



Neural Networks

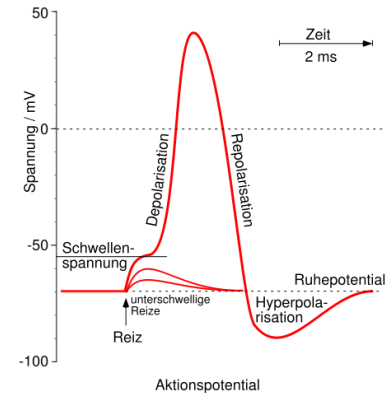
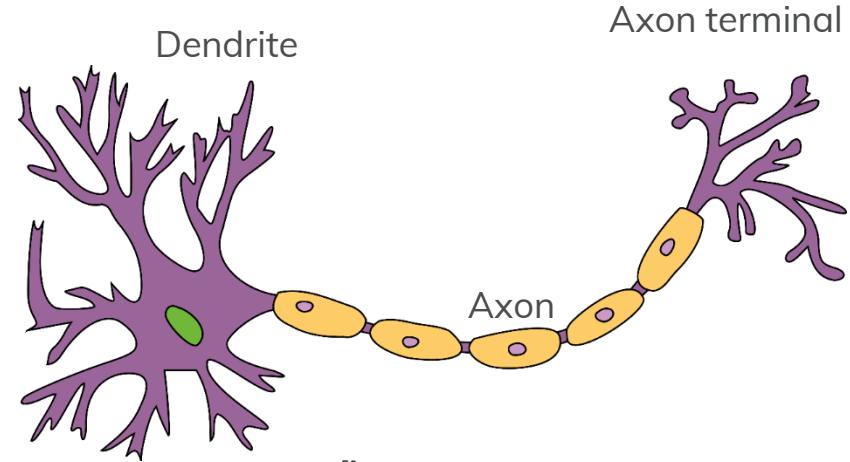
Programming for Data Science

Examination dates

- For urgent matters: 07.02.2020
- However, we recommend the 03.04.2020.

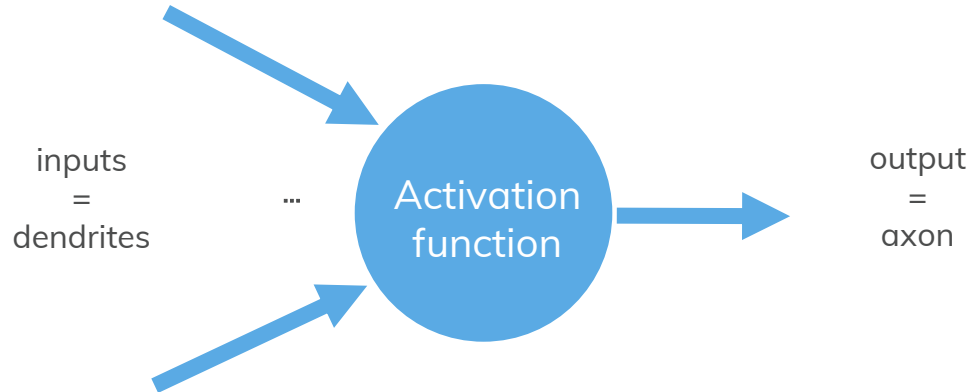
Biological Basics

- Neuron/nerve cell is the structural and functional building block of nervous system
- Cell specializes in receiving, processing and transmitting excitation
- Dendrites receive excitation from other cells
- Axon transmits excitation to other cells
- Transmission: internal electric, external chemical (neurotransmitter, synapses)
- Incoming excitation summed at axon hillock
- If threshold potential is exceeded, action potential is released (all-or-nothing)
- Analog/Digital Converter



