

Neuronal Networks

Formulae Collection

Marvin Ritter
marvin.ritter@gmail.com

17th September 2015

Introduction

This is my personal collection of formulae in the field of neural networks. I started it in preparation to my exam on the neural nets at Karlsruhe Institute of Technology (KIT). Even so the content will be similar to the course it is neither limited to it nor linked with the course in any means. *This is not official script and may contain errors and lag completeness.*

Corrections, supplements (or wishes) and links to good sources/ papers are very welcome. Just mail me to marvin.ritter@gmail.com.

Contents

1	Unsupervised Vector Quantization	3
1.1	Loss Functions for Clustering	3
1.2	k-Means	4
1.3	Fuzzy k-Means	4
1.4	LBG-Algorithm	5
2	Supervised Vector Quantization	6
2.1	Distortion Measures	6
2.2	Learning Vector Quantization	7
2.2.1	LVQ1	7
2.2.2	OLVQ	8
2.2.3	LVQ2.1	8
2.2.4	LVQ3	9
2.3	Self-Organizing Maps	9
3	Multi-Layer Perceptrons	10
3.1	Backpropagation	10
3.1.1	Notation	10
3.1.2	Weight Updates in Output Layer	11
3.1.3	Weight Outputs for Hidden Layers	11
3.1.4	Backpropagation using Matrix Notation	12
3.2	Error Functions	13
3.2.1	Mean-Squared Error	13
3.2.2	Cross Entropy Error	14
3.3	Activation Functions	14
3.3.1	Step Function	14
3.3.2	Sigmoid	15
3.3.3	Softmax Function	15
3.3.4	Hyperbolic Tangent Function	16
3.3.5	Linear Function	16
3.3.6	Rectified Linear Unit	17
3.3.7	Softplus	17
3.3.8	Maxout	17
4	Preprocessing	19
4.1	Data Normalization	19
	Acronyms	21

Chapter 1

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 (1.3). 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.

1.1 Loss Functions for Clustering

- Intra-class scatter: $W(C) = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(j)=k} d(x_i, x_j)$

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

- 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 $B(C)$

1.2 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

$$J = \sum_{k=1}^K \sum_{C(i)=k} d(x_i, \mu_k) \quad (1.1)$$

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 $W(C)$

$$W(C) = \frac{1}{2} \sum_{k=1}^K \sum_{C(i)=k} \sum_{C(j)=k} \|x_i - x_j\|^2 = \sum_{k=1}^K N_k \sum_{C(i)=k} \|x_i - \mu_k\|^2 \quad (1.2)$$

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.

1.3 Fuzzy k-Means

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:

$$J_m = \sum_{i=1}^N \sum_{k=1}^K u_{ik} d(x_i, \mu_k) \quad (1.3)$$

1. **Compute degree of membership:**

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

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.**

1.4 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:

$$\mu_1(0) = \frac{1}{N} \sum_{i=1}^N x_i \quad (1.4)$$

$$t = 0 \quad (1.5)$$

2. **Double codebook:**

$$\mu_j(t+1) = (1 + \epsilon) \mu_j(t) \quad \text{for } j = 1, \dots, 2^t \quad (1.6)$$

$$\mu_{2*j+1}(t+1) = (1 - \epsilon) \mu_j(t) \quad \text{for } j = 1, \dots, 2^t \quad (1.7)$$

$$t = t + 1 \quad (1.8)$$

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

Chapter 2

Supervised Vector Quantization

2.1 Distortion Measures

TODO

Most common is the squared-error distortion:

$$d(x, \hat{x}) = \sum_{i=1}^N |x_i - \hat{x}_i|^2 \quad (2.1)$$

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

$$d(x, \hat{x}) = \left(\sum_{i=1}^N |x_i - \hat{x}_i|^\nu \right)^{\frac{1}{\nu}} = \|x - \hat{x}\|_\nu \quad (2.2)$$

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

$$d(x, \hat{x}) = \sum_{i=1}^N |x_i - \hat{x}_i|^\nu \quad (2.3)$$

The holder Norm (2.2) is a distance and fulfills the triangle inequality¹, the ν^{th} -law distortion not.

