# VL Deep Learning for Natural Language Processing

2. Neural Networks

*Prof. Dr. Ralf Krestel*

*AG Information Profiling and Retrieval*

# Summary of Previous Session

- Organization of DL4NLP lecture
  - Weekly exercises
  - Three assignments (required for exam participation; gets bonus points)

- Deep learning has revolutionized machine learning
  - No feature engineering
  - Large datasets necessary
  - Works very well for „perception" tasks
  - DL is also a hype: it cannot solve all problems

- Machine learning learns from data (un-)supervised
  - Most methods are based on statistics
  - Can be seen as transforming input to output
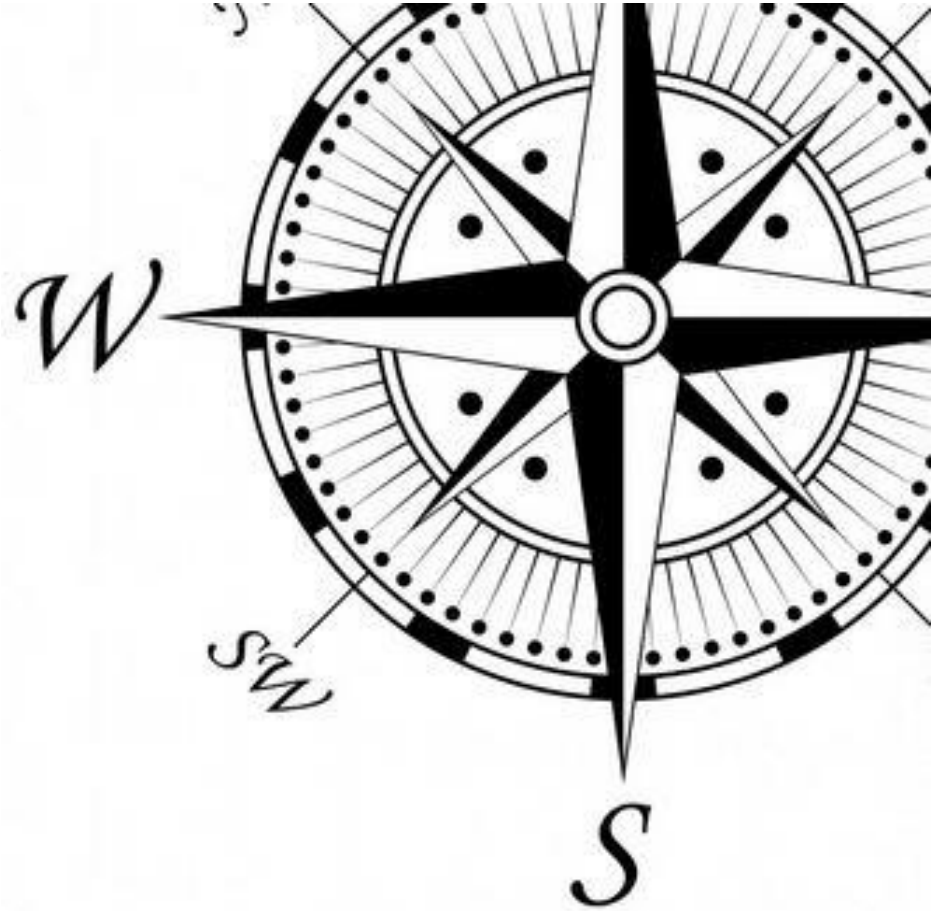
# Learning Goals for this Chapter

- Know how the perceptron works
- Explain the need for multiple layers
- Understand gradient-based optimization
- Describe the components of deep neural networks
- Understand and apply the backpropagation algorithm
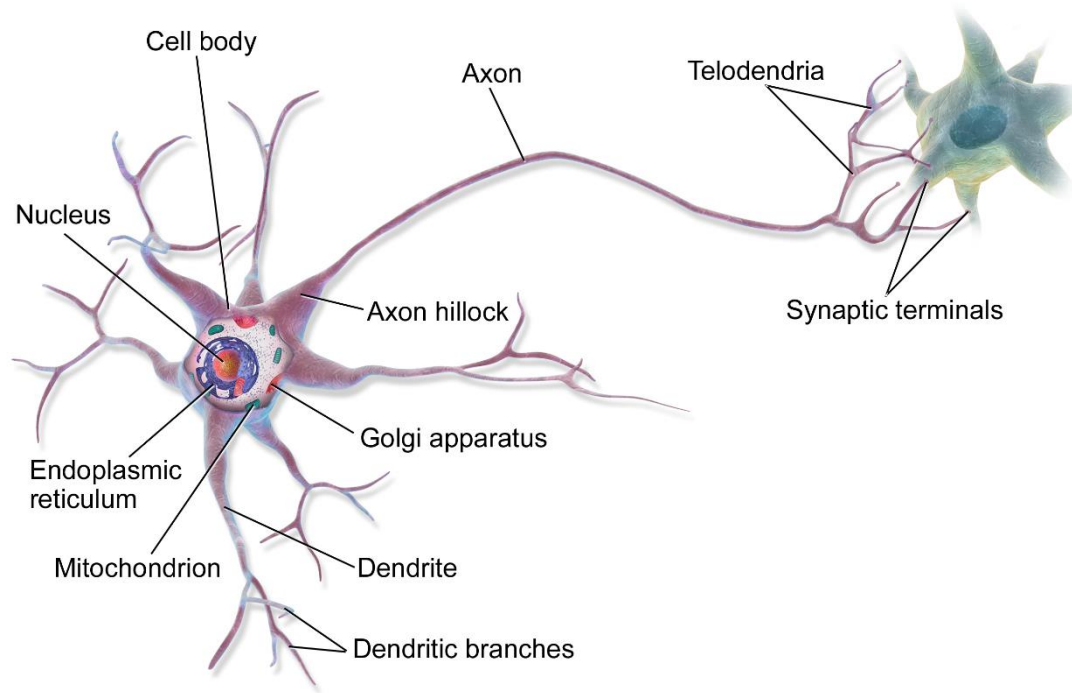
- Relevant chapters
  - P2, P3
  - S3 (2021) https://www.youtube.com/watch?v=X0Jw4kgaFlg

# Topics Today

# This is in your Brain

Cell body

Axon

Telodendria

Nucleus

Axon hillock

Synaptic terminals

Endoplasmic
reticulum

Golgi apparatus

Mitochondrion

Dendrite

Dendritic branches

Output value

y

a

Non-linear transform

σ

z

Weighted sum

Σ

Weights

bias

$w_1$ $w_2$ $w_3$ b

Input layer

$x_1$ $x_2$ $x_3$ +1

# Neural Unit

- Take weighted sum of inputs, plus a bias:

$$z = b + \sum_i w_i x_i$$

- Or in vector notation:

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

- Instead of just using z, we'll apply a nonlinear activation function f:
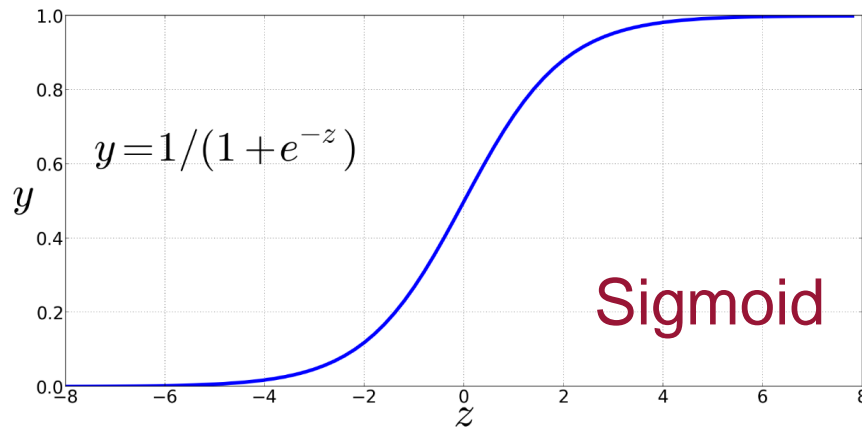
$$y = a = f(z)$$