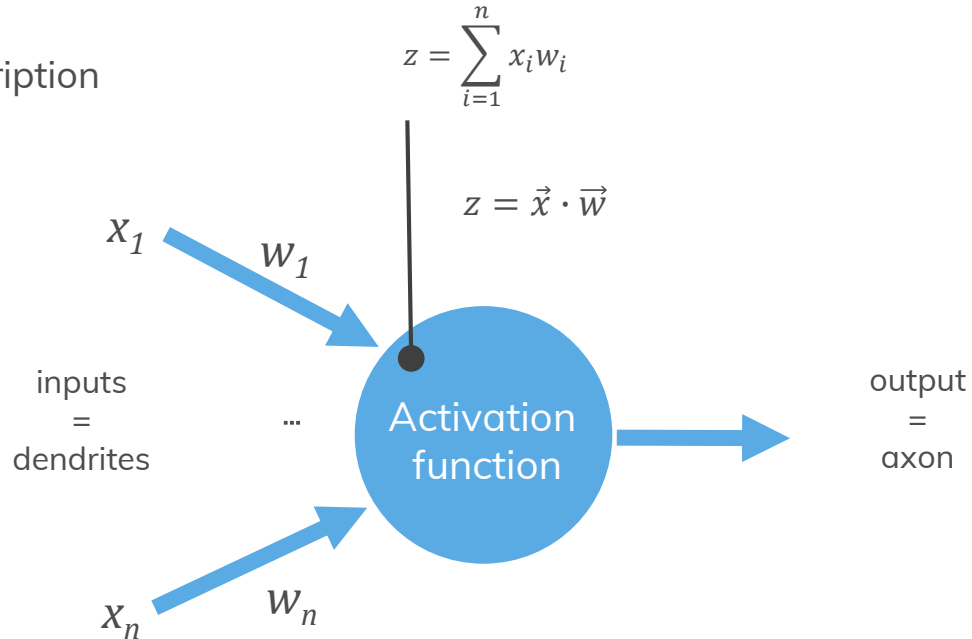
Basics

- Core element is the artificial neuron → Perceptron
- Published 1958



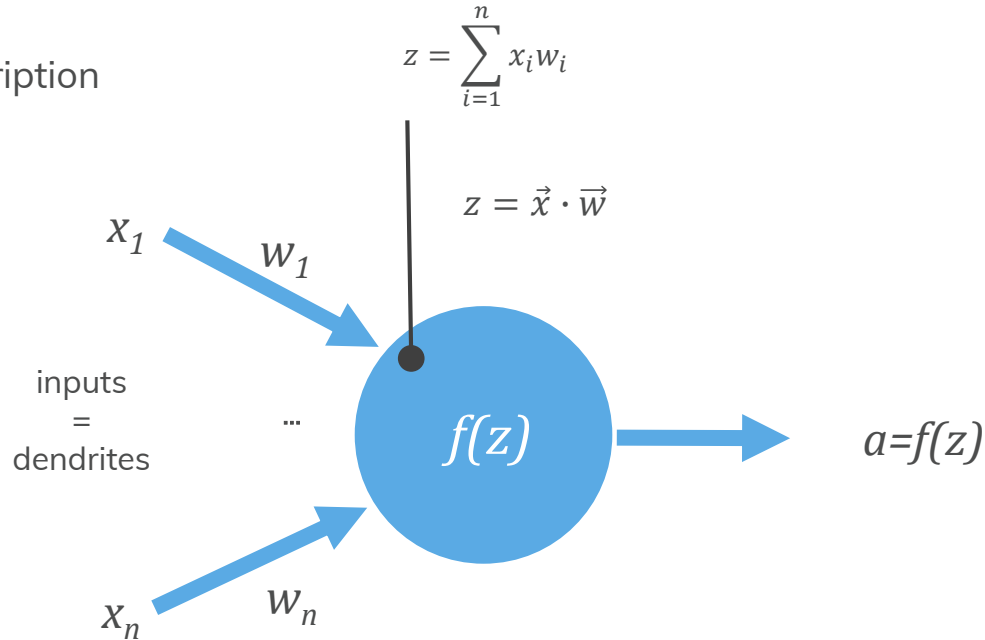
Basics

- Mathematical description



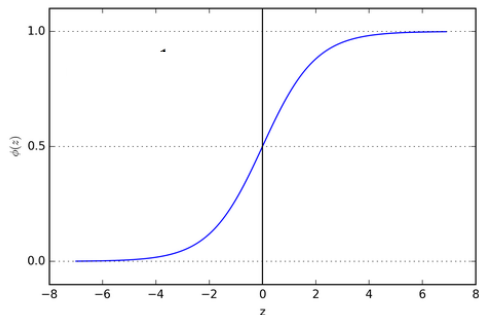
Basics

- Mathematical description



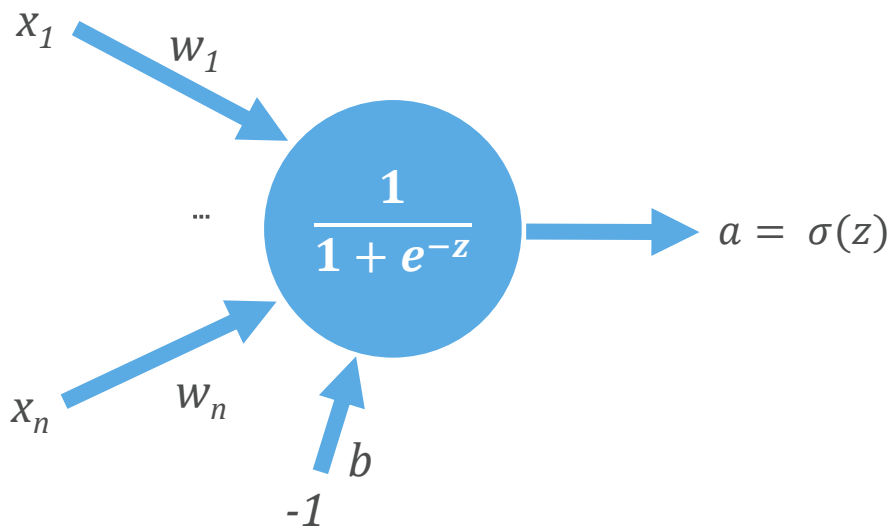
Sigmoid Neuron

- Smoothed Perceptron
- Basic structure is kept
- Inputs and outputs range from 0 to 1
