

FliK Modul 2020

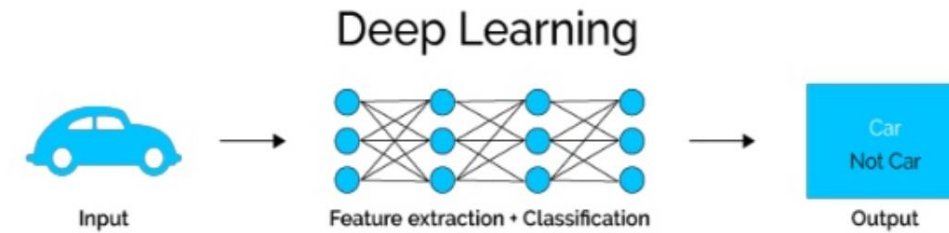
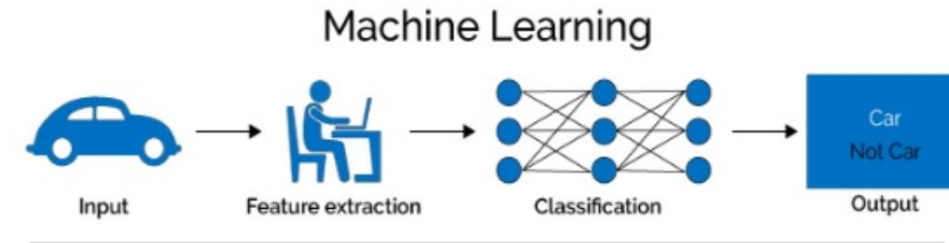
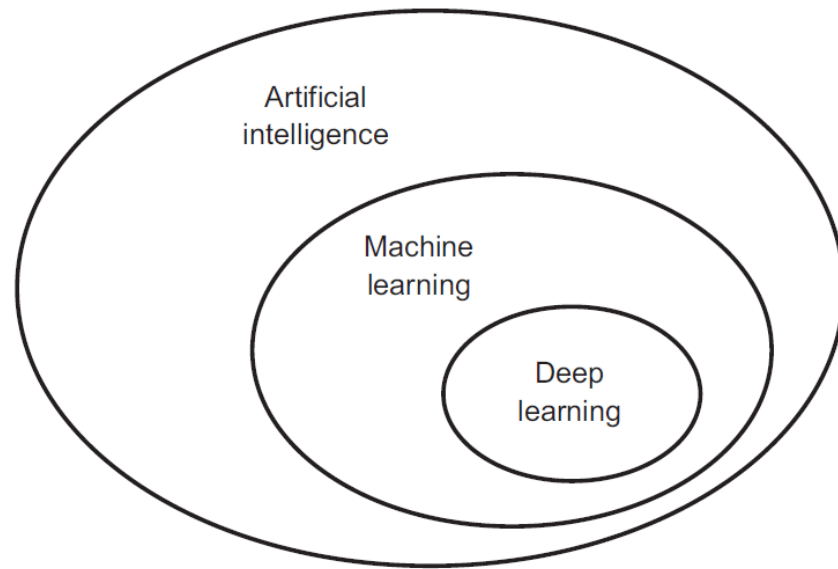
# Neural Networks

Steffen Seitz, Marvin Arnold & Markus Fritzsche

Prof. Ronald Tetzlaff

Dresden, 19-23.10.

# Neural Networks & Deep Learning



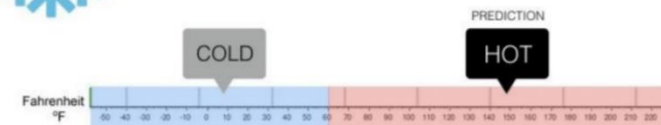
## Regression

What is the temperature going to be tomorrow?