All three and many others, as the weighted-squares distortion and the quadratic distortion, depend on the difference $x - \hat{x}$. We call them can be described as $d(x, \hat{x}) = L(x - \hat{x})$. A distortion not having this form is the one by Itakura, Saito and Chaffee,

$$d(x, \hat{x}) = (x - \hat{x})^T R(x) (x - \hat{x}) \quad (2.4)$$

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

, where $R(x)$ is the autocorrelation matrix, see [Linde et al., 1980] for details.

2.2 Learning Vector Quantization

Literature: [Kohonen, 2001, 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

$$s(t) = \begin{cases} +1 & \text{if } x \text{ and } m_c \text{ belong to the same class} \\ -1 & \text{if } x \text{ and } m_c \text{ belong to different classes} \end{cases} \quad (2.5)$$

2.2.1 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

$$m_c(t+1) = m_c(t) + \eta(t)(x - m_c(t)) \quad (2.6)$$

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

$$m_c(t+1) = m_c(t) - \eta(t)(x - m_c(t)) \quad (2.7)$$

Leave the other codebook vectors unchanged

$$m_i(t+1) = m_i(t) \text{ for } i \neq c \quad (2.8)$$

Or in a compressed form:

$$m_i(t+1) = m_i(t) + \eta(t)s(t)\delta_{ci}(x(t) - m_i(t)) \quad (2.9)$$

2.2.2 OLVQ

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

$$\eta_c(t) = \frac{\eta_c(t-1)}{1 + s(t)\eta_c(t-1)} \quad (2.10)$$

Modify update rule from LVQ1.

$$m_i(t+1) = [1 - \eta_i(t)s(t)\delta_{ci}]m_i(t) + \eta_i(t)s(t)\delta_{ci}x(t) \quad (2.11)$$

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

2.2.3 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 .

$$\min\left\{\frac{d(x, m_i)}{d(x, m_j)}, \frac{d(x, m_j)}{d(x, m_i)}\right\} > \frac{1-w}{1+w} \quad (2.12)$$

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.

$$\begin{aligned} m_i(t+1) &= m_i(t) - \eta(t)[x(t) - m_i(t)] \\ m_j(t+1) &= m_j(t) + \eta(t)[x(t) - m_j(t)] \end{aligned} \quad (2.13)$$

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

2.2.4 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:

$$m_k(t+1) = m_k(t) + \epsilon\eta(t)[x(t) - m_k(t)] \quad \text{for } k \in i, j \quad (2.14)$$

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

2.3 Self-Organizing Maps

Sometimes also called self-organizing feature maps.

Chapter 3

Multi-Layer Perceptrons

3.1 Backpropagation

3.1.1 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. $\phi(y) = \sigma(y) = \frac{1}{1+e^{-y}}$)
- $x_{ij}^{(l)}$: input i of neuron j in layer l
- $w_{ij}^{(l)}$: weight of connection from neuron i in layer $l-1$ to neuron j in layer l
- $z_j^{(l)} = \sum_i w_{ij}^{(l)} * a_i^{(l-1)}$ with $a_i^{(0)} = x$
- $a_j^{(l)} = \phi(z_j^{(l)})$
- $o = a^{(L)}$ is the output of the neural network
- η : learning rate (usually $\eta < 1$)
- $b^{(l)}$: bias unit layer l , either 1 if there is a weight for it, or without weight if $b^{(l)}$ is adjusted directly during training
- $\delta_j^{(l)}$: error of neuron j in layer l