- Activation thresholds are modelled as constant input / bias
- Activation function → Sigmoid function



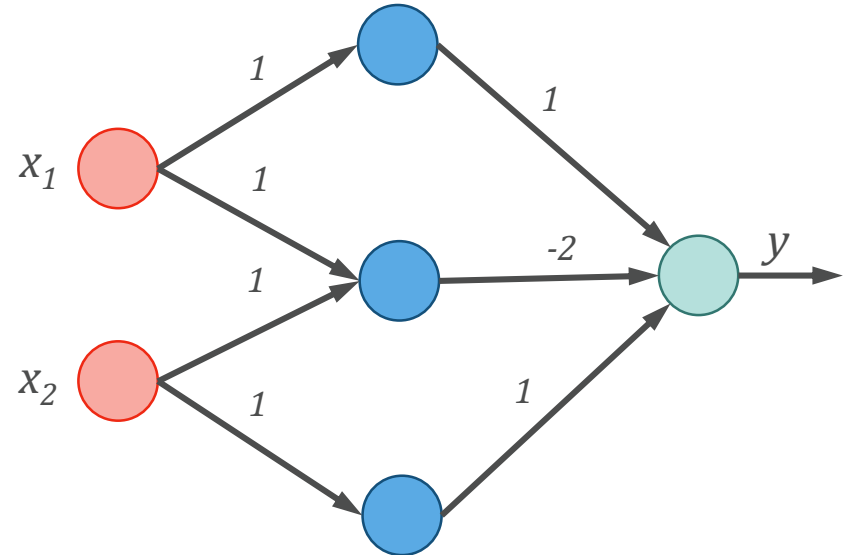
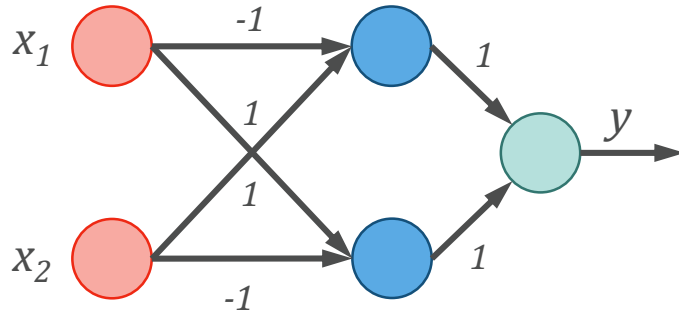
$$z = b + \sum_{i=1}^n x_i w_i$$

$$z = b + \vec{x} \cdot \vec{w}$$



Multilayer Perceptron

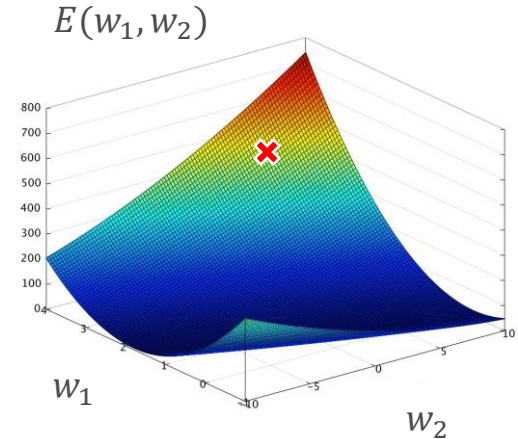
- Input layer
- Hidden layer
- Output layer
- Arbitrary Combinations