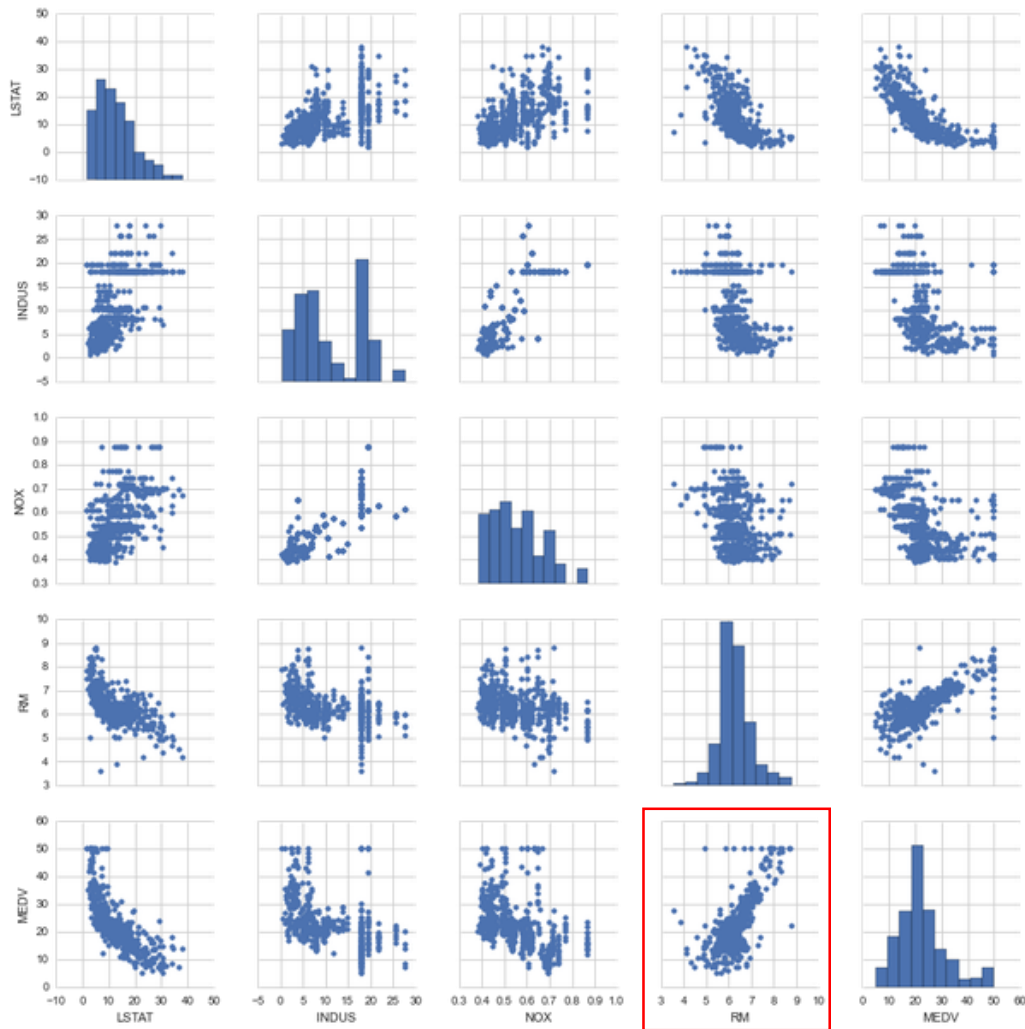


## Classification

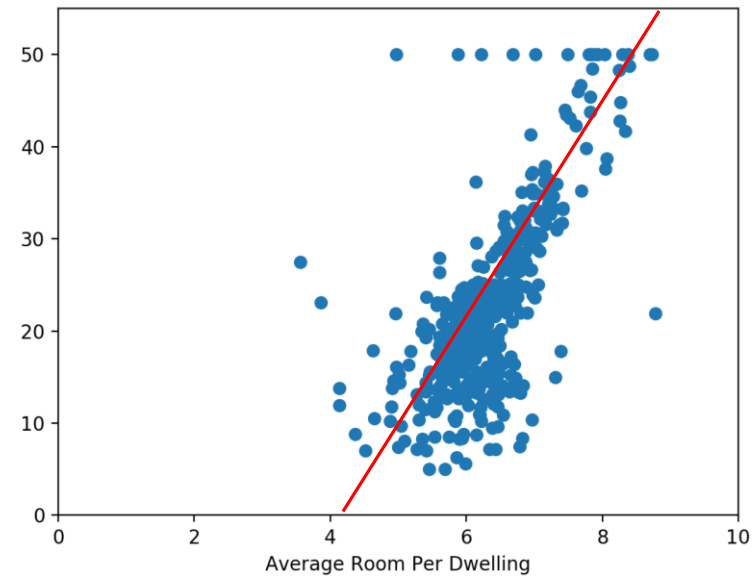
Will it be Cold or Hot tomorrow?



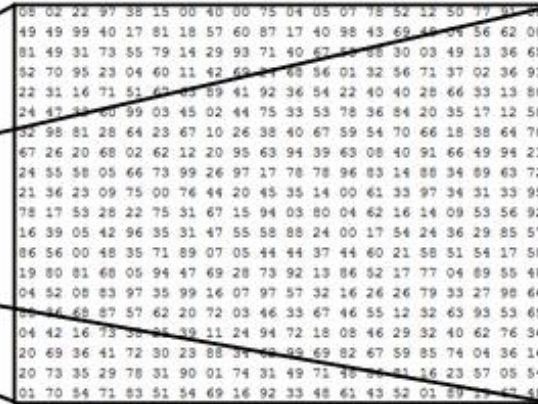
# Regression: Boston Housing Dataset



<i>CRIM</i>	Per capita crime rate by town
<i>ZN</i>	Proportion of residential land zoned for lots over 25,000 ft <sup>2</sup>
<i>INDUS</i>	Proportion of nonretail business acres per town
<i>CHAS</i>	Charles River dummy variable (= 1 if tract bounds river; = 0 otherwise)
<i>NOX</i>	Nitric oxide concentration (parts per 10 million)
<i>RM</i>	Average number of rooms per dwelling
<i>AGE</i>	Proportion of owner-occupied units built prior to 1940
<i>DIS</i>	Weighted distances to five Boston employment centers
<i>RAD</i>	Index of accessibility to radial highways
<i>TAX</i>	Full-value property-tax rate per \$10,000
<i>PTRATIO</i>	Pupil/teacher ratio by town
<i>B</i>	$1000(B_k - 0.63)^2$ where $B_k$ is the proportion of blacks by town
<i>LSTAT</i>	% Lower status of the population
<i>MEDV</i>	Median value of owner-occupied homes in \$1000s



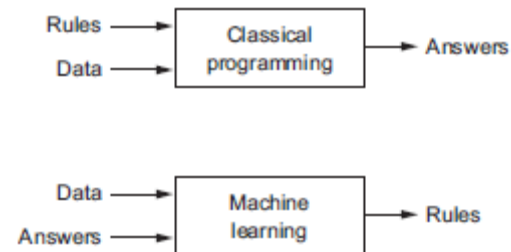
# Classification:

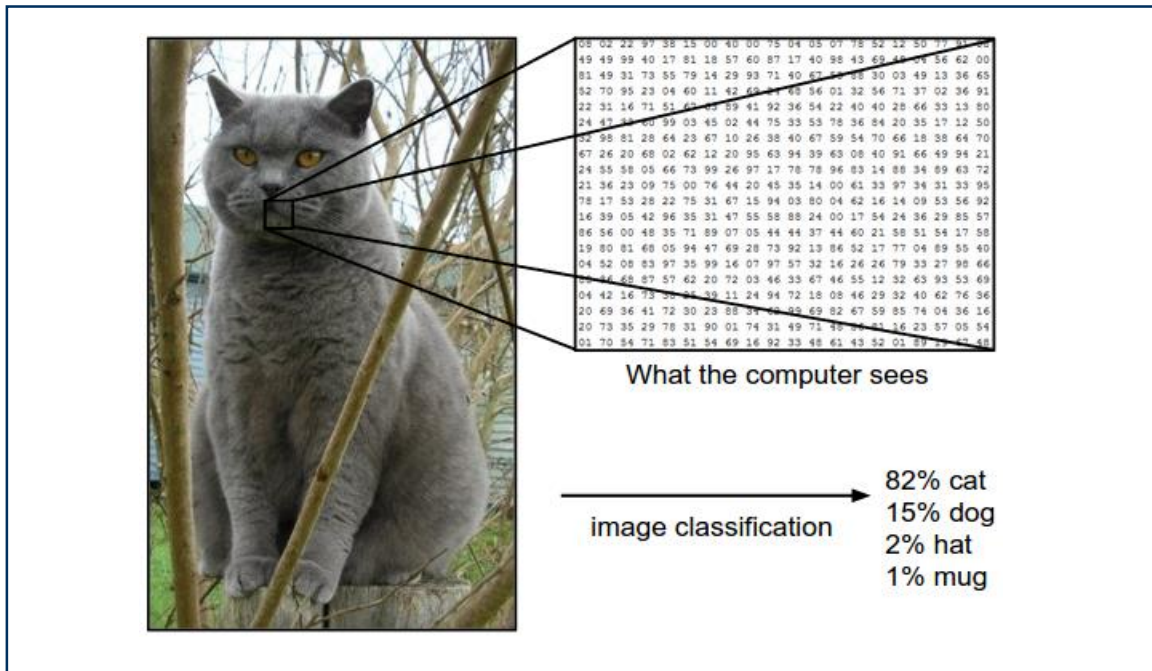


What the computer sees

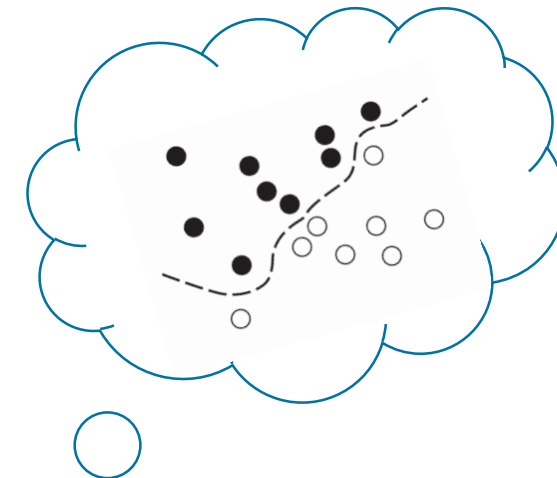


- Cat detection is a complex task!
- Neural Networks try to unfold the complexity

