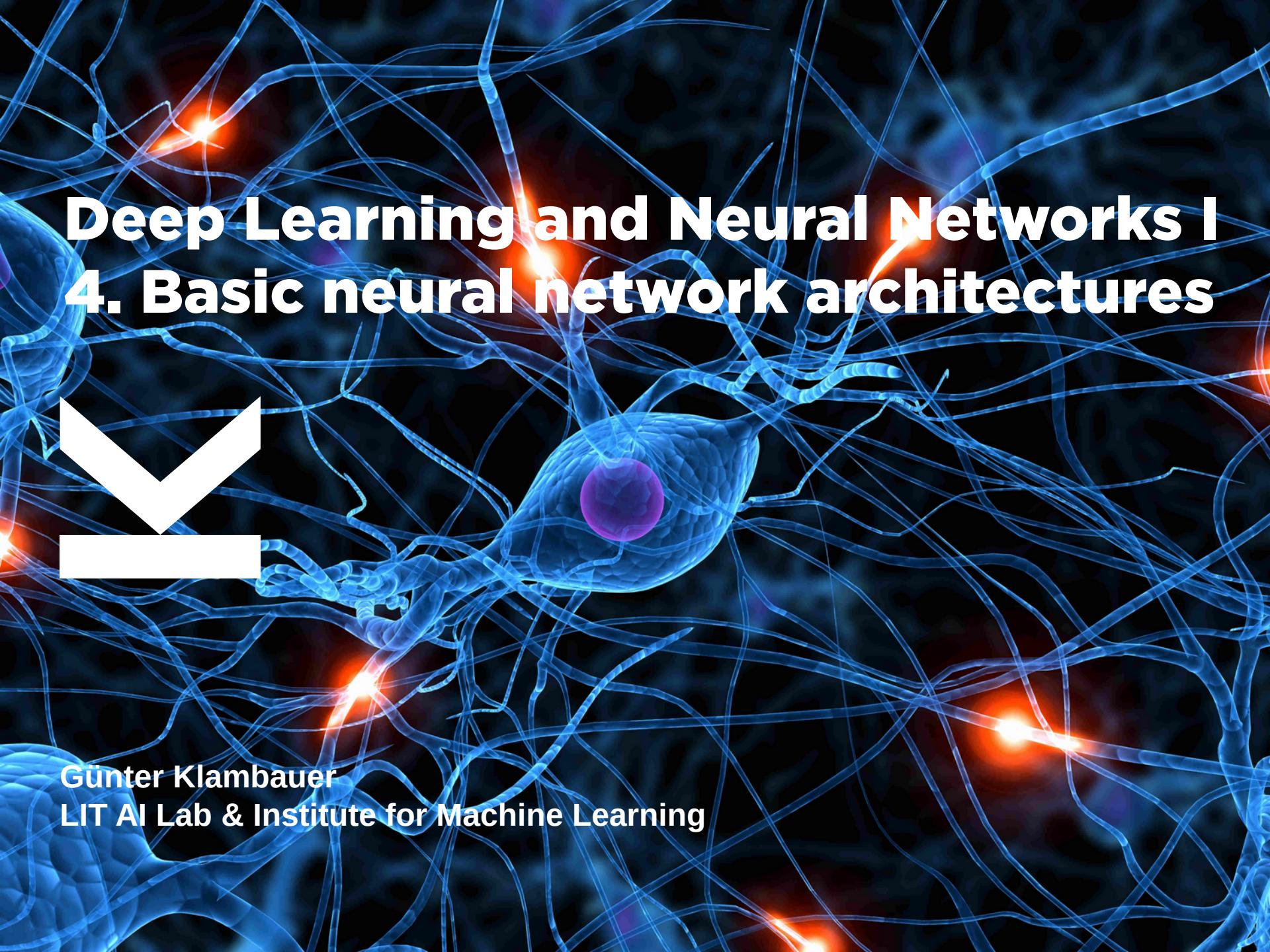


**JOHANNES KEPLER
UNIVERSITY LINZ**



Deep Learning and Neural Networks I

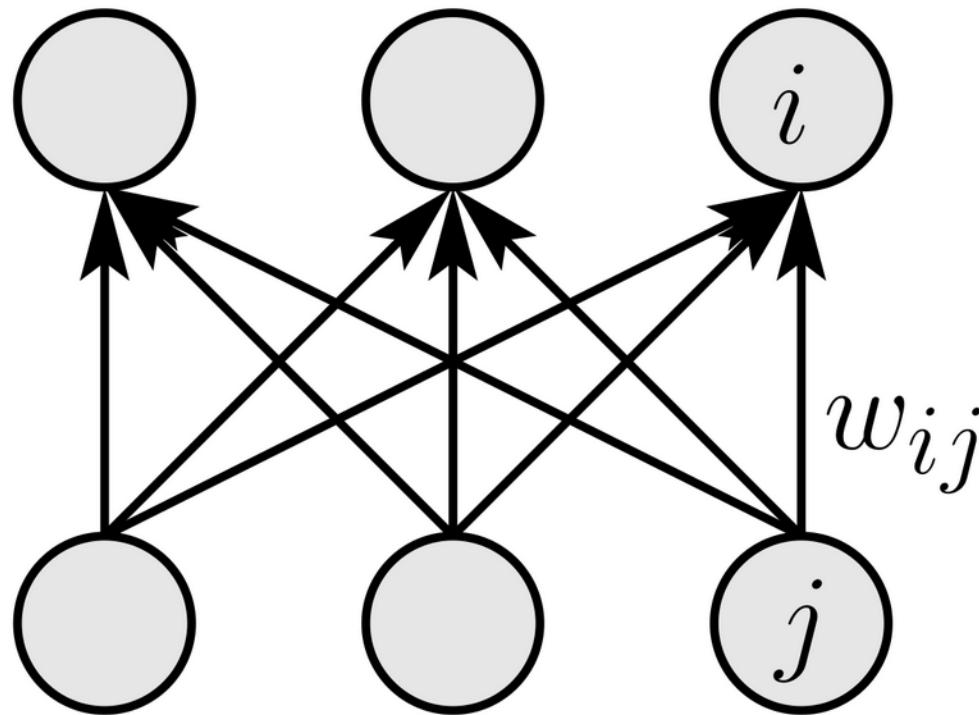
4. Basic neural network architectures



Günter Klambauer
LIT AI Lab & Institute for Machine Learning

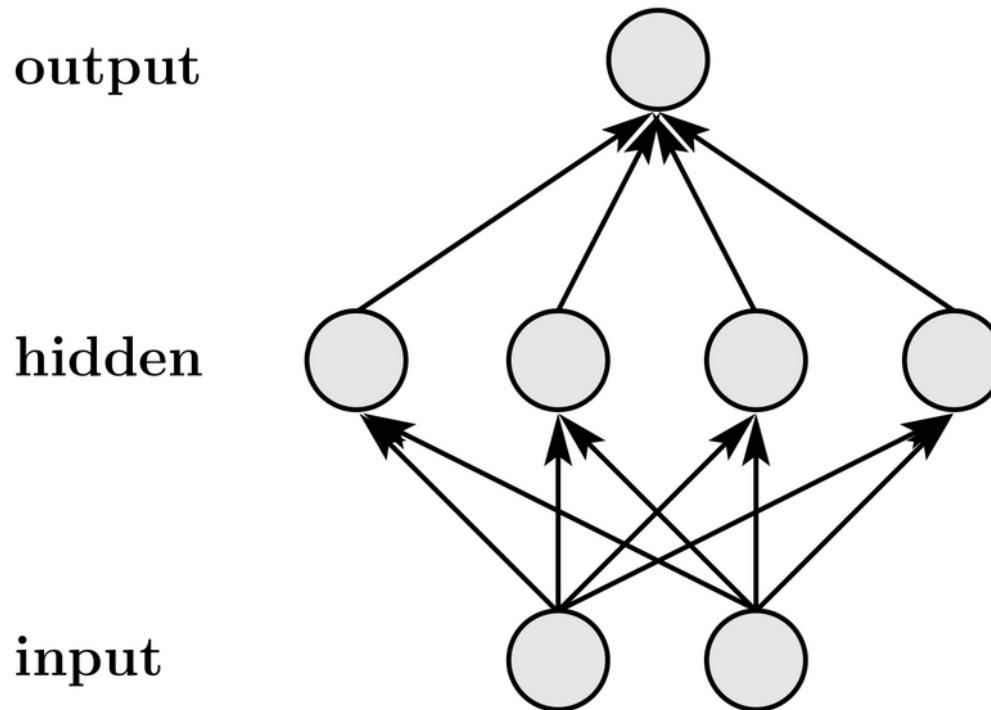
This material, no matter whether in printed or electronic form, may be used for personal and non-commercial educational use only. Any reproduction of this material, no matter whether as a whole or in parts, no matter whether in printed or in electronic form, requires explicit prior acceptance of the authors.