# Non-Linear Activation Function: Sigmoid

- E.g. sigmoid (aka logistic) function:

$$y = s(z) = \frac{1}{1 + e^{-z}}$$

$$y = 1/(1 + e^{-z})$$

Sigmoid

Final function the neural unit is computing:

$$y = s(\boldsymbol{w} \cdot \boldsymbol{x} + b) = \frac{1}{1 + \exp(-(\boldsymbol{w} \cdot \boldsymbol{x} + b))}$$

# An Example

- Suppose a unit has:

$$\mathbf{w} = [0.2, 0.3, 0.9]$$
$$b = 0.5$$

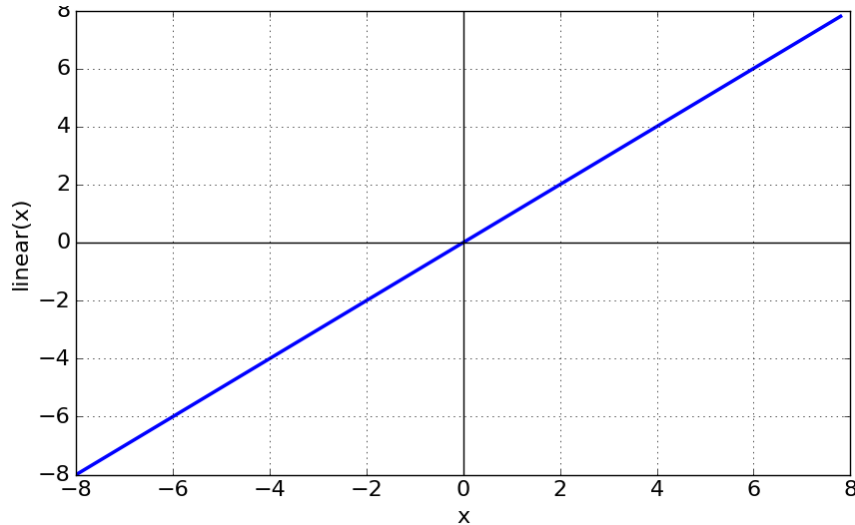- What happens with input **x**:

$$\mathbf{x} = [0.5, 0.6, 0.1]$$
$$y = ?$$

$$y = sigmoid(\mathbf{w} \cdot \mathbf{x} + b)$$
$$= \frac{1}{1 + \exp(-(\mathbf{w} \cdot \mathbf{x} + b))}$$
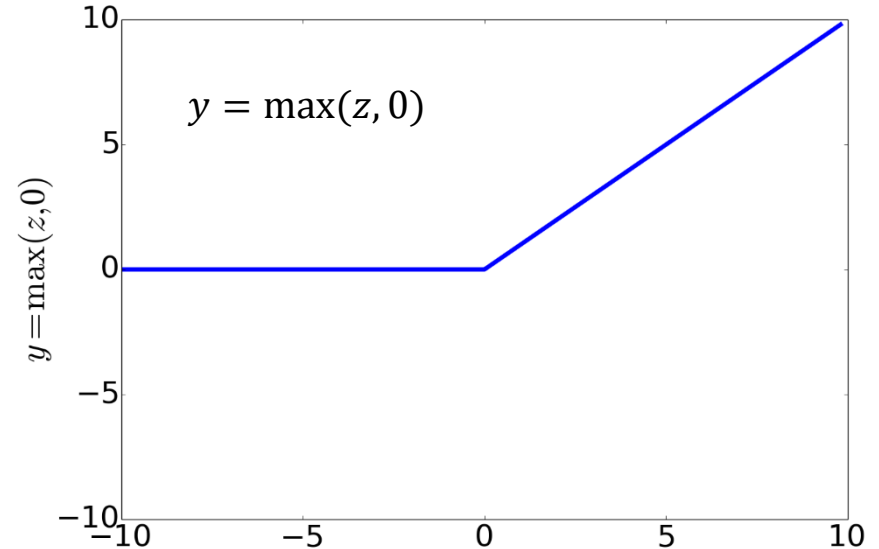$$= \frac{1}{1 + \exp(-(0.2 * 0.5 + 0.3 * 0.6 + 0.9 * 0.1 + 0.5))}$$
$$= \frac{1}{1 + \exp(-0.87)} = 0.70$$
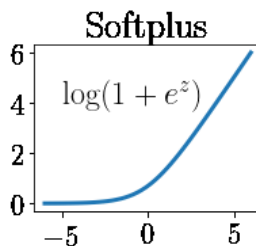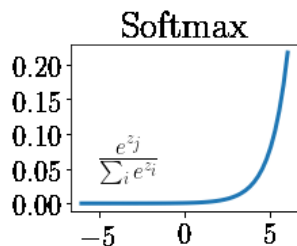
# Activation Functions Besides Sigmoid

Most Common:

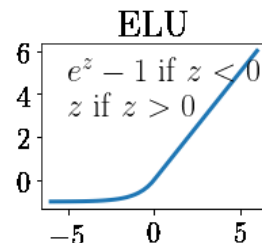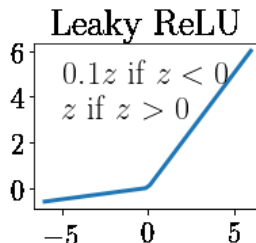

$y = \max(z, 0)$

Linear

ReLU
Rectified Linear Unit

# Most Common Non-Linear Activation Functions

# Topics Today

# Simple Linear Classification

- Given a data point $x = \begin{pmatrix} x_1 \\ x_2 \end{pmatrix}$, predict to which class $\hat{y} = \{-1, 1\}$ this sample belongs
- Model $M(b, w_1, w_2) = M_{\boldsymbol{\theta}}$
  - $M_{\boldsymbol{\theta}}(x) = f(x; \boldsymbol{\theta}) = \hat{y}$
  - $\hat{y} = \begin{cases} 1 & if\ b + w_1 x_1 + w_2 x_2 > 0 \\ -1 & otherwise \end{cases}$
  - Add 1 to all $x$ for simplification:
    - $x = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_0 \\ w_1 \\ w_2 \end{bmatrix}$
    - $\hat{y} = \begin{cases} 1 & if\ \sum_j w_j x_j > 0 \\ -1 & otherwise \end{cases}$
  - $\hat{y} = \text{sign}(\mathbf{w}^T x)$



$w_1 x_1 + w_2 x_2 > -w_o$

$w_1 x_1 + w_2 x_2 \leq -w_o$

# The Perceptron

- Model:

$$\widehat{\boldsymbol{y}} = \text{sign}(\boldsymbol{w}^T \boldsymbol{x})$$

- Activation function of output layer:

$$sign(\boldsymbol{w}^T \boldsymbol{x})$$

# Loss Function

$$L_{0/1}(\hat{y}, y) = \begin{cases} 0, & \hat{y} = y \\ 1, & \hat{y} \neq y \end{cases}$$

$$L_{0/1}(\boldsymbol{w}) = \frac{1}{N} \sum_{i=1}^{N} [y_j == sign(\boldsymbol{w}^T \boldsymbol{x}_i)]$$

$$L_p(\boldsymbol{w}) = \frac{1}{N} \sum_{i=1}^{N} \max(0, -y_i \boldsymbol{w}^T \boldsymbol{x}_i)$$

$$L_p(\boldsymbol{w}) = \sum_{\substack{i=missclassified \\ Samples}} -y_i \boldsymbol{w}^T \boldsymbol{x}_i$$