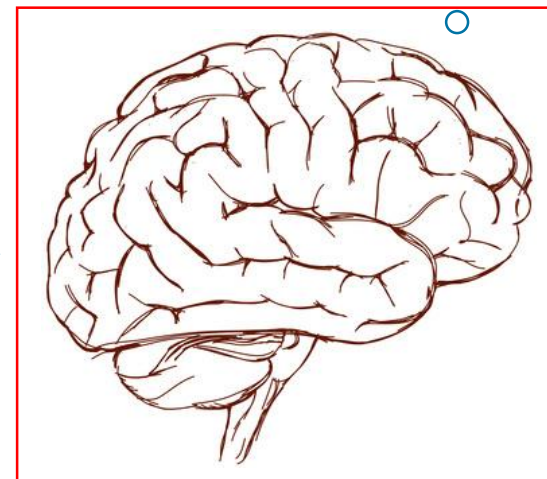




## Decision Boundary



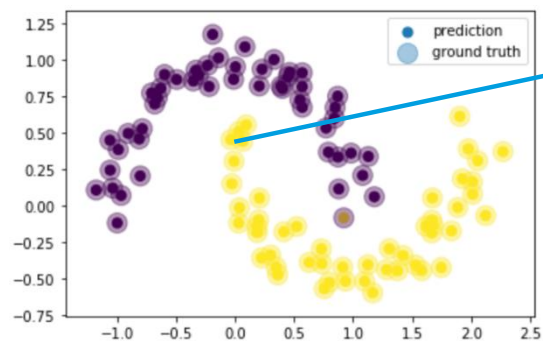
## Network



## Class decision



## Data



## Vektor/Matrix

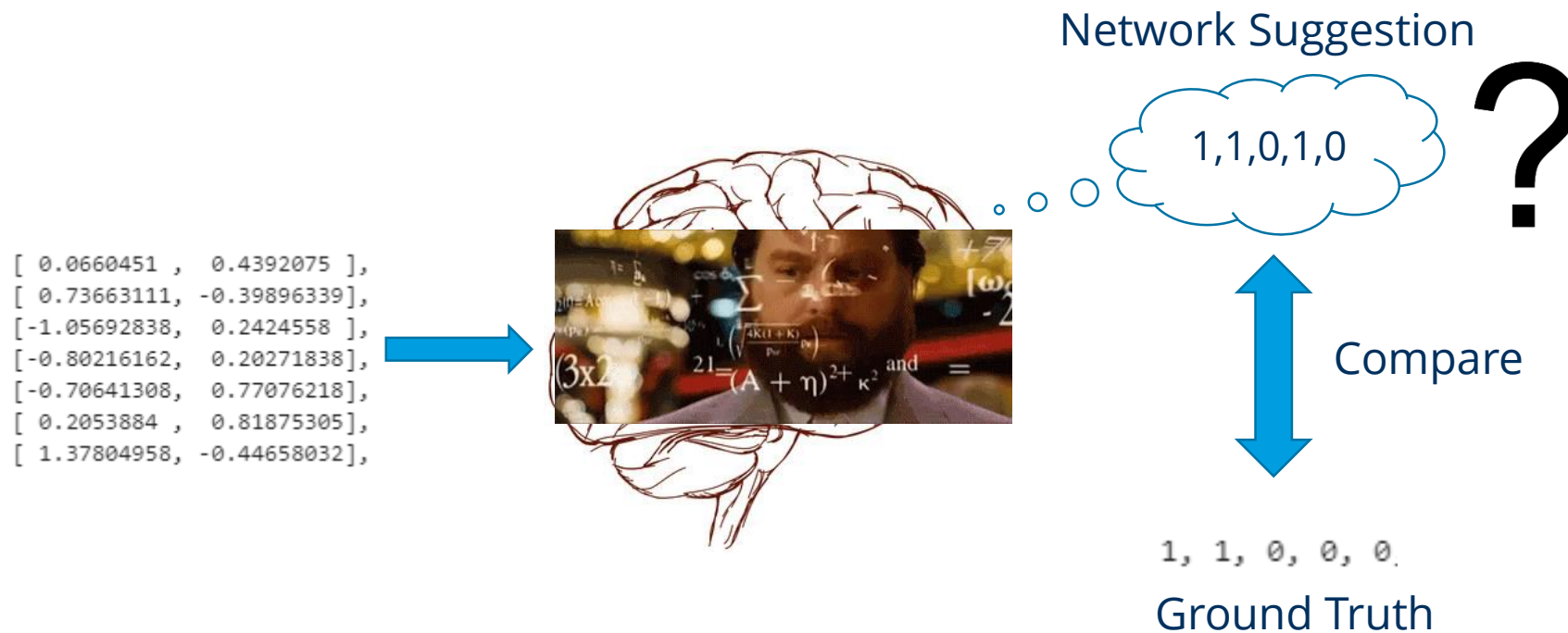
```
[ 0.0660451 ,  0.4392075 ],
[ 0.73663111, -0.39896339],
[-1.05692838,  0.2424558 ],
[-0.80216162,  0.20271838],
[-0.70641308,  0.77076218],
[ 0.2053884 ,  0.81875305],
[ 1.37804958, -0.44658032],
```



# Multilayer Perceptron

## Loss functions

A loss function is a grade (**Metric**) how good we have been doing our task.  
Mathematically speaking, a metric is a **measure of „Distance“**.



# Multilayer Perceptron

## Loss functions

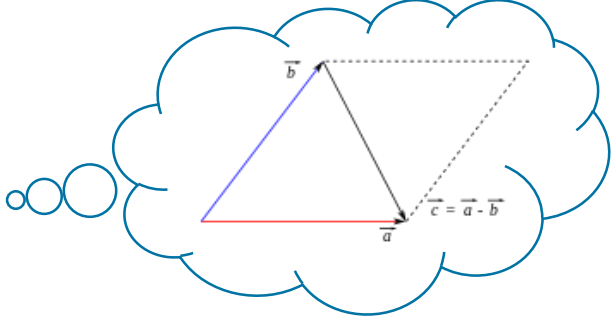
### MSE (Naive approach)

Computes the mean square Error which can be defined as the **distance** of **two points** in a vector space.

$$\text{MSE} = \frac{1}{N} \sum_{i=1}^N (y_i - a)^2$$

Ground truth  $y_i$  (green box), Suggestion  $a$  (blue box)

Mean over all (possible) classes



### Crossentropy

The (mean) cross entropy is a measure of **difference** between **two discrete probability distributions**. (created by Claude Shannon)

$$C = -\frac{1}{N} \sum_{i=1}^N [y_i \ln a + (1 - y_i) \ln (1 - a)]$$

Mean over all (possible) classes

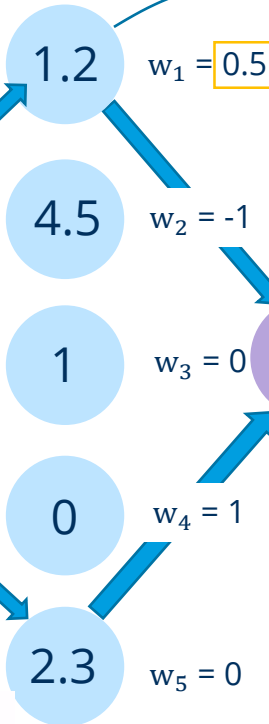
Ground truth  $y_i$  (green box), Suggestion  $a$  (blue box)

The naive approach of applying **MSE** for regression is ok. But (later) we will discuss a reason why classification **crossentropy** is the way to go in **nearly every** Situation!

# Multilayer Perceptron

## Forward Pass

Input vector:  $\begin{bmatrix} 0.0660451, 0.4392075, 0.7366311, -0.39896339, -1.05692838, 0.2424558, -0.80216162, 0.20271838, -0.70641308, 0.77076218, 0.2053884, 0.81875305, 1.37804958, -0.44658032 \end{bmatrix}$



$$output = \sigma \left( \sum_{j=1} \omega_j x_j + b \right)$$

$b = 0$

$z$

$$output = \sigma(1.2 \times 0.5 - 4.5 \times 1) = 0.05$$

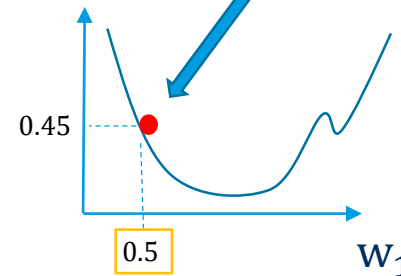
Target:  $\begin{bmatrix} 1, 1, 0, 0 \end{bmatrix}$