Gradient descent

- \rightarrow minimize an error function E with multiple parameters
- Example: neuron with two inputs $\rightarrow E(w_1, w_2)$
- Calculate gradient $\nabla E(w_1, w_2) \rightarrow$ points into direction of biggest increase
- Negative gradient points into direction of biggest reduction
- Gradient is a vector containing the partial derivatives of the error function

$$\nabla E(w_1, \dots, w_n) = \left\langle \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right\rangle$$



Gradient descent

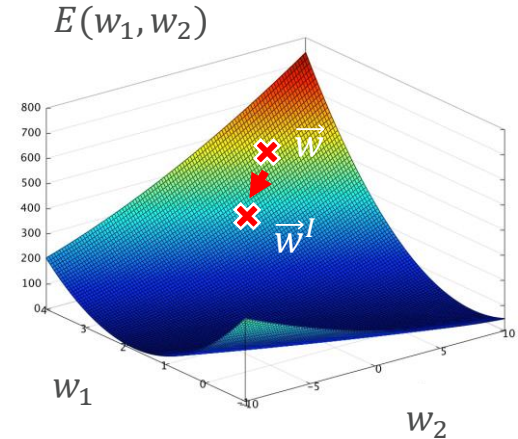
- Gradient is a vector containing the partial derivatives of the error function

$$\nabla E(w_1, \dots, w_n) = \left\langle \frac{\partial E}{\partial w_1}, \dots, \frac{\partial E}{\partial w_n} \right\rangle$$

- Using the gradient, new weights are calculated

$$\vec{w}^I = \vec{w} - \alpha \nabla E(\vec{w})$$

- α step size \rightarrow learning rate
- Step-by-step adjustment until gradient = 0 or smaller than threshold \rightarrow convergence
- Risk of local minimum \rightarrow initialization



Gradient descent

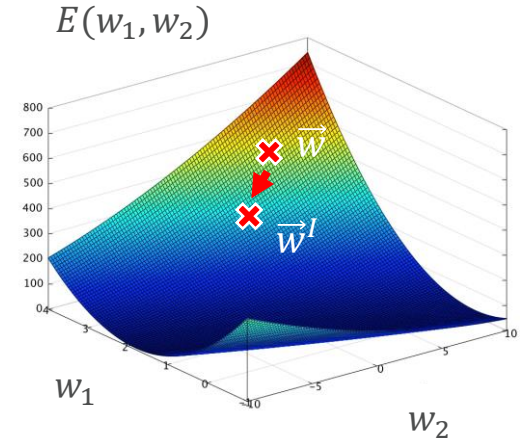
- Gradient is used to calculate new weights

$$\vec{w}^I = \vec{w} - \alpha \nabla E(\vec{w})$$

- Simplified form, actually:

$$\vec{w}^I = \vec{w} - \alpha \nabla \frac{1}{n} \sum_{x=1}^n (t_x - a_x)$$

- Consider all n training inputs during calculation of average error
- → costly, slow



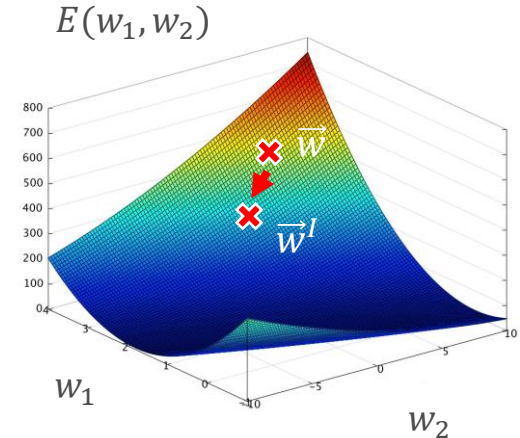
Gradient descent

$$\vec{w}^I = \vec{w} - \alpha \nabla \frac{1}{n} \sum_{x=1}^n (t_x - a_x)$$

- Stochastic Gradient Descent \rightarrow approximation
- Instead of n inputs only consider one

$$\vec{w}^I = \vec{w} - \alpha \nabla (t_x - a_x)$$

- Mathematically not optimal, but “good enough”
- Path less direct due to noise/outliers
- More iterations until convergence, BUT greatly reduce cost per iteration