- Optimization
  - Gradient-based $\rightarrow \boldsymbol{w} = \boldsymbol{w} + \eta y_i \boldsymbol{x}_i$
  - For missclassified sample $x_i$
  - $\frac{\partial L(\boldsymbol{w})}{\partial \boldsymbol{w}_j} = -y_i x_{ij}$ $\qquad \nabla L(\boldsymbol{w}) = -y_i \boldsymbol{x}_i$



http://ttsuchi.github.io/2015/08/26/perceptron.html

# Learning Rate for the Perceptron



- $\boldsymbol{w} = \boldsymbol{w} + \eta y_i \boldsymbol{x}_i$
- Different learning rates
  - $\eta = 0.2, \eta = 0.5, \eta = 1.0$

# Perceptron Training I

- Postive samples
  - $x_0 = \begin{bmatrix} 1 \\ 3 \end{bmatrix}$ $x_1 = \begin{bmatrix} 5 \\ 4 \end{bmatrix}$
- Negative samples
  - $x_2 = \begin{bmatrix} 6 \\ 3 \end{bmatrix}$ $x_3 = \begin{bmatrix} 3 \\ 1 \end{bmatrix}$
- Initial weights
  - $w_0 = 0; w_1 = -1; w_2 = 1$
- Extended data points

$$x_0 = \begin{bmatrix} 1 \\ 1 \\ 3 \end{bmatrix} \quad x_1 = \begin{bmatrix} 1 \\ 5 \\ 4 \end{bmatrix} \quad x_2 = \begin{bmatrix} 1 \\ 6 \\ 3 \end{bmatrix} \quad x_3 = \begin{bmatrix} 1 \\ 3 \\ 1 \end{bmatrix} \qquad y = \begin{bmatrix} 1 \\ 1 \\ -1 \\ -1 \end{bmatrix} \quad w = \begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix}$$

- Loss function

$$L = \max(0, -y_i w^T x_i) \qquad \nabla L = \begin{cases} 0 & if \ y_i w^T x_i > 0 \\ -y_i x_i & otherwise \end{cases}$$

# Perceptron Training II

- $\boldsymbol{x}_0 = \begin{bmatrix} 1 \\ 1 \\ 3 \end{bmatrix}$ $\boldsymbol{x}_1 = \begin{bmatrix} 1 \\ 5 \\ 4 \end{bmatrix}$ $\boldsymbol{x}_2 = \begin{bmatrix} 1 \\ 6 \\ 3 \end{bmatrix}$ $\boldsymbol{x}_3 = \begin{bmatrix} 1 \\ 3 \\ 1 \end{bmatrix}$ $\boldsymbol{y} = \begin{bmatrix} 1 \\ 1 \\ -1 \\ -1 \end{bmatrix}$ $\boldsymbol{w} = \begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix}$

- $L_i = \max(0, -y_i \boldsymbol{w}^T \boldsymbol{x}_i)$ $\qquad \nabla L_i = \begin{cases} 0 & if\ y_i \boldsymbol{w}^T \boldsymbol{x}_i > 0 \\ -y_i \boldsymbol{x_i} & otherwise \end{cases}$ $\quad \eta = 0.5$

- $L_{x_0} = \max(0, -2) = 0$

- $L_{x_1} = \max(0, 1) = 1 \longrightarrow \boldsymbol{w} = \boldsymbol{w} + \eta y_i \boldsymbol{x_i} = \begin{bmatrix} 0 \\ -1 \\ 1 \end{bmatrix} + 0.5 \cdot 1 \cdot \begin{bmatrix} 1 \\ 5 \\ 4 \end{bmatrix} = \begin{bmatrix} 0.5 \\ 1.5 \\ 3 \end{bmatrix}$

- $L_{x_2} = \max(0, 18.5) = 18.5 \longrightarrow \boldsymbol{w} = \boldsymbol{w} - \eta y_i \boldsymbol{x_i} = \begin{bmatrix} 0.5 \\ 1.5 \\ 3 \end{bmatrix} - 0.5 \cdot 1 \begin{bmatrix} 1 \\ 6 \\ 3 \end{bmatrix} = \begin{bmatrix} 0 \\ -1.5 \\ 1.5 \end{bmatrix}$

- $L_{x_3} = \max(0, -3) = 0 \qquad L_{x_0} = \max(0, -3) = 0 \longrightarrow \ ...$

# Perceptron

- How do you need to set the weights of a perceptron to seperate these points:

Class +1: $a = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$ $b = \begin{bmatrix} 0 \\ 1 \end{bmatrix}$ $c = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$

Class -1: $d = \begin{bmatrix} 0 \\ 0 \end{bmatrix}$

- Which function does this perceptron then compute?



**Start** **5** **4** **3** **2** **1** **End**

# History of the Perceptron

- Introduced by Frank Rosenblatt in **1958** [R58]

- Opponent: Marvin Minsky
  - Showed in **1969** that **XOR-Problem** cannot be solved with a perceptron [MP69]
  - The assumption that an extended perceptron could solve the problem proofed to be wrong.
    - AI winter

- Rumelhart and McClelland developed the multi-layer perceptron in **1986**
  - Two-layer perceptrons can represent all Boolean functions
    - Including XOR
  - Hard to train

# The XOR Problem

- Can neural units compute simple functions of input?

|        | AND |     |        | OR  |     |        | XOR |     |
| ------ | --- | --- | ------ | --- | --- | ------ | --- | --- |
| x1     | x2  | y   | x1     | x2  | y   | x1     | x2  | y   |
| 0      | 0   | 0   | 0      | 0   | 0   | 0      | 0   | 0   |
| 0      | 1   | 0   | 0      | 1   | 1   | 0      | 1   | 1   |
| 1      | 0   | 0   | 1      | 0   | 1   | 1      | 0   | 1   |
| 1      | 1   | 1   | 1      | 1   | 1   | 1      | 1   | 0   |

- Only a certain class of functions can be computed!
  -> Linear functions!

# Perceptrons are Linear Classifiers

- Perceptron equation given $x_1$ and $x_2$, is the equation of a line

$$w_1 x_1 + w_2 x_2 + b = 0$$

  - in standard linear format:  $x_2 = (-w_1/w_2)x_1 + (-b/w_2)$

- This line acts as a decision boundary
  - 0 if input is on one side of the line
  - 1 if on the other side of the line

# Decision Boundaries



a)  $x_1$ AND $x_2$

b)  $x_1$ OR $x_2$

c)  $x_1$ XOR $x_2$

XOR is not a linearly separable function!

# Solution to the XOR Problem

- XOR **can't** be calculated by a single perceptron
- XOR **can** be calculated by a layered network of units.