3.1.2 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 as 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:

$$\begin{aligned} \frac{\partial E}{\partial w_{ij}^{(L)}} &= \frac{\partial E}{\partial z_j^{(L)}} \frac{\partial z_j^{(L)}}{\partial w_{ij}^{(L)}} = \frac{\partial E}{\partial a_j^{(L)}} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \frac{\partial z_j^{(L)}}{\partial w_{ij}^{(L)}} && \text{(chain rule)} \\ \frac{\partial z_j^{(L)}}{\partial w_{ij}^{(L)}} &= \frac{\partial \sum_i w_{ij}^{(L)} * x_i^{(L-1)}}{\partial w_{ij}^{(L)}} = a_i^{(L-1)} && (3.1) \\ \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} &= \frac{\partial \phi(z_j^{(L)})}{\partial z_j^{(L)}} = (\phi(z_j^{(L)})(1 - \phi(z_j^{(L)}))) = o_j(1 - o_j) \\ &&& \text{(sigmoid derivative, see 3.3.2)} \\ \frac{\partial E}{\partial a_j^{(L)}} &= \frac{\partial E}{\partial o_j} = \frac{\partial \frac{1}{2} \sum_{i=1}^k (t_i - o_i)^2}{\partial o_j} = (o_j - t_j) && \text{(MSE derivative, see 3.2.1)} \end{aligned}$$

putting everything together

$$\frac{\partial E}{\partial w_{ij}^{(L)}} = (o_j - t_j) o_j (1 - o_j) a_i^{(L-1)} \quad (3.2)$$

We can now define the error for neuron j as

$$\begin{aligned} \delta_j^{(L)} &= \frac{\partial E}{\partial z_j^{(L)}} \\ &= \frac{\partial E}{\partial a_j^{(L)}} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \\ &= (o_j - t_j) o_j (1 - o_j) \end{aligned} \quad (3.3)$$

And our weight update $\Delta w_{ij}^{(L)}$ and new weight $\hat{w}_{ij}^{(L)}$

$$\Delta w_{ij}^{(L)} = -\eta \frac{\partial E}{\partial w_{ij}^{(L)}} = -\eta \delta_j^{(L)} a_i^{(L-1)} \quad (3.4)$$

$$\hat{w}_{ij}^{(L)} = w_{ij}^{(L)} + \Delta w_{ij}^{(L)} \quad (3.5)$$

3.1.3 Weight Outputs for Hidden Layers

For hidden layers we need to propagate the error back to the neuron j in layer l . Intuitively we put in the error in the output layer and use the weights to propa-

ate it backwards.

$$\delta_j^{(l)} = \frac{\partial E}{\partial z_j^{(l)}} = \frac{\partial E}{\partial a_j^{(l)}} \frac{\partial a_j^{(l)}}{\partial z_j^{(l)}} \quad (3.6)$$

$$\frac{\partial E}{\partial a_j^{(l)}} = \sum_k \frac{\partial E}{\partial z_k^{(l+1)}} \frac{\partial z_k^{(l+1)}}{\partial a_j^{(l)}} \quad (3.7)$$

$$= \sum_k \frac{\partial E}{\partial z_k^{(l+1)}} \frac{\partial \sum_i w_{ik}^{(l+1)} * a_i^{(l)}}{\partial a_j^{(l)}} \quad (3.8)$$

$$= \sum_k \frac{\partial E}{\partial z_k^{(l+1)}} w_{jk}^{(l+1)} \quad (3.9)$$

$$= \sum_k \delta_k^{(l+1)} w_{jk}^{(l+1)} \quad (3.10)$$

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

$$\delta_j^{(l)} = a_j^{(l)}(1 - a_j^{(l)}) \sum_k \delta_k^{(l+1)} w_{jk}^{(l+1)} \quad (3.11)$$

And finally our weight updates for hidden neurons

$$\Delta w_{ij}^{(l)} = -\eta \frac{\partial E}{\partial w_{ij}^{(l)}} \quad (3.12)$$

$$= -\eta \frac{\partial E}{\partial z_j^{(l)}} \frac{\partial z_j^{(l)}}{\partial w_{ij}^{(l)}} \quad (3.13)$$

$$= -\eta \delta_j^{(l)} a_i^{(l-1)} \quad (3.14)$$

$$\hat{w}_{ij}^{(l)} = w_{ij}^{(l)} + \Delta w_{ij}^{(l)} \quad (3.15)$$

3.1.4 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 be incomplete and wrong

