon w_{ij} - \eta' \frac{\partial E}{\partial w_{ij}} \quad (7.7)$$

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:

$$E_{\text{weight decay}} = E + \frac{2\epsilon}{\eta} \mathbf{w}^T \mathbf{w} \quad (7.8)$$

Bibliography

- [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.
- [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.
- [Zeiler, 2012] Zeiler, M. D. (2012). Adadelata: 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: Verify that *Rprop* and *Momentum* are really the same or add differences 21