XOR

| x1 | x2 | y |
|----|----|---|
| 0  | 0  | 0 |
| 0  | 1  | 1 |
| 1  | 0  | 1 |
| 1  | 1  | 0 |

ReLU $y_1$

ReLU $h_1$ $h_2$

1    -2    0

1    1    1    1    0    -1

+1

Hidden layer

$x_1$    $x_2$

+1

# The Hidden Representation h



a) The original $x$ space

b) The new (linearly separable) $h$ space

(With learning:  hidden layers will learn to form useful representations)

# Topics Today

1. Neural Network Unit
2. The XOR Problem
3. **Feedforward Neural Networks**
4. Gradient-Based Optimization
5. Backprop(agation Algorithm)
6. Summary

# Multilayer Perceptron

- Multilayer perceptrons (MLP) or multilayer neural networks or feedforward neural networks

    - Generalization/Overfitting

    - High-dimensional
        $\rightarrow$ Needs a lot of training examples

    - Hard to optimize

    - **But**: With two layers, all Boolean functions can be modeled!

# Multilayer Neural Networks

- Activation function: sigmoid or tanh



- Recap: compact representation
  - Omit bias/intercept terms
  - Node = input tensor or result of an activation function
  - Edge = parameters, which map input layer to output layer

# Deep Feedforward Networks

- Activation function: relu
  - **Re**ctified **l**inear **u**nit



- XOR-Problem mit Deep Feedforward Networks:
  - $f(\boldsymbol{x}; \boldsymbol{W}, \boldsymbol{w}) = \boldsymbol{w}^T \max(0, \boldsymbol{W}^T \boldsymbol{x})$

# XOR-Problem: Deep Feedforward Network

- XOR-Problem with deep feedforward networks:
  - $f(x; W, w) = w'^T \max(0, W'^T x + c) + b$
  - $W' = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$ $c = \begin{bmatrix} 0 \\ -1 \end{bmatrix}$ $w' = \begin{bmatrix} 1 \\ -2 \end{bmatrix}$ $b = 0$
  - $W = \begin{bmatrix} 0 & -1 \\ 1 & 1 \\ 1 & 1 \end{bmatrix}$ $w = \begin{bmatrix} 0 \\ 1 \\ -2 \end{bmatrix}$ $X = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{bmatrix}$

**Design matrix: One sample per row**



- Batch processing
  - Compute the output of the network for all four samples simultaneously
    - Forward pass
  - $XW = \begin{bmatrix} 0 & -1 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$ $H' = \max(0, XW) = \begin{bmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 0 \\ 2 & 1 \end{bmatrix}$ $H = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 2 & 1 \end{bmatrix}$ $Hw = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$

# Binary Logistic Regression as a 1-layer Network

(we don't count the input layer in counting layers!)

Output layer
(σ node)

σ

$y = \sigma(w \cdot x + b)$

(y is a scalar)

w

(vector)

$w_1$

$w_n$

b

(scalar)

Input layer
vector x

$x_1$

$x_n$

+1

# Multinomial Logistic Regression as a 1-layer Network

Fully connected single layer network



$y = \text{softmax}(Wx + b)$

(y is a vector)

Output layer
(**softmax nodes**)

W

(matrix)

b

(vector)

Input layer
scalars

# Softmax: a Generalization of Sigmoid

- For a vector $z$ of dimensionality $k$, the softmax is:

$$\text{softmax}(z) = \left[ \frac{\exp(z_1)}{\sum_{i=1}^{k} \exp(z_i)}, \frac{\exp(z_2)}{\sum_{i=1}^{k} \exp(z_i)}, ..., \frac{\exp(z_k)}{\sum_{i=1}^{k} \exp(z_i)} \right]$$

- Example:

$$\text{softmax}(z_i) = \frac{\exp(z_i)}{\sum_{j=1}^{k} \exp(z_j)} \quad 1 \le i \le k$$

$$z = [0.6, 1.1, -1.5, 1.2, 3.2, -1.1]$$

$$\text{softmax}(z) = [0.055, 0.090, 0.006, 0.099, 0.74, 0.010]$$

# Two-Layer Network with Scalar Output

Output layer
(σ node)

hidden units
(σ node)

Input layer
(vector)

$y = \sigma(z)$    (y is a scalar)

$z = Uh$

U

$$h = \boldsymbol{\sigma}(\boldsymbol{Wx + b})$$

Could be ReLU or tanh

$W_{ji}$

W

b

$x_1$   $x_n$   +1

# Two-Layer Network with Softmax Output

Output layer
(σ node)

hidden units
(σ node)

Input layer
(vector)



$y = \text{softmax}(z)$

$z = Uh$       (y is a vector)

$h = \boldsymbol{\sigma}(\boldsymbol{W}\boldsymbol{x} + \boldsymbol{b})$

Could be ReLU or tanh

U

W

b

$x_1$

$x_n$

+1

# Multi-layer Notation



$$y = a^{[2]}$$

$$a^{[2]} = g^{[2]}(z^{[2]}) \quad \text{sigmoid or softmax}$$

$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$

$$a^{[1]} = g^{[1]}(z^{[1]}) \quad \text{ReLU}$$

$$z^{[1]} = W^{[1]}a^{[0]} + b^{[1]}$$

$W^{[2]}$

$b^{[2]}$

$W^{[1]}$

$b^{[1]}$

$a^{[0]}$

$x_1$   $i$   $x_n$   +1

# Multi Layer Notation

$$z^{[1]} = W^{[1]}a^{[0]} + b^{[1]}$$
$$a^{[1]} = g^{[1]}(z^{[1]})$$
$$z^{[2]} = W^{[2]}a^{[1]} + b^{[2]}$$
$$a^{[2]} = g^{[2]}(z^{[2]})$$
$$\hat{y} = a^{[2]}$$

**for** $i$ **in** $1..n$
$$z^{[i]} = W^{[i]}a^{[i-1]} + b^{[i]}$$
$$a^{[i]} = g^{[i]}(z^{[i]})$$
$$\hat{y} = a^{[n]}$$

# Topics Today

# What Happens in a Neural Network?

# What Happens in a Layer?

$$\texttt{output = f(x)= f(dot(W, input) + b)}$$

- W and b are tensors.
  - W are the **weights** or **parameters** of a layer and need to be trained / learned.
  - b can be merged into W by using a constant 1 as pseudo input
- At the beginning, weights are randomly initialized.
  - The network has not learned any meaningful representation or transormation and therefore no good mapping from input to output.

$$w_{0,1}^0 = b_{0,1}^0 \qquad f(x) = relu(x) = \max(0, x)$$

$$f(x) = softmax(x) = \frac{e^x}{\sum e^x}$$

$$P(y = 0|x)$$

$$P(y = 9|x)$$

$$W^0 = \begin{bmatrix} w_{0,1}^0 & \cdots & w_{28,1}^0 \\ \vdots & \ddots & \vdots \\ w_{0,512}^0 & \cdots & w_{28,512}^0 \end{bmatrix}$$

# How to Train a Network?

1. Randomly choose $k$ training samples $x$ (**mini batch**).
   – Alternative:
     o Batch (all training samples)
     o Stochastisic (one training sample)

2. Compute the network's output $\hat{y}$ for input $x$.

3. Compute the **loss** of the network on the batch, i.e. the discrepancy between the predicition $\hat{y}$ and the actual value $y$

$$Loss = L(\hat{y}, y)$$

4. Update the weights in a way that reduces the loss a little

$$f(x) = relu(x) = \max(0, x)$$

$$f(x) = softmax(x) = \frac{e^x}{\sum e^x}$$

# Loss Functions

- Mean absolute error
  - $MAE = \frac{\sum_{i=1}^{n}|y_i - \hat{y}_i|}{n}$
- Mean squared error
  - $MSE = \frac{\sum_{i=1}^{n}|y_i - \hat{y}_i|^2}{n}$
- Binary cross entropy
  - $BCE = -(y_i \log(\hat{y}_i)$
    $+ (1 - y_i) \log(1 - \hat{y}_i))$
- Hinge loss
  - $Hinge = \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1)$
- Categorical cross entropy
  - $CE = -\sum_C \boldsymbol{y}_i \log(\boldsymbol{\hat{y}}_i)$
- Kullback Leibler divergence
  - $KL(P, Q) = \sum_{x \in X} P(x) \log \frac{P(x)}{Q(x)}$