$$A^{(0)} = X \quad (3.16)$$

$$Z^{(l)} = A^{(l-1)} * W^{(l)} \quad (3.17)$$

$$A^{(l)} = \phi(Z^{(l)}) \quad (3.18)$$

$$O = A^{(L)} \quad (3.19)$$

$$E = \frac{1}{2} \sum (T - O) \circ (T - O) \quad (3.20)$$

$$\delta^{(L)} = (O - T) \circ O \circ (1 - O) \quad (3.21)$$

$$\delta^{(l)} = A^{(l)}(1 - A^{(l)})W^{(l+1)}\delta^{(l+1)}\Delta W^{(l)} = -\eta\delta^{(l)}A^{(l-1)} \quad (3.22)$$

$$\hat{W}^{(l)} = W^{(l)} + \Delta W^{(l)} \quad (3.23)$$

3.2 Error Functions

Sometimes also called objective functions or loss-functions.

3.2.1 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.

$$E_{\text{MSE}} = \frac{1}{n} \sum_{i=1}^n (t_i - o_i)^2 \quad (3.24)$$

This would derive to

$$\begin{aligned} \frac{\partial E_{\text{MSE}}}{\partial o_j} &= \frac{1}{n} \sum_{i=1}^n \frac{\partial (t_i - o_i)^2}{\partial o_j} \\ &= \frac{1}{n} \frac{\partial (t_j - o_j)^2}{\partial o_j} \\ &= \frac{1}{n} 2(t_j - o_j)(-1) \\ &= \frac{2}{n} (o_j - t_j) \end{aligned} \quad (3.25)$$

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

$$E_{\text{MSE}} = \frac{1}{2} \sum_{i=1}^k (t_i - o_i)^2 \quad (3.26)$$

$$\frac{\partial E_{\text{MSE}}}{\partial o_j} = -(t_j - o_j) = o_j - t_j \quad (3.27)$$

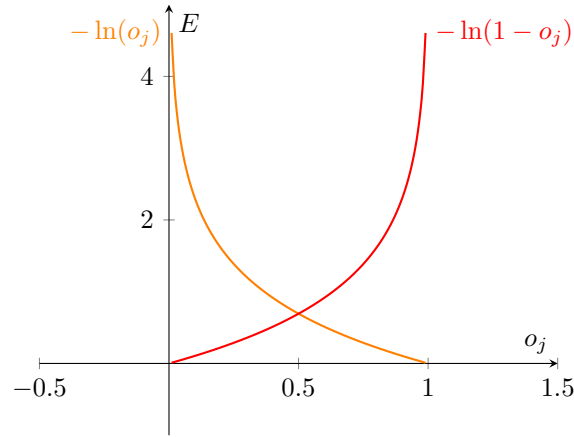


Figure 3.1: Error produced for $t_j = 1$ (orange) and for $t_j = 0$ (red). o_j must be in range $(0, 1]$, otherwise the logarithm can be complex or undefined or the error gets rediculously large.

3.2.2 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.

$$E_{CE} = - \sum_{i=1}^k (t_i \log(o_i) + (1 - t_i) \log(1 - o_i)) \quad (3.28)$$

This will derive to:

$$\begin{aligned} \frac{\partial E_{CE}}{\partial o_j} &= - \frac{\partial}{\partial o_j} \sum_{i=1}^k (t_i \log(o_i) + (1 - t_i) \log(1 - o_i)) \\ &= - \frac{\partial}{\partial o_j} t_j \log(o_j) - \frac{\partial}{\partial o_j} (1 - t_j) \log(1 - o_j) \\ &= - \frac{t_j}{o_j} + \frac{1 - t_j}{1 - o_j} \\ &= \frac{1 - t_j}{1 - o_j} - \frac{t_j}{o_j} \end{aligned} \quad (3.29)$$

3.3 Activation Functions

3.3.1 Step Function

$$\phi(x) = \begin{cases} 1 & \text{if } x > 0 \\ 0 & \text{if } x \leq 0 \end{cases} \quad (3.30)$$

The derivative is always 0.