$$J_{total}(output) = \frac{1}{2}(target - output)^2$$

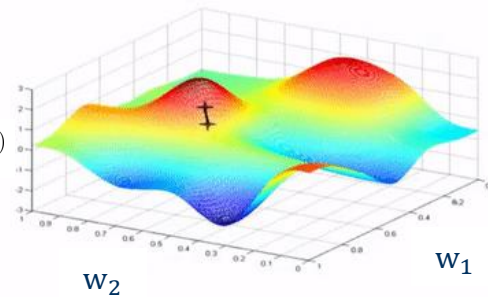
(MSE loss)

$$= \frac{1}{2}(1 - 0.05)^2 = 0.45125$$

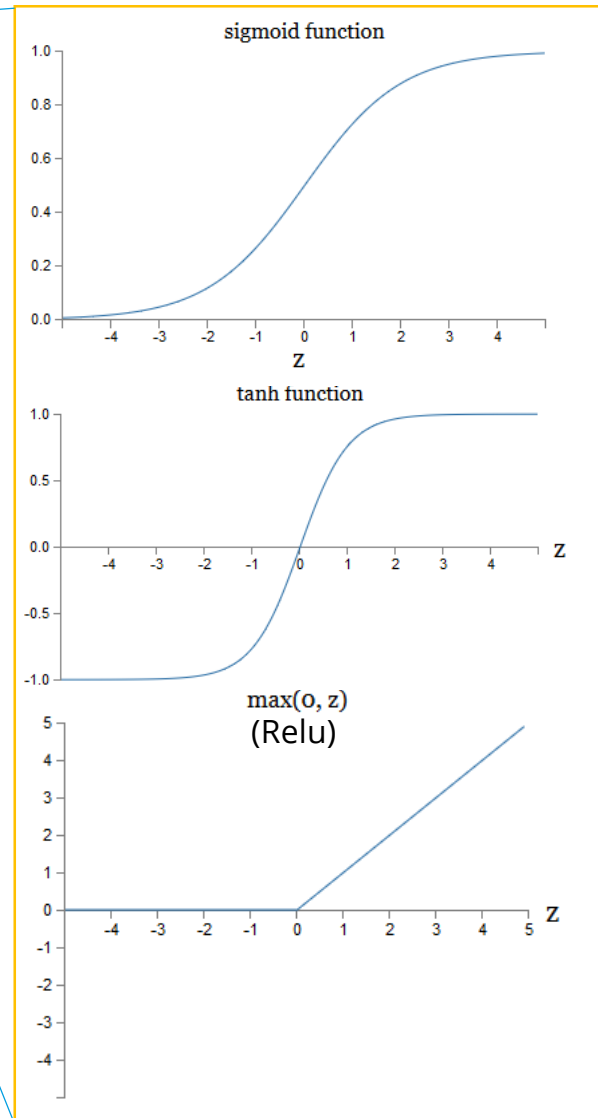
$J_{total}(output)$



For all weights:  
Sum over all  
training data



## Activation functions





# Multilayer Perceptron

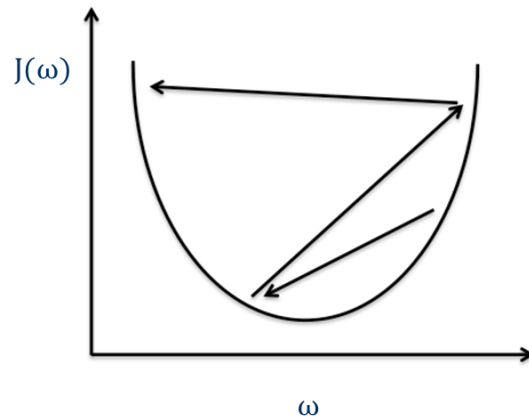
## Backward Pass

Update rule:

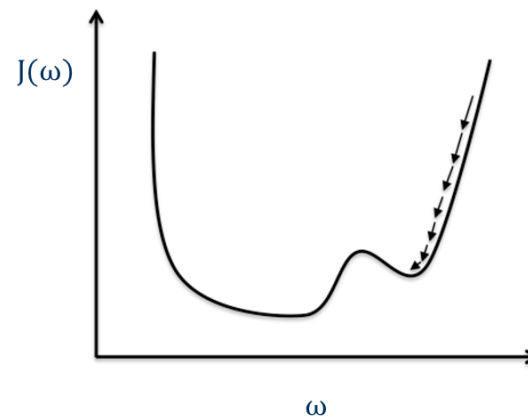
$$w_1(\text{new}) = w_1(\text{old}) - \alpha \cdot \frac{\partial J_{\text{total}}(\text{output})}{\partial w_1}$$

Gradient Descent update

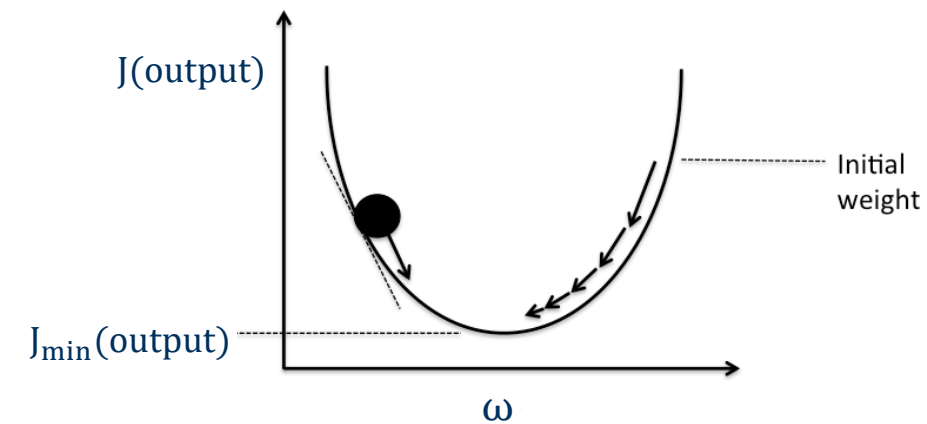
Learning rate



Large learning rate: Overshooting



Small learning rate: Many iterations until convergence and trapping in local minima