- Which loss functions for
  1. Regression problem
     - Input tensor $\boldsymbol{x} \rightarrow \mathbb{R}$
  2. Binary classification
     - Input tensor $\boldsymbol{x} \rightarrow [0,1]$ (can also be x $\rightarrow [-1,1]$)
  3. (Single-label), multi-class classification
     - Input tensor $\boldsymbol{x} \rightarrow \boldsymbol{y}$, with $y =$binary vector and $|y| = 1$
  4. Multi-label, (multi-class) classification
     - Input tensor $\boldsymbol{x} \rightarrow \boldsymbol{y}$, with $y =$binary vector

# How to Update the Weights I?

- Naive appraoch:
1. Choose a training **sample** randomly
2. Compute the **loss**
3. Choose a weight randomly
4. Update this weight **randomly**, keep all other weights fixed
5. Compute the **loss** again
   - If loss lower, keep weight and go to step 3.
   - If loss higher, go back to step 4.
   - If all weights are updated, go back to step 1.
   - If all samples were used for training, repeat whole process (i times)

- Way too **inefficient**!
  - Theoretically possible
  - Possible: Finding only **local minimum**

# How to Update the Weights II?

- Analytical appraoch:
  1. Compute the derivative / gradient of the loss function $L(W)$
  2. Identify all $W^*$ with $L'(W^*) = 0$
  3. Compute the loss for all $W^*$
  4. Select the $W^*$ for which $L(W^*)$ is minimal
     - Globle minimum

- Not **efficiently** solvable!
  - If more than a few weights are involved
  - Typically: millions of weights!

# Stochastic Gradient Descent (SGD)

- (true) **SGD**
  - Just like naive approach, only updates of weights not random
  - Updates are based on derivative / gradient
  - Stochastic, since samples are chosen randomly from training set



Learning rate (step size)

# Batch and Mini-Batch SGD

- **Batch SGD**
  - Similar to SGD, but instead of using only one sample, compute loss on all training samples
    - Updates of weights much more accurate
    - Computation much more **expensive**
- **Mini-batch SGD**
  - Compromise between looking at all samples and only one sample
  - Simulaneously evaluating a small set of samples
    - Typically 8, 16, 32, 64, 128 or 256
- **Learning rate** is an important (hyper-) parameter
  - Variations:
    - Adaptive learning rate
      - Higher order derivatives (**momentum**)

# The Backpropagation Algorithm

- SGD needs the derivative / gradient of the loss function for each weight
- Typically, a NN consists of many tensor operations

  **f(W1, W2, W3) = a(W1, b(W2, c(W3)))**

  - For each single one it is easy to compute the gradient
- **Chain rule**:

$$\text{f} = u \circ v : V \to \mathbb{R}$$
$$(u \circ v)'(x_0) = u'\big(v(x_0)\big) \cdot v'(x_0)$$

- **Backpropagation Algoritm**
  - Application of chain rule to compute gradient of NN
  - Start with loss at the last (output) layer of the network and compute backwards the proportion that each weight contributed to this loss (backpropagation).
  - Implemented in Keras using symbolic differentiation
    - A gradient function for the chain of derivatives maps network parameter values (weights) to the respective gradients.

# Topics Today

"I think you should be more explicit here in step two."

# Computation of Gradients

- **Forward pass**
  - For feedforward neural nets
    - o Computation of $\hat{y}$ given input $x$
  - During training: Additionally computation of
    - o Error / Loss function $J(\boldsymbol{\theta})$
  - Batch processing possible
    - o Simultaneously computing $J(\boldsymbol{\theta})$ for multiple input samples $X$
- **Backpropagation algorithm** (Rumelhart et al., 1986)
  - short: backprop
  - Propagation of the error back through the network to compute the gradients
  - **Backward pass**
  - Actual learning is done via gradient descent

$x \longrightarrow h \longrightarrow o$

$\theta = x, y, w$

# Computational Graphs

Logistic regression $\hat{y} = \sigma(\boldsymbol{x}^T \boldsymbol{w} + b)$



Fully connected feed forward network

$$\hat{\boldsymbol{y}} = softmax\{\boldsymbol{h}\boldsymbol{W}^{(2)}\}$$
$$= softmax\{\max\{0, \boldsymbol{x}\boldsymbol{W}^{(1)}\}\boldsymbol{W}^{(2)}\}$$

Input: 2d, i.e. 2 features

output: Probabilities for each of the three exclusive classes

# Chain Rule

- Given $g, f: \mathbb{R} \to \mathbb{R}$
  - $y = g(x)$
  - $z = f(g(x)) = f(y)$

- Chain rule

$$\frac{dz}{dx} = \frac{dz}{dy}\frac{dy}{dx}$$

- In case of vectors $\quad g: \mathbb{R}^m \to \mathbb{R}^n \qquad f: \mathbb{R}^n \to \mathbb{R}$
  - $\boldsymbol{x} \in \mathbb{R}^m; \boldsymbol{y} \in \mathbb{R}^n; \boldsymbol{y} = g(\boldsymbol{x}); z = f(\boldsymbol{y})$

$$\frac{\partial z}{\partial x_i} = \sum_{j=1}^{n} \frac{\partial z}{\partial y_j}\frac{\partial y_j}{\partial x_i}$$

$n \times m$ **Jacobi-Matrix**

$$\nabla_{\boldsymbol{x}} z = \left(\frac{\partial \boldsymbol{y}}{\partial \boldsymbol{x}}\right)^T \nabla_{\boldsymbol{y}} z$$

$n \times m$ **Jacobi-Matrix**

# Graphical Representation

- For epoche 1 to k:
  - For each training sample / batch of samples:
    - Forward pass
    - Computation of loss function
    - Backward pass
    - Update of weights (gradient descent)

# The Algorithm: Forward Pass

- Input:
  - Network depth $l$
  - $\boldsymbol{W}^{(i)}, i \in \{1, \ldots, l\}$
  - $\boldsymbol{b}^{(i)}, i \in \{1, \ldots, l\}$
  - $\boldsymbol{x}$ input data
  - $\boldsymbol{y}$ target data
- Output
  - Value of the loss function at position $\boldsymbol{x}$



$\boldsymbol{h}^{(0)} = \boldsymbol{x}$
$\boldsymbol{for}\ k = 1, \ldots, l\ \boldsymbol{do}$
$\quad \boldsymbol{a}^{(k)} = \boldsymbol{b}^{(k)} + \boldsymbol{W}^{(k)} \boldsymbol{h}^{(k-1)}$
$\quad \boldsymbol{h}^{(k)} = f(\boldsymbol{a}^{(k)})$
$\boldsymbol{end\ for}$
$\widehat{\boldsymbol{y}} = \boldsymbol{h}^{(l)}$
$J = L(\widehat{\boldsymbol{y}}, \boldsymbol{y}) + \lambda \Omega(\boldsymbol{\theta})$