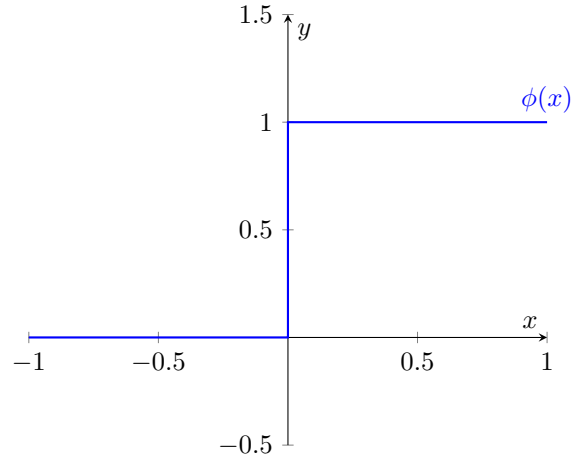


Figure 3.2: Step Function

3.3.2 Sigmoid

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

$$\sigma(x) = \frac{1}{1 + e^{-\beta x}} \quad (3.31)$$

And it's derivative:

$$\begin{aligned} \frac{d\sigma(x)}{dx} &= \frac{d}{dx} (1 + e^{-\beta x})^{-1} \\ &= (-1)(1 + e^{-\beta x})^{-2} (-\beta e^{-\beta x}) \\ &= \beta (1 + e^{-\beta x})^{-2} e^{-\beta x} \\ &= \beta \frac{1}{1 + e^{-\beta x}} \frac{e^{-\beta x}}{1 + e^{-\beta x}} \\ &= \beta \sigma(x) \frac{e^{-\beta x} + 1 - 1}{1 + e^{-\beta x}} \\ &= \beta \sigma(x) \left(1 - \frac{1}{1 + e^{-\beta x}}\right) \\ &= \beta \sigma(x) (1 - \sigma(x)) \end{aligned} \quad (3.32)$$

3.3.3 Softmax Function

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

$$a_i = \phi(z_i) = \frac{e^{z_i}}{\sum_j e^{z_j}} \quad (3.33)$$

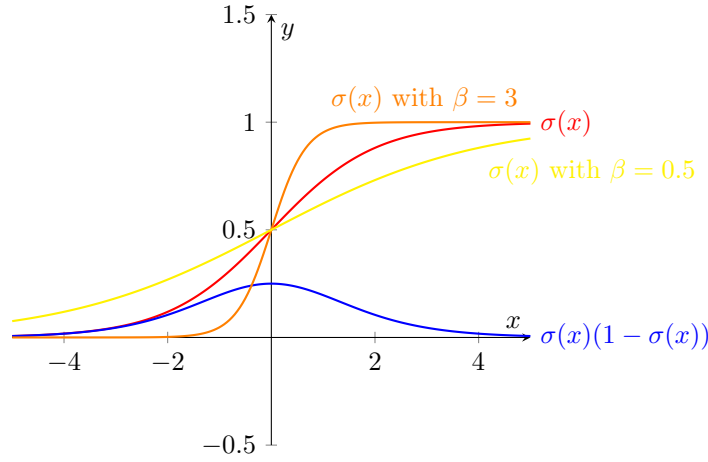


Figure 3.3: Sigmoid Function (red) and its derivative (blue).

And it's derivative:

$$\begin{aligned}
 \frac{\partial \phi(z_i)}{\partial z_i} &= \frac{\left(\frac{\partial}{\partial z_i} e^{z_i} \right) \left(\sum_j e^{z_j} \right) - e^{z_i} \left(\frac{\partial}{\partial z_i} \left(\sum_j e^{z_j} \right) \right)}{\left(\sum_j e^{y_j} \right)^2} \\
 &= \frac{e^{y_i} \sum_j e^{y_j} - e^{y_i} e^{y_i}}{\left(\sum_j e^{y_j} \right)^2} \\
 &= \frac{e^{y_i} \sum_j e^{y_j}}{\left(\sum_j e^{y_j} \right)^2} - \frac{e^{y_i} e^{y_i}}{\left(\sum_j e^{y_j} \right)^2} \\
 &= \frac{e^{y_i}}{\sum_j e^{y_j}} - \left(\frac{e^{y_i}}{\sum_j e^{y_j}} \right)^2 \\
 &= \phi(y_i) - \phi(y_i)^2 \\
 &= \phi(y_i)(1 - \phi(y_i))
 \end{aligned} \tag{3.34}$$