Neural network with adaptive weights



Artificial neural networks: units and weights. The weight w_{ij} gives the weight, connection strength, or synaptic weight from unit j to unit i

Artificial neural networks: 3-layer network



- Artificial neural networks: a 3-layered net with an input, hidden, and output layer

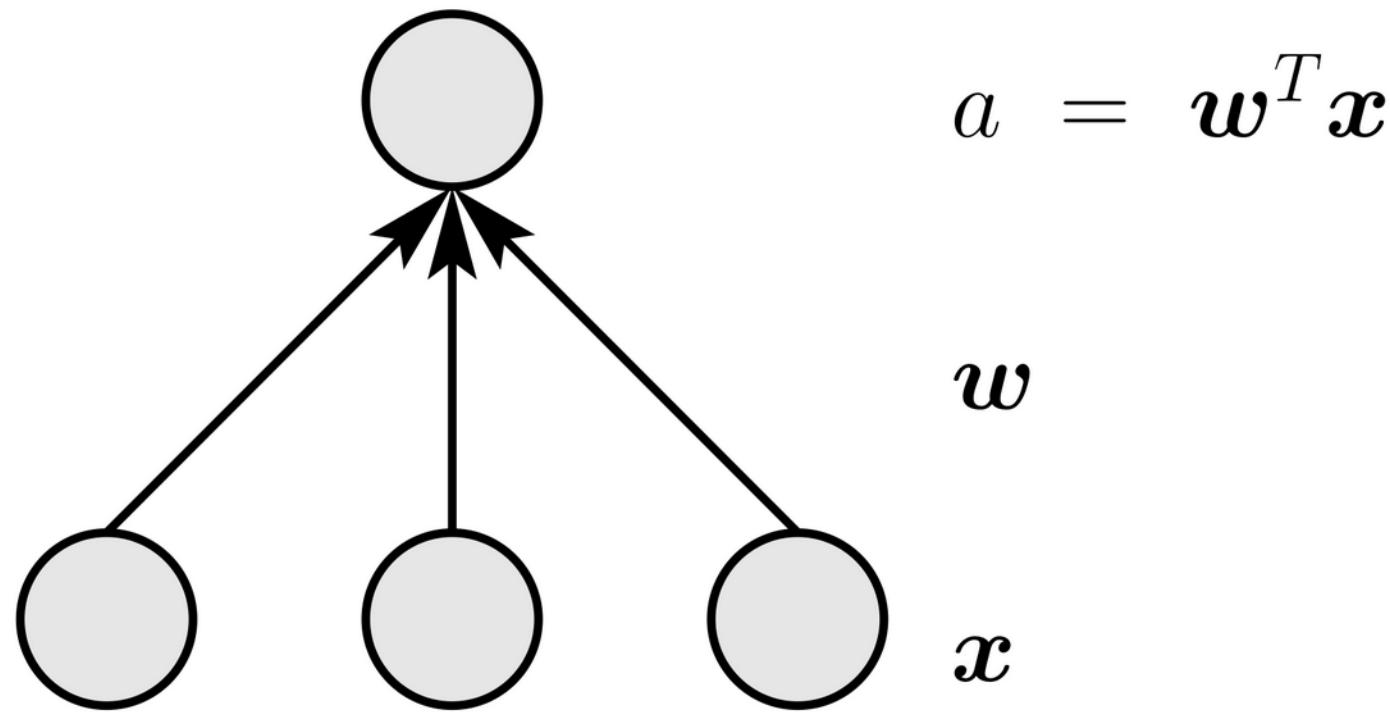
Data for supervised learning: regression