Regularization;
$\theta = $ **all weights** $+$ **bias terms**

# The Algorithm: Backward Pass

- Output: The gradients of all activations $\boldsymbol{a}^{(k)}$
- Afterwards, update the weights
  - E.g. using gradient descent

$$\boldsymbol{g} \leftarrow \nabla_{\widehat{\boldsymbol{y}}} J = \nabla_{\widehat{\boldsymbol{y}}} L(\widehat{\boldsymbol{y}}, \boldsymbol{y})$$

$$\boldsymbol{for} \; k = l, l-1, \dots, 1 \; \boldsymbol{do}$$

$$\boldsymbol{g} \leftarrow \nabla_{\boldsymbol{a}^{(k)}} J = \boldsymbol{g} \odot f'(\boldsymbol{a}^{(k)})$$

$$\nabla_{\boldsymbol{b}^{(k)}} J = \boldsymbol{g} + \lambda \nabla_{\boldsymbol{b}^{(k)}} \Omega(\boldsymbol{\theta})$$

$$\nabla_{\boldsymbol{W}^{(k)}} J = \boldsymbol{g} \boldsymbol{h}^{(k-1)T} + \lambda \nabla_{\boldsymbol{W}^{(k)}} \Omega(\boldsymbol{\theta})$$

$$\boldsymbol{g} \leftarrow \nabla_{\boldsymbol{h}^{(k-1)}} J = \boldsymbol{W}^{(k)T} \boldsymbol{g}$$

$$\boldsymbol{end} \; \boldsymbol{for}$$

**Gradient of output layer**

**Gradient before non-linear activation; $\odot$ element-wise multiplication**

**Gradients of weights and bias terms of layer $k$**

**Propagation of the gradients to the activations of the next lower layer**

# Walk-Through Example: Iris Dataset

```
>>> from sklearn.datasets import load_iris
>>> data = load_iris()
>>> list(data.target_names)
['setosa', 'versicolor', 'virginica']
```



**Binary cross-entropy:**

$$L = BCE(\hat{y}, y) = -(y\,log(\hat{y}) + (1-y)\,log(1-\hat{y}))$$

**Cross-entropy:**

$$L = CE(\hat{y}, y) = -\sum_C y_i\,log(\hat{y}_i)$$

**Cross-entropy considering also negative samples:**

$$L = CE(\hat{y}, y) = -\sum_C y_i\,log(\hat{y}_i)\,(1-y_i)\,log(1-\hat{y}_i)$$

# Walk-Through Example: Forward Pass

- Initializing the weights — **Different Initializations possible**

$$\boldsymbol{W}^{(1)} = \begin{bmatrix} 0.1 & -0.2 \\ 0.4 & 0.6 \end{bmatrix} \quad \boldsymbol{b}^{(1)} = \begin{bmatrix} 0.2 \\ -0.3 \end{bmatrix} \quad \boldsymbol{W}^{(2)} = \begin{bmatrix} 0.3 & -0.3 \\ 0.2 & 0.1 \\ 0.3 & -0.2 \end{bmatrix} \quad \boldsymbol{b}^{(2)} = \begin{bmatrix} -0.1 \\ -0.5 \\ -0.4 \end{bmatrix}$$

- First training sample $\boldsymbol{x}^{(1)} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1.8 \\ 0.5 \end{bmatrix}$ $\quad \boldsymbol{y}^{(1)} = \begin{bmatrix} 1.0 \\ 0.0 \\ 0.0 \end{bmatrix}$ $\quad \boldsymbol{h}^{(0)} = \begin{bmatrix} 1.8 \\ 0.5 \end{bmatrix}$

  - $\boldsymbol{a}^{(1)} = \begin{bmatrix} 0.2 \\ -0.3 \end{bmatrix} + \begin{bmatrix} 0.1 & -0.2 \\ 0.4 & 0.6 \end{bmatrix} \begin{bmatrix} 1.8 \\ 0.5 \end{bmatrix} = \begin{bmatrix} 0.28 \\ 0.72 \end{bmatrix}$
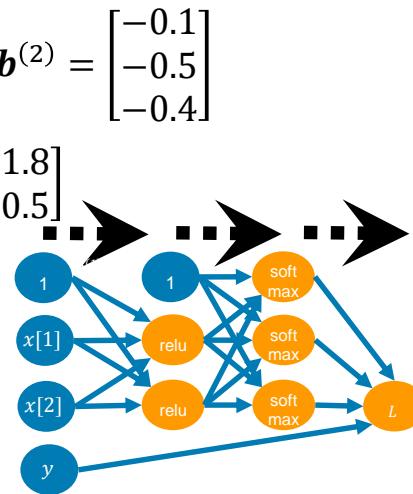
  - $\boldsymbol{h}^{(1)} = relu\left(\begin{bmatrix} 0.28 \\ 0.72 \end{bmatrix}\right) = \begin{bmatrix} 0.28 \\ 0.72 \end{bmatrix}$

  - $\boldsymbol{a}^{(2)} = \begin{bmatrix} -0.1 \\ -0.5 \\ -0.4 \end{bmatrix} + \begin{bmatrix} 0.3 & -0.3 \\ 0.2 & 0.1 \\ 0.3 & -0.2 \end{bmatrix} \begin{bmatrix} 0.28 \\ 0.72 \end{bmatrix} = \begin{bmatrix} -0.23 \\ -0.37 \\ -0.46 \end{bmatrix}$

  - $\boldsymbol{o}^{(2)} = \text{softmax}\left(\begin{bmatrix} -0.23 \\ -0.37 \\ -0.46 \end{bmatrix}\right) = \begin{bmatrix} 0.38 \\ 0.32 \\ 0.30 \end{bmatrix} = \widehat{\boldsymbol{y}}$ 

    **Considers also negative samples**

  - $J = CE\left(\begin{bmatrix} 0.38 \\ 0.32 \\ 0.30 \end{bmatrix}, \begin{bmatrix} 1.0 \\ 0.0 \\ 0.0 \end{bmatrix}\right) = -(\log(0.38) + \log(1 - 0.32) + \log(1 - 0.30)) = 0.74$



$$\boldsymbol{h}^{(0)} = \boldsymbol{x}$$
$$for\ k = 1, \dots, l\ \boldsymbol{do}$$
$$\quad \boldsymbol{a}^{(k)} = \boldsymbol{b}^{(k)} + \boldsymbol{W}^{(k)} \boldsymbol{h}^{(k-1)}$$
$$\quad \boldsymbol{h}^{(k)} = f(\boldsymbol{a}^{(k)})$$
$$end\ for$$
$$\widehat{\boldsymbol{y}} = \boldsymbol{h}^{(l)} = \boldsymbol{o}^{(l)}$$
$$J = L(\widehat{\boldsymbol{y}}, \boldsymbol{y}) + \lambda \Omega(\boldsymbol{\theta})$$

# Walk-Through Example: Backward Pass I

- Loss for $x^{(1)} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 1.8 \\ 0.5 \end{bmatrix}$

$$L = CE\left(\begin{bmatrix} 0.38 \\ 0.32 \\ 0.30 \end{bmatrix}, \begin{bmatrix} 1.0 \\ 0.0 \\ 0.0 \end{bmatrix}\right) = 0.74$$

- Reminder: chain rule

$$\frac{\partial L}{\partial W^{(1)}} = \frac{\partial L}{\partial o_{out}} \frac{\partial o_{out}}{\partial o_{in}} \frac{\partial o_{in}}{\partial h_{out}} \frac{\partial h_{out}}{\partial h_{in}} \frac{\partial h_{in}}{\partial W^{(1)}}$$