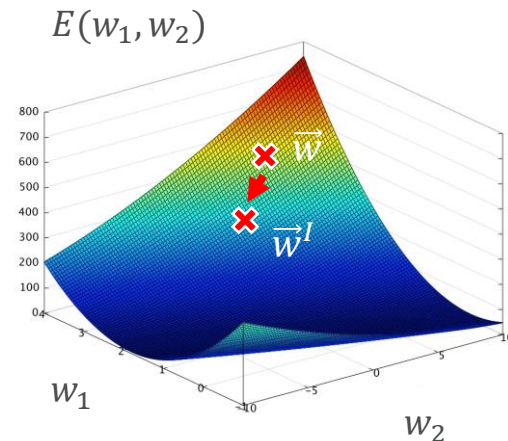
Gradient descent

$$\vec{w}^I = \vec{w} - \alpha \nabla \frac{1}{n} \sum_{x=1}^n (t_x - a_x)$$

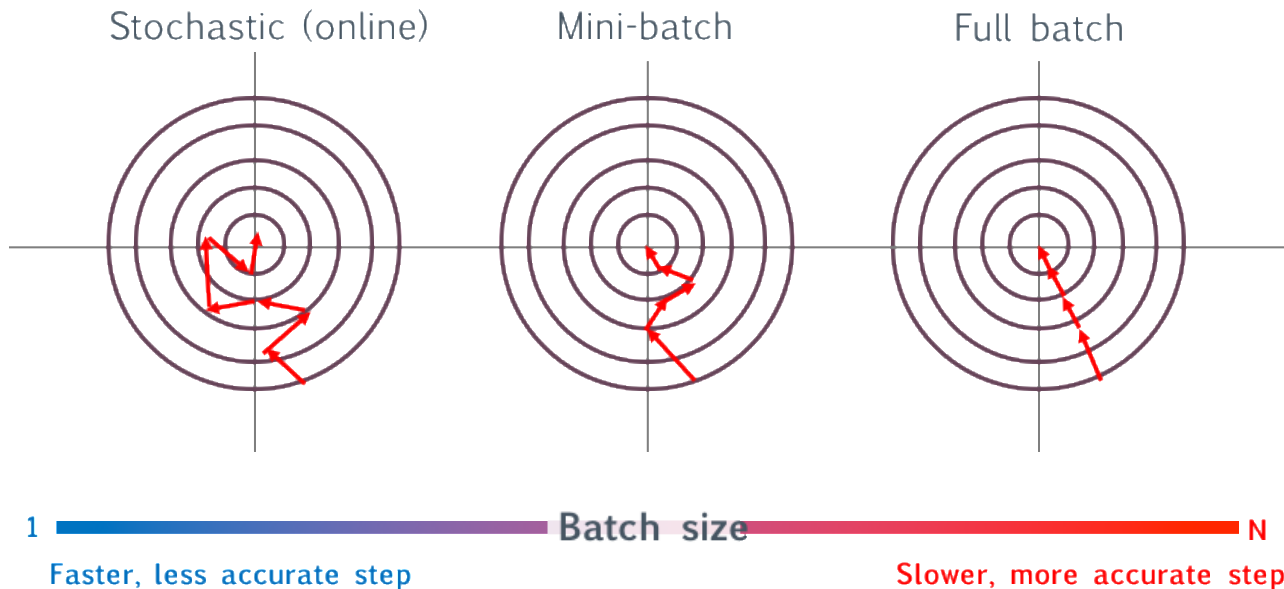
- Mini Batch Gradient Descent \rightarrow compromise
- Instead of n inputs use m with $m < n$

$$\vec{w}^I = \vec{w} - \alpha \nabla \frac{1}{m} \sum_{x=1}^m (t_x - a_x)$$

- “Best of both worlds“
- Most common approach for training
- Rule of thumb: $16 \leq m \leq 256$



Gradient descent - approaches



Expansions

- Alternatives to gradient descent
 - Newton's method
 - Conjugate gradient method
 - Quasi-Newton method
 - Levenberg-Marquardt Algorithm (only sum-of-squared errors)
- In general these methods converge faster due to usage of Hessian and Jacobian matrices
- Calculation of these matrices requires additional storage space, therefore these methods are not suitable for large networks with thousands of parameters
- Ranking with regards to storage requirements and speed
- Gradient descent < Conjugate gradient < Newton < Quasi Newton < Levenberg-Marquardt



speed/storage

Backpropagation

- Goal: Calculate $\frac{\partial E}{\partial w_{ij}^l} \rightarrow$ Chain rule

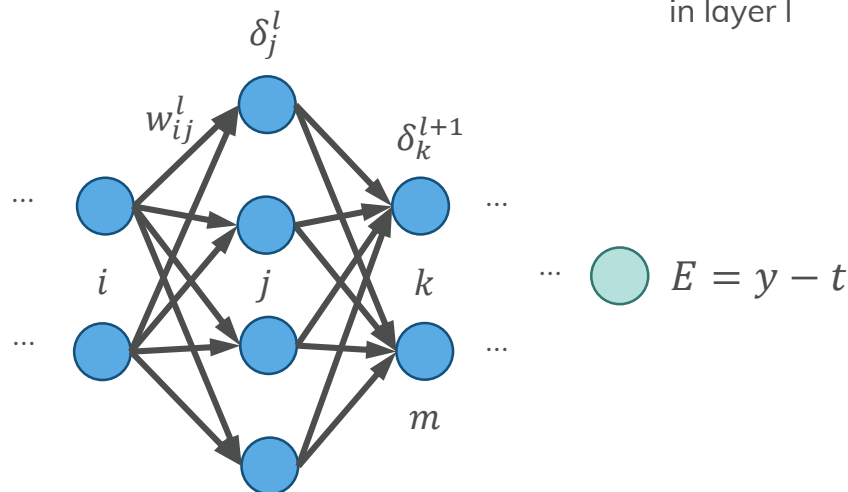
$$\frac{\partial E}{\partial w_{ij}^l} = \frac{\partial E}{\partial a_j^l} \cdot \frac{\partial a_j^l}{\partial z_j^l} \cdot \frac{\partial z_j^l}{\partial w_{ij}^l}$$

$$\frac{\partial E}{\partial w_{ij}^l} = \dots = \underbrace{(y - t) \cdot \sigma(z_j^l)(1 - \sigma(z_j^l))}_{\text{Error } \delta_j^l} \cdot a_i^{l-1}$$

$$\delta_j^l = \left(\sum_{k=1}^m w_{jk}^{l+1} \delta_k^{l+1} \right) \cdot \sigma(z_j^l)(1 - \sigma(z_j^l))$$

- Update weights: $w_{ij}^l = w_{ij}^l - \alpha \delta_j^l a_i^{l-1}$

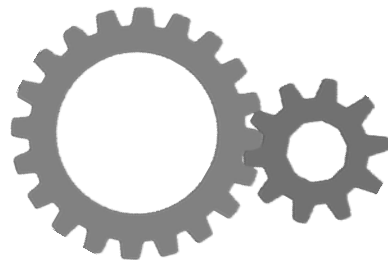
w_{ij}^l weight from i-th neuron
in layer l-1
to j-th neuron
in layer l



Backpropagation – Problems

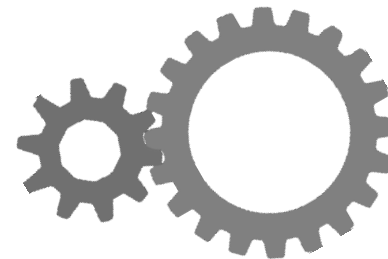
Vanishing Gradient

- Backpropagation over multiple layers leads to smaller and smaller changes \rightarrow greatly decreases learning speed
- $\sigma'(x) \leq 0,25$ and $0 < w < 1$
- Repeated multiplication of small values \rightarrow product gets smaller