X			y
x^1	0.13	0.03	0.35
x^2	0.15	0.14	0.57
	0.79	0.19	0.87
	0.48	0.28	1.21
	0.93	0.43	0.48
x^n	0.43	0.41	0.70
	0.62	0.63	-0.44
	0.69	0.69	-1.05
	0.19	0.79	-1.29
	0.23	0.85	-1.11
x^N	0.11	0.99	0.32

- Our supervised data set is:
 - Samples are rows in the data matrix \mathbf{X} :
$$\mathbf{X} = (\mathbf{x}^1, \dots, \mathbf{x}^N)$$
 - Features are columns of \mathbf{X} :
 - Single features, e.g.: $x_{21} = 0.15$
 - Feature vector, e.g.: $\mathbf{x}_{\cdot 1}$
 - Scalar label for each data point:

$$\mathbf{y} = (y^1, \dots, y^N)^T$$

Artificial neural networks: linear neuron



- Inputs directly connected to output
- Input: x adaptive weights: w output (activation): a

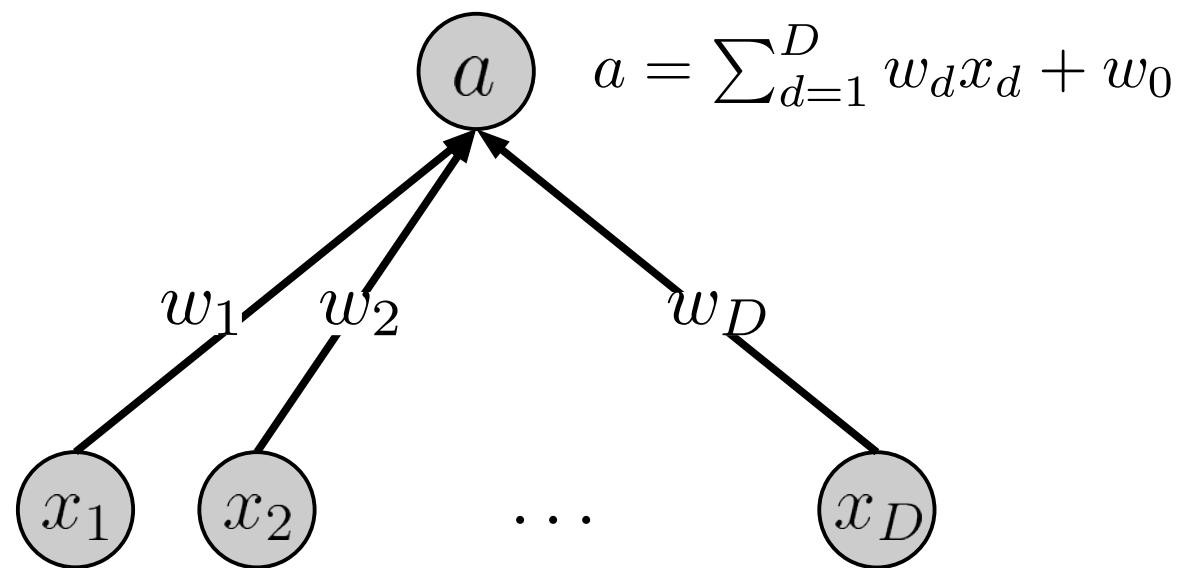
Overview

- 4. Basic neural network architectures
 - 4.1 Linear model: Neuron
 - 4.2 Perceptron
 - 4.3 Linear neuron with sigmoid activation: logistic regression
 - 4.4 Linear neuron with softmax activation: softmax regression

4.1 Linear neural network: model

4.1 Linear neural network: model

$$g(\mathbf{x}; \mathbf{w}) = a = \mathbf{w}^T \mathbf{x}$$



- Corresponds to a linear model as used for linear regression

Linear neural network: additive noise

- We assume that the target variable y is given by a deterministic function $g(\mathbf{x}, \mathbf{w})$ with additive Gaussian noise such that

$$y = g(\mathbf{x}, \mathbf{w}) + \epsilon,$$

where ϵ follows $\epsilon \sim \mathcal{N}(0, \sigma)$.

- Thus we have the distribution of the target variable

$$p(y \mid \mathbf{x}, \mathbf{w}, \sigma) = \mathcal{N}(y \mid g(\mathbf{x}, \mathbf{w}), \sigma)$$

- As a prediction we could give the conditional mean

$$\mathbb{E}(y \mid \mathbf{x}) = \int y p(y \mid \mathbf{x}) dy = g(\mathbf{x}, \mathbf{w})$$

Linear neural network: likelihood

- Now consider a dataset of inputs $\mathbf{X} = \{\mathbf{x}^1, \dots, \mathbf{x}^N\}$ with corresponding target values (labels) y^1, \dots, y^N
Assuming that the data points are drawn independently from each other, we obtain the following likelihood function:

$$p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \sigma) = \prod_{n=1}^N \mathcal{N}(y^n \mid g(\mathbf{x}^n, \mathbf{w}), \sigma)$$

$$p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \sigma) = \prod_{n=1}^N \mathcal{N}(y^n \mid \mathbf{w}^T \mathbf{x}^n, \sigma)$$

$$p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \sigma) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y^n - \mathbf{w}^T \mathbf{x}^n)^2}{2\sigma^2}}$$

Linear NN: residual, loss function

- We take the logarithm of the likelihood function which does not change the location of the minimum:

$$p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \sigma) = \prod_{n=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y^n - \mathbf{w}^T \mathbf{x}^n)^2}{2\sigma^2}}$$

$$\log p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \sigma) = -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{\sigma^2} S(\mathbf{w}, \mathbf{X}, \mathbf{y})$$

$$S(\mathbf{w}, \mathbf{X}, \mathbf{y}) = \frac{1}{2} \sum_{n=1}^N \underbrace{(y^n - \mathbf{w}^T \mathbf{x}^n)^2}_{:= \epsilon^n}$$

- The function: $L(y, g(\mathbf{x}; \mathbf{w})) = \frac{1}{2}(y - g(\mathbf{x}; \mathbf{w}))^2$
is called *squared loss*

Linear NN: candidate parameters

- Assume we have a candidate parameter vector $\tilde{\mathbf{w}}$:

$$S(\tilde{\mathbf{w}}) = \frac{1}{2} \sum_{n=1}^N (\epsilon^n)^2 = \frac{1}{2} \sum_{n=1}^N (y^n - \tilde{\mathbf{w}}^T \mathbf{x}^n)^2 = \frac{1}{2} (\mathbf{y} - \mathbf{X}\tilde{\mathbf{w}})^T (\mathbf{y} - \mathbf{X}\tilde{\mathbf{w}}).$$

(we have dropped the dependency on data and labels)

- We can determine the *sum of residuals*
 - Low sum of residuals is equivalent to high likelihood
- **Learning:**
 - We can check different candidate parameter vectors and determine $S(\tilde{\mathbf{w}})$
 - Maximizing the likelihood of our data (under the assumption of Gaussian noise), is equivalent to minimizing the squared error

Linear NN: learning

- Let us *learn* the prediction task, i.e., find the parameters that maximize the likelihood or minimize the quadratic loss on the given data (*training set*):

$$\hat{\boldsymbol{w}} = \arg \min_{\tilde{\boldsymbol{w}}} S(\tilde{\boldsymbol{w}})$$

- This approach is called *empirical risk minimization* (ERM)

$$R_{\text{emp}}(\tilde{\boldsymbol{w}}, \mathbf{X}, \mathbf{y}) = S(\tilde{\boldsymbol{w}}).$$

- Does this problem have a solution? How many?
- How could we find a solution?

Linear NN: learning

- Let us *learn* the prediction task, i.e., find the parameters that maximize the likelihood or minimize the quadratic loss on the given