- Gradient of loss function:

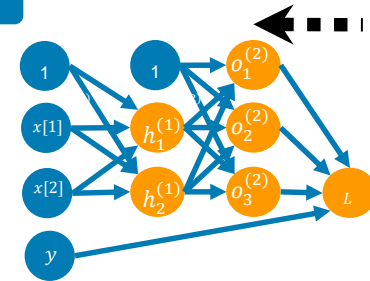- $\dfrac{\partial L}{\partial o_{out}} = \begin{bmatrix} \frac{\partial L}{\partial o_{1,out}} \\ \frac{\partial L}{\partial o_{2,out}} \\ \frac{\partial L}{\partial o_{3,out}} \end{bmatrix} = \begin{bmatrix} -1 \cdot (1 \cdot \frac{1}{0.38} + (1-1) \cdot \frac{1}{1-0.38}) \\ -1 \cdot (0 \cdot \frac{1}{0.32} + (1-0) \cdot \frac{1}{1-0.32}) \\ -1 \cdot (0 \cdot \frac{1}{0.30} + (1-0) \cdot \frac{1}{1-0.30}) \end{bmatrix} = \begin{bmatrix} -2.63 \\ -1.47 \\ -1.43 \end{bmatrix}$

**Partial derivative of cross-entropy:**

$$\frac{\partial L}{\partial \hat{y}_i} = -1 \cdot (y_i \frac{1}{\hat{y}_i} + (1-y_i) \frac{1}{1-\hat{y}_i})$$

$x_1 \rightarrow h^{(1)} \rightarrow o^{(2)} \rightarrow L$

$w^{(1)} \quad h_{in}^{(1)} \quad h_{out}^{(1)} \quad w^{(2)} \quad o_{in}^{(2)} \quad o_{out}^{(2)} \quad L$

**Hadamard product**

$g \leftarrow \nabla_{\hat{y}} J = \nabla_{\hat{y}} L(\hat{y}, y)$
$for\ k = l, l-1, \ldots, 1\ do$
$\quad g \leftarrow \nabla_{a^{(k)}} J = g \odot f'(a^{(k)})$
$\quad \nabla_{b^{(k)}} J = g + \lambda \nabla_{b^{(k)}} \Omega(\theta)$
$\quad \nabla_{W^{(k)}} J = g h^{(k-1)T} + \lambda \nabla_{W^{(k)}} \Omega(\theta)$
$\quad g \leftarrow \nabla_{h^{(k-1)}} J = W^{(k)T} g$
$end\ for$

# Walk-Through Example: Backward Pass II



- Gradient of the output of the output layer

$$\frac{\partial o_{out}}{\partial o_{in}} = \begin{bmatrix} \dfrac{\partial o_{1,out}^{(2)}}{\partial o_{1,in}^{(2)}} \\ \dfrac{\partial o_{2,out}^{(2)}}{\partial o_{2,in}^{(2)}} \\ \dfrac{\partial o_{3,out}^{(2)}}{\partial o_{3,in}^{(2)}} \end{bmatrix} = \begin{bmatrix} \dfrac{e^{-0.23} \cdot (e^{-0.37} + e^{-0.46})}{(e^{-0.23} + e^{-0.37} + e^{-0.46})^2} \\ \dfrac{e^{-0.37} \cdot (e^{-0.23} + e^{-0.46})}{(e^{-0.23} + e^{-0.37} + e^{-0.46})^2} \\ \dfrac{e^{-0.46} \cdot (e^{-0.23} + e^{-0.37})}{(e^{-0.23} + e^{-0.37} + e^{-0.46})^2} \end{bmatrix} = \begin{bmatrix} 0.23 \\ 0.22 \\ 0.21 \end{bmatrix}$$

**Partrial derivative of softmax:**

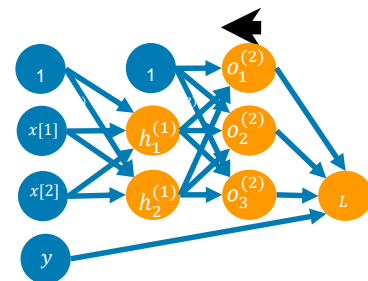$$\frac{\partial o_{i,out}^{(2)}}{\partial o_{i,in}^{(2)}} = \frac{e^{o_{i,in}^{(2)}} \cdot \sum_{j \neq i} e^{o_{j,in}^{(2)}}}{\left( \sum_j e^{o_{j,in}^{(2)}} \right)^2}$$

$$g \leftarrow \nabla_{\hat{y}} J = \nabla_{\hat{y}} L(\hat{y}, y)$$
$$\mathbf{for}\ k = l, l-1, \dots, 1\ \mathbf{do}$$
$$\quad g \leftarrow \nabla_{a^{(k)}} J = g \odot f'(a^{(k)})$$
$$\quad \nabla_{b^{(k)}} J = g + \lambda \nabla_{b^{(k)}} \Omega(\theta)$$
$$\quad \nabla_{W^{(k)}} J = g h^{(k-1)T} + \lambda \nabla_{W^{(k)}} \Omega(\theta)$$
$$\quad g \leftarrow \nabla_{h^{(k-1)}} J = W^{(k)T} g$$
$$\mathbf{end\ for}$$

- Gradient of the input of the output layer with respect to weights $W^{(2)}$

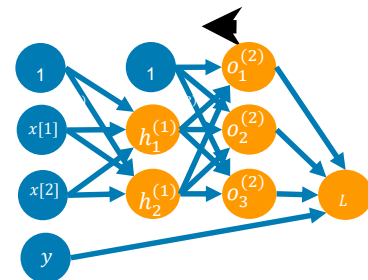$$\frac{\partial o_{1,in}^{(2)}}{\partial W_{1,1}^{(2)}} = \frac{\partial o_{2,in}^{(2)}}{\partial W_{1,2}^{(2)}} = \frac{\partial o_{3,in}^{(2)}}{\partial W_{1,3}^{(2)}} = 0.28 \qquad \frac{\partial o_{1,in}^{(2)}}{\partial W_{2,1}^{(2)}} = \frac{\partial o_{2,in}^{(2)}}{\partial W_{2,2}^{(2)}} = \frac{\partial o_{3,in}^{(2)}}{\partial W_{2,3}^{(2)}} = 0.72$$

**Partial derivative:**

$$\frac{\partial o_{1,in}^{(2)}}{\partial W_{1,1}^{(2)}} = \frac{\partial (h_1^{(1)} W_{1,1}^{(2)} + h_2^{(1)} W_{2,1}^{(2)} + b_1^{(2)})}{\partial W_{1,1}^{(2)}} = h_{1,out}^{(1)}$$

$$g \leftarrow \nabla_{\widehat{y}} J = \nabla_{\widehat{y}} L(\widehat{y}, y)$$
$$\textbf{for } k = l, l-1, \dots, 1 \textbf{ do}$$
$$\qquad g \leftarrow \nabla_{a^{(k)}} J = g \odot f'(a^{(k)})$$
$$\qquad \nabla_{b^{(k)}} J = g + \lambda \nabla_{b^{(k)}} \Omega(\theta)$$
$$\qquad \nabla_{W^{(k)}} J = g h^{(k-1)T} + \lambda \nabla_{W^{(k)}} \Omega(\theta)$$
$$\qquad g \leftarrow \nabla_{h^{(k-1)}} J = W^{(k)T} g$$
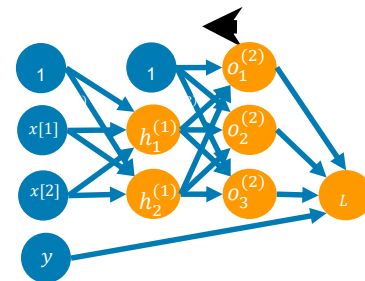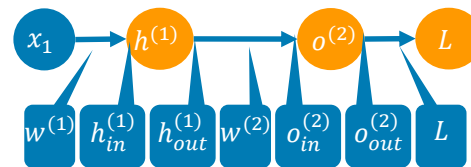$$\textbf{end for}$$