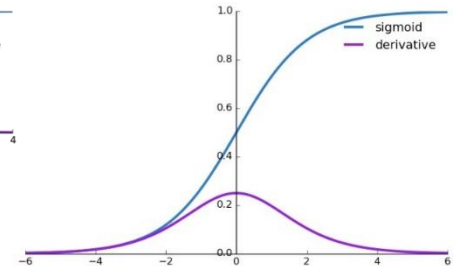
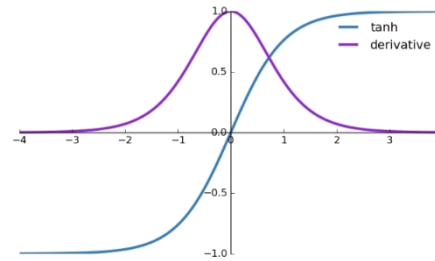
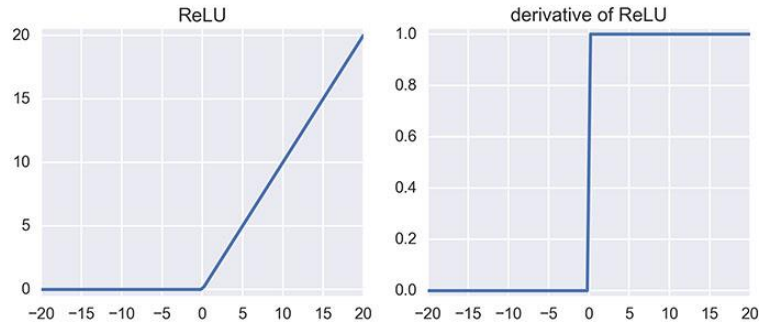
Exploding Gradient

- $w_j \sigma'(z_j) > 1 \rightarrow$ too fast and unstable learning in the front layers
- Basic problem: gradients depend on all subsequent layers \rightarrow more layers mean more instability



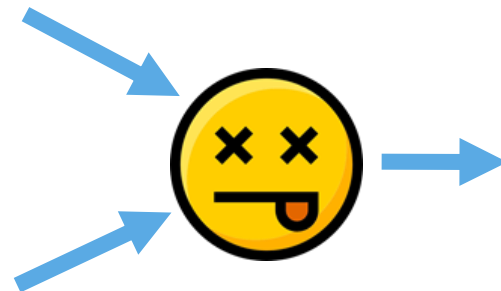
Backpropagation – Activation functions

- Alternatives to sigmoid
- Hyperbolic tangent (tanh): scaled sigmoid $\tanh(z) = 2\sigma(2z) - 1$, $\tanh'(x) \leq 1$
- Rectified Linear Unit (ReLU): $A(z) = \max(0, z)$



Backpropagation – Activation functions

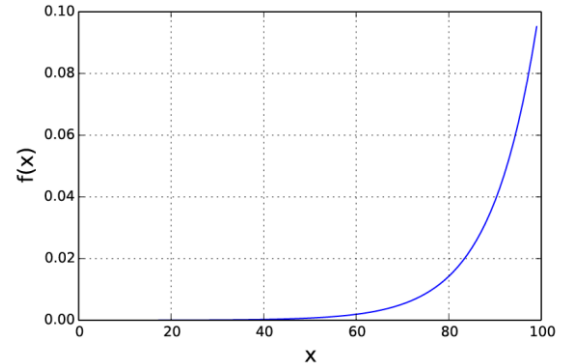
- Rectified Linear Unit (ReLU): $A(z) = \max(0, z)$
- Pros:
 - no “vanishing gradients”
 - selective neuron activation
 - easy to calculate
 - good performance
- Cons:
 - inflates activation \rightarrow negative effect on learning
 - Prone to overfitting
 - “dead neuron” problem: if weights of an neuron lead to $z_j \leq 0$, the neuron can never be activated again
- \rightarrow different variants



Backpropagation – Activation functions

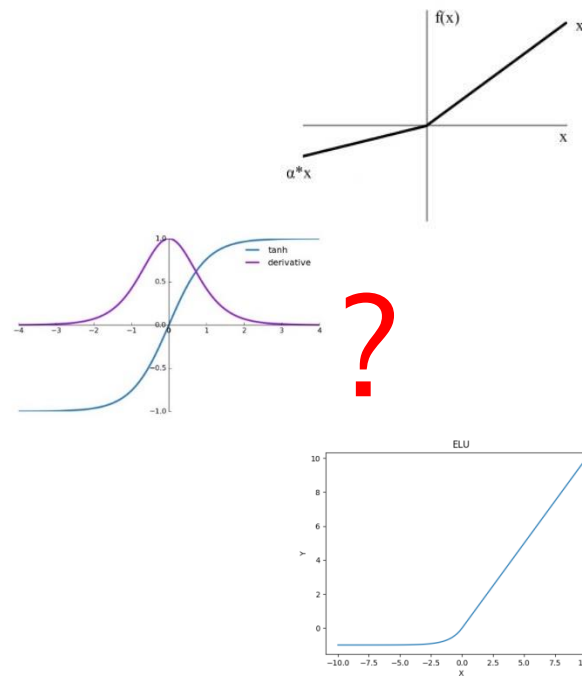
- Softmax
 - Transforms output into probabilities
 - → ideal for classification
 - Sigmoid allows only two classes
 - Generally for output layer
-
- z is input vector to output layer
 - K is number of output neurons

$$\text{softmax}(z)_j = \frac{e^{z_j}}{\sum_{k=1}^K e^{z_k}}$$



Backpropagation – Activation functions

- What function do I use?
- “Depends on the problem“
- Sigmoid and Tangent suited for classification
- Currently ReLu is considered the best default
- However: Try!
- Why use activation functions at all?
 - Provide non-linearity → make hidden layers useful
 - Linear function of linear functions is linear



Training – Regularization

- L2 regularization (weight decay)
- Add penalty term to error function

$$E = \frac{1}{2n}(t - y)^2 + \frac{\lambda}{2n} \sum_w w^2 \quad \lambda > 0, n \text{ number of training samples}$$

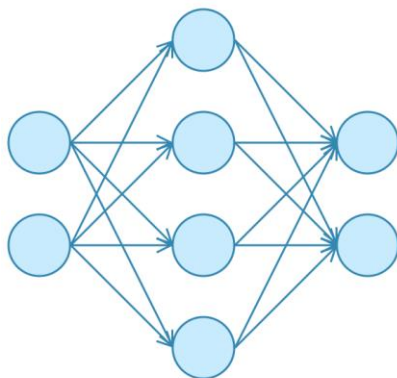
- Favors small weights, big weights have to provide strong improvement
- Trade-off between error minimization and small weights
- Controlled using λ , large $\lambda \rightarrow$ small weights, small $\lambda \rightarrow$ error minimization

- L1 regularization:

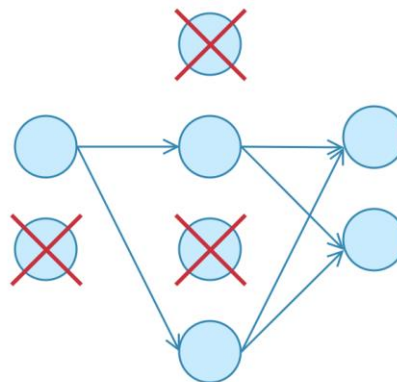
$$E = \frac{1}{2n}(t - y)^2 + \frac{\lambda}{n} \sum_w |w|$$

Training – Dropout

- Currently most popular method to prevent overfitting
- In each iteration of trainings, some neurons are “dropped“
- Inputs and output are ignored for these neurons
- Dropout rate typically 20%-50%



No Dropout



With Dropout

Step 0

- Get the Boston housing price dataset from keras.
- Use the predefined split into training and test data.

Step 1

- Use a funnel MLP* with width (w_1) 512 and depth (l_{\max}) 2 to predict the house prices. ($w_{l+1} = \frac{w_l}{2}$)
- Other parameters: RELU activations, MSE, the Adam optimizer, 100 epochs and a batch size of 32.
- Evaluate the NN on the test data.

Step 2

- Min-max normalize the features. Then, train the same MLP from step 1. Compare the errors.

Step 3

- Use a grid search* to find the best combination of w_1 and l_{\max} for the normalized data.

*use your own implementation

Package suggestions

R

- (data.table)
- keras

python3

- numpy
- pandas
- keras

Exercise Appointment

We compare and discuss the results

- Tuesday, 14.01.2020,
- Consultation: 09.01.2020,
- Please prepare your solutions! Send us your code!

If you have questions, please mail us:

claudio.hartmann@tu-dresden.de Orga + Code + R

lucas.woltmann@tu-dresden.de Tasks + Python