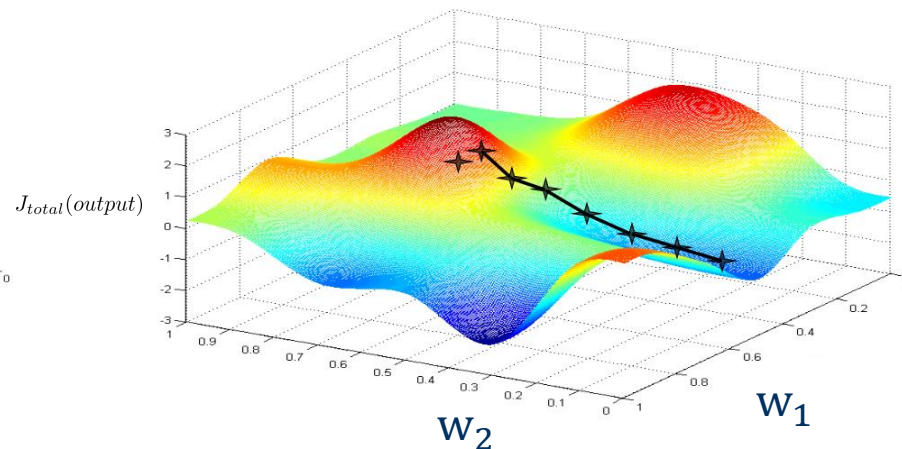
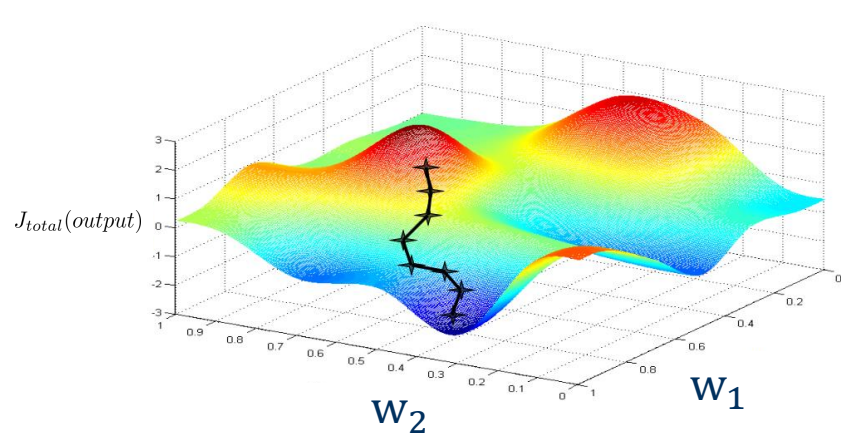
By chain rule we know for the weights that:

$$\frac{\partial \text{input}_{o1,1}}{\partial w_1} \cdot \frac{\partial \text{output}}{\partial \text{input}_{o1,1}} \cdot \frac{\partial J_{\text{total}}(\text{output})}{\partial \text{output}} = \frac{\partial J_{\text{total}}(\text{output})}{\partial w_1}$$

$$J_{\text{total}}(\text{output}) = \frac{1}{2}(\text{target} - \text{output})^2 \quad (\text{MMSE Loss})$$

$$\text{output} = \frac{1}{1 + e^{-\text{input}_{o1,1}}} \quad (\text{Sigmoid})$$

$$\text{input}_{o1,1} = w_1 \cdot \text{output}_{h1,1} + w_2 \cdot \text{output}_{h1,2} + \dots$$



## Initialisation matters!

*"Glorot" init:*

$$\sigma(w) = 1/n$$

$$\mu(w) = 0$$

*"He" init (Relu):*

$$\sigma(w) = 2/n$$

$$\mu(w) = 0$$

*Bias is set to zero.*

n ... Number of Neurons

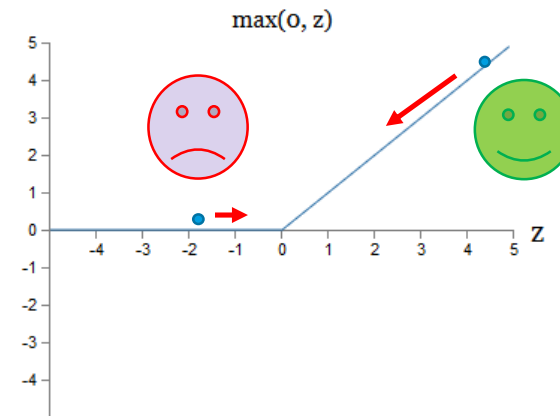
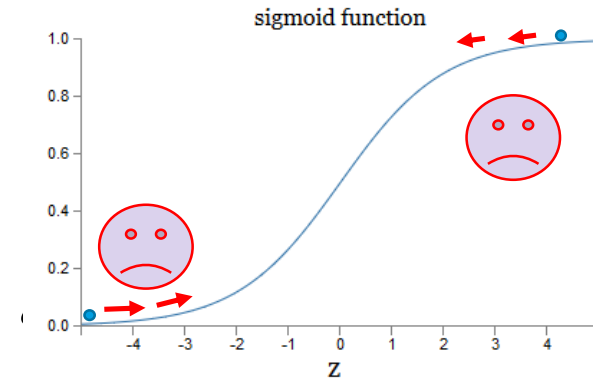
# Neuron Saturation