$$\frac{\partial L}{\partial W^{(2)}} = \frac{\partial L}{\partial o_{out}} \frac{\partial o_{out}}{\partial o_{in}} \frac{\partial o_{in}}{\partial W^{(2)}} = \begin{bmatrix} \frac{\partial L}{\partial o_{1,out}} \frac{\partial o_{1,out}}{\partial o_{1,in}} \frac{\partial o_{1,in}}{\partial W_{1,1}^{(2)}} & \frac{\partial L}{\partial o_{1,out}} \frac{\partial o_{1,out}}{\partial o_{1,in}} \frac{\partial o_{1,in}}{\partial W_{2,1}^{(2)}} \\ \frac{\partial L}{\partial o_{2,out}} \frac{\partial o_{2,out}}{\partial o_{2,in}} \frac{\partial o_{2,in}}{\partial W_{1,2}^{(2)}} & \frac{\partial L}{\partial o_{2,out}} \frac{\partial o_{2,out}}{\partial o_{2,in}} \frac{\partial o_{2,in}}{\partial W_{2,2}^{(2)}} \\ \frac{\partial L}{\partial o_{3,out}} \frac{\partial o_{3,out}}{\partial o_{13in}} \frac{\partial o_{3,in}}{\partial W_{1,3}^{(2)}} & \frac{\partial L}{\partial o_{3,out}} \frac{\partial o_{3,out}}{\partial o_{13in}} \frac{\partial o_{3,in}}{\partial W_{2,3}^{(2)}} \end{bmatrix}$$

$$\frac{\partial L}{\partial W^{(2)}} = \begin{bmatrix} -2.63 \cdot 0.23 \cdot 0.28 & -2.63 \cdot 0.23 \cdot 0.72 \\ -1.47 \cdot 0.22 \cdot 0.28 & -1.47 \cdot 0.22 \cdot 0.72 \\ -1.43 \cdot 0.21 \cdot 0.28 & -1.43 \cdot 0.21 \cdot 0.72 \end{bmatrix}$$

- Gradient descent with learning rate $\lambda = 0.5$ results in new weights:

$$W^{(2)} = \begin{bmatrix} 0.3 & -0.3 \\ 0.2 & 0.1 \\ 0.3 & -0.2 \end{bmatrix} \quad \nabla_{W^{(2)}} L = \begin{bmatrix} -0.17 & -0.44 \\ -0.09 & -0.23 \\ -0.08 & -0.22 \end{bmatrix}$$

$$W'^{(2)} = W^{(2)} - \lambda \nabla_{W^{(2)}} L = \begin{bmatrix} 0.39 & -0.08 \\ 0.25 & 0.22 \\ 0.34 & -0.09 \end{bmatrix}$$
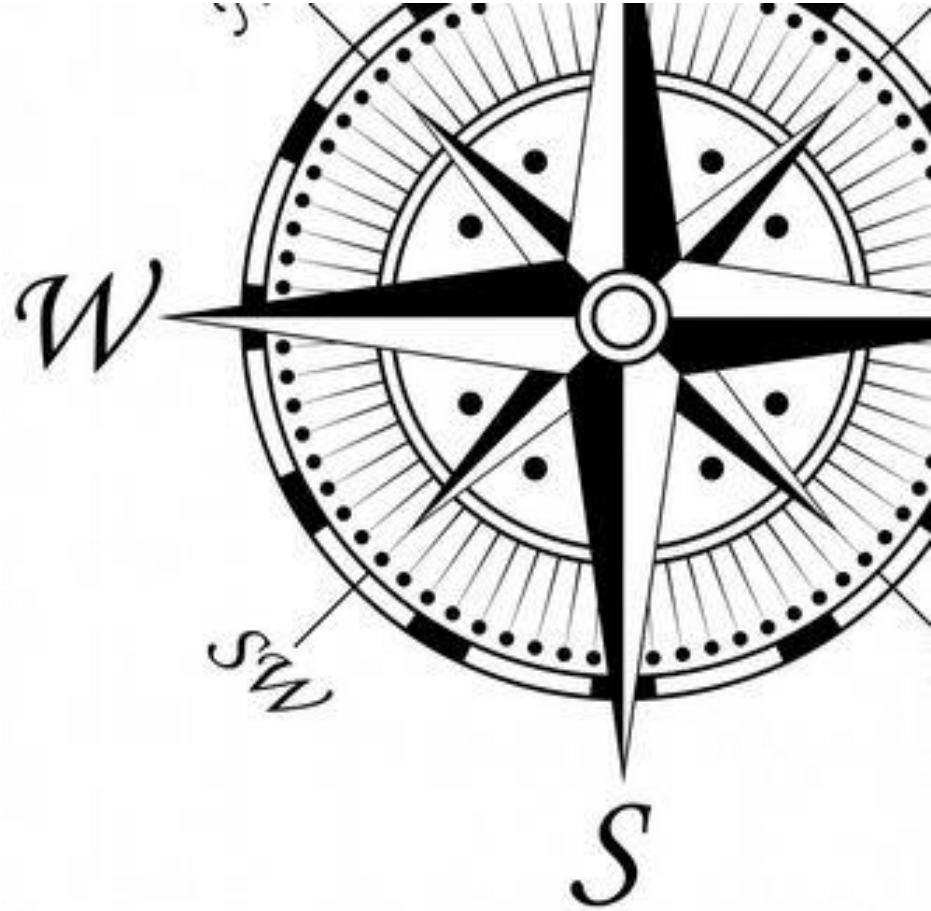


$$g \leftarrow \nabla_{\hat{y}} J = \nabla_{\hat{y}} L(\hat{y}, y)$$
$$\textbf{for } k = l, l-1, \dots, 1 \textbf{ do}$$
$$\quad g \leftarrow \nabla_{a^{(k)}} J = g \odot f'(a^{(k)})$$
$$\quad \nabla_{b^{(k)}} J = g + \lambda \nabla_{b^{(k)}} \Omega(\theta)$$
$$\quad \nabla_{W^{(k)}} J = g h^{(k-1)T} + \lambda \nabla_{W^{(k)}} \Omega(\theta)$$
$$\quad g \leftarrow \nabla_{h^{(k-1)}} J = W^{(k)T} g$$
$$\textbf{end for}$$

# Topics Today

# Summary

- Will be filled out during next session!

# Learning Goals for this Chapter

- Know how the perceptron works
- Explain the need for multiple layers
- Understand gradient-based optimization
- Describe the components of deep neural networks
- Understand and apply the backpropagation algorithm

- Relevant chapters
  - P2, P3
  - S3 (2021) https://www.youtube.com/watch?v=X0Jw4kgaFlg

# References

- [R58] Rosenblatt, F. (1958). The perceptron: a probabilistic model for information storage and organization in the brain. *Psychological review*, *65*(6), 386.

- [MP69] Minsky, M., & Papert, S. (1969). An introduction to computational geometry. *Cambridge tiass., HIT*, *479*(480), 104.