3.3.4 Hyperbolic Tangent Function

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

$$\phi(x) = \tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} \tag{3.35}$$

And it's derivative:

$$\frac{d\phi(x)}{dx} = 1 - \tanh(x) = \frac{(e^x - e^{-x})^2}{(e^x + e^{-x})^2} \tag{3.36}$$

3.3.5 Linear Function

$$\phi(x) = x \tag{3.37}$$

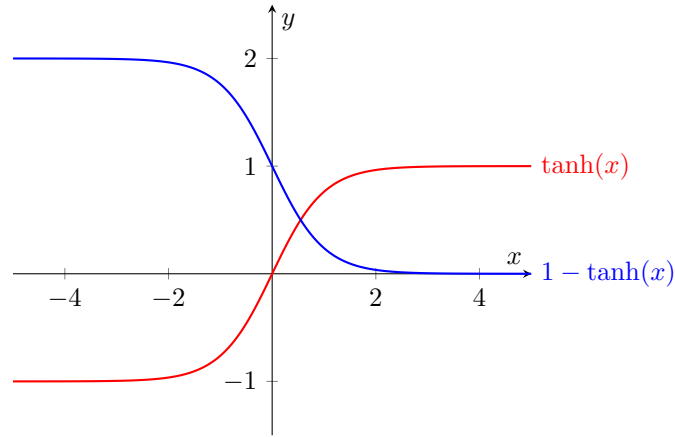


Figure 3.4: Hyperbolic Tangent Function (red) and its derivative (blue).

3.3.6 Rectified Linear Unit

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

$$\phi(x) = \max(0, x) \quad (3.38)$$

And it's derivative:

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

3.3.7 Softplus

Smoothed version of ReLU.

$$\phi(x) = \ln(1 + e^x) \quad (3.39)$$

And it's derivative:

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

3.3.8 Maxout

Outputs the maximum of its inputs

$$a_j^{(l)} = \max_i z_i^{(l)}, \quad (j-1)g + 1 \leq i \leq jg \quad (3.40)$$

$$z_i^{(l)} = \sum_k x_k^{(l-1)} w_{ki}^{(l)} + b^{(l)} \quad (3.41)$$

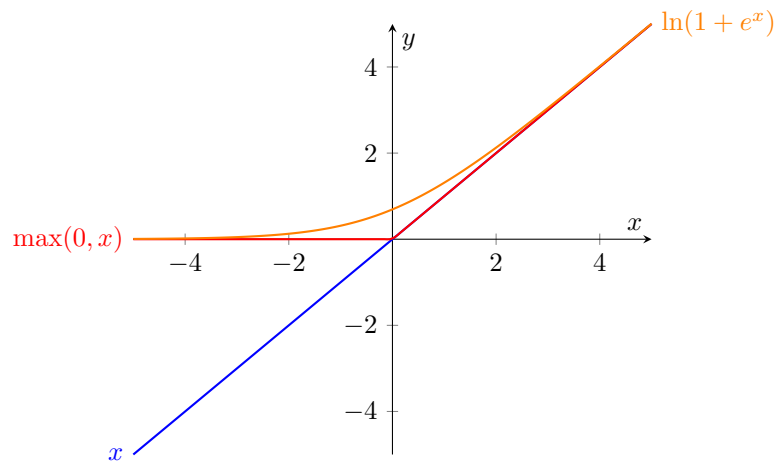


Figure 3.5: Linear Function (blue), Rectified Linear Unit (red), Softplus (orange)

Chapter 4

Preprocessing

4.1 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),

$$\hat{x}_i = \frac{x_i - \min(x)}{\max(x) - \min(x)} \quad (4.1)$$

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

$$\hat{x}_i = \frac{x_i - \text{mean}(x)}{\text{std}(x)} = \frac{x_i - \mu}{\sigma} \quad (4.2)$$

, 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.

Bibliography

- [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.

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