By chain rule we know our knowledge update is:

$$\frac{\partial \text{input}_{o_{1,1}}}{\partial w_1} \cdot \frac{\partial \text{output}}{\partial \text{input}_{o_{1,1}}} \cdot \frac{\partial J_{\text{total}}(\text{output})}{\partial \text{output}} = \frac{\partial J_{\text{total}}(\text{output})}{\partial w_1}$$

$$\text{output} = \frac{1}{1 + e^{-\text{input}_{o_{1,1}}}} \quad (\text{Sigmoid})$$

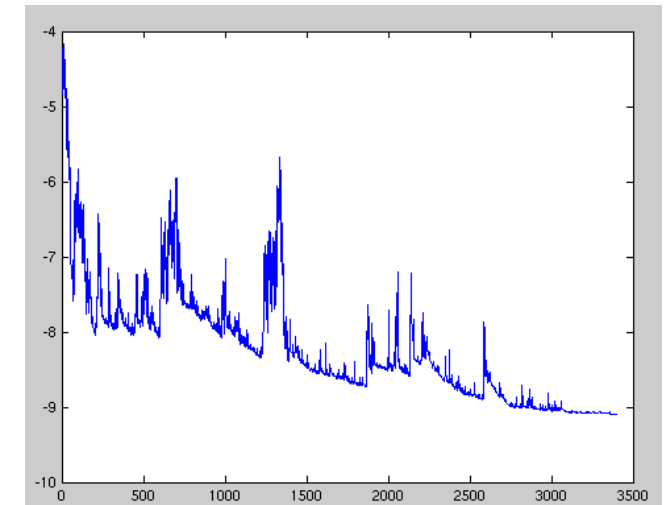
So if the input of our neuron get's very large (Sigmoid & ReLU) or very low (Sigmoid) our **gradient will vanish**. This is sometimes referred as a „**dying neuron**“.



# Optimizer – The “learning” Backbone

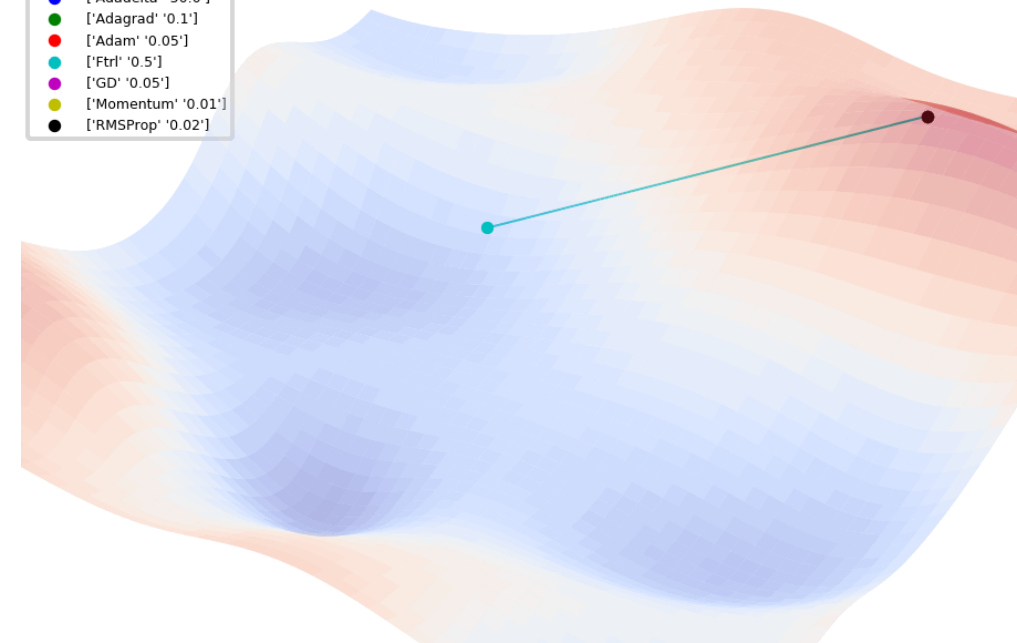
## Update rules

$$w_1(new) = w_1(old) - \alpha \cdot \frac{\partial J_{total}(output)}{\partial w_1}$$



Stochastic Gradient Descent  
(Single-Batch GD)

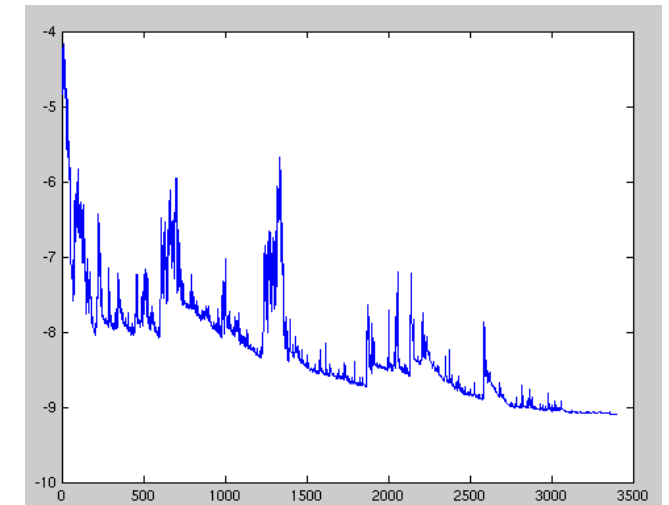
- ['Adadelta' '50.0']
- ['Adagrad' '0.1']
- ['Adam' '0.05']
- ['Ftrl' '0.5']
- ['GD' '0.05']
- ['Momentum' '0.01']
- ['RMSProp' '0.02']



# Optimizer – The “learning” Backbone

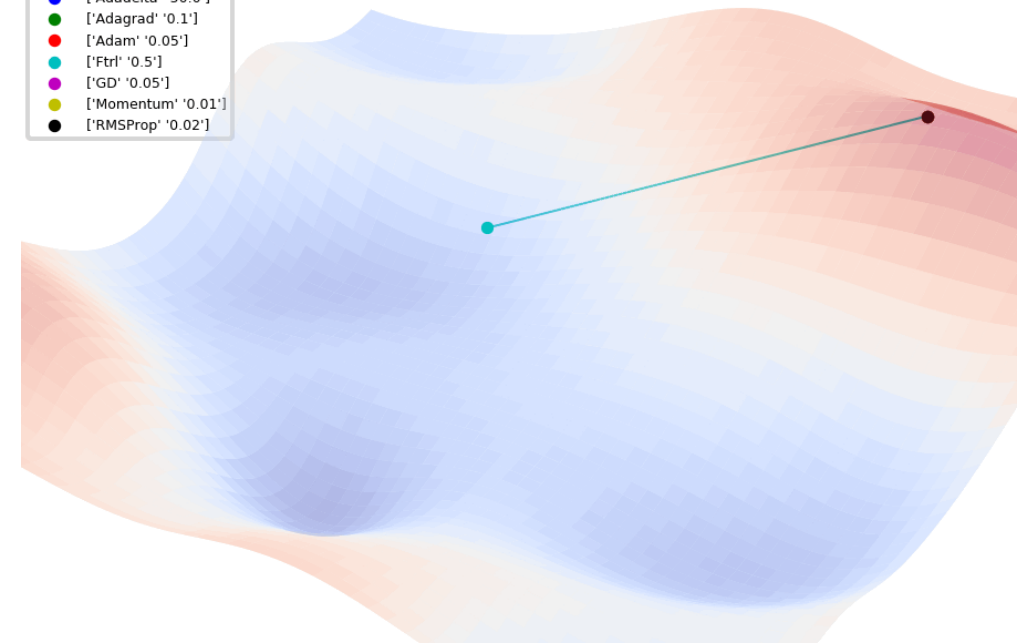
## Update rules

$$w_1(new) = w_1(old) - \alpha \cdot \frac{\partial J_{total}(output)}{\partial w_1}$$



Stochastic Gradient Descent  
(Single-Batch GD)

- ['Adadelta' '50.0']
- ['Adagrad' '0.1']
- ['Adam' '0.05']
- ['Ftrl' '0.5']
- ['GD' '0.05']
- ['Momentum' '0.01']
- ['RMSProp' '0.02']

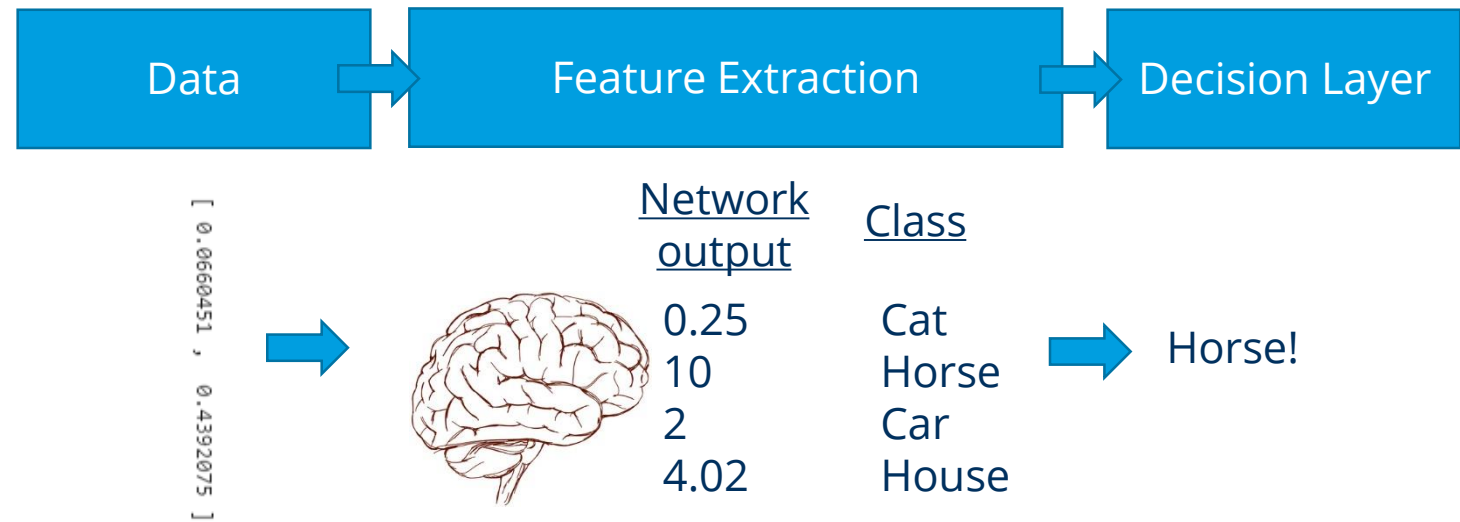




# Multilayer Perceptron

## Decision Layers

How to decide which output to print?  
There are two possible solutions here.



### Argmax(x)

(Trained) Cross-entropie loss „x“	Corresponding Class	<u>Argmax(x)</u>
0.25	Cat	0
10	Horse	10
2	Car	0
4.02	House	0

Naive approach. Just take the maximum of the output vector. → Downside: **No interpretability**

### Softmax(x)

(Trained) Log-Loss „x“	Class	<u>Softmax(x)</u>	<u>Argmax(Softmax(x))</u>
0.25	Cat	0.0105	→ Horse!
10	Horse	0.8	
2	Car	0.085	
4.02	House	0.1	

Better: Interpret output as „pseudo“ probability first. The output sums up to 1 now. Then argmax! Typically used together with a Log-loss or Crossentropy

# 3. Exercise

Let's train our first classifier from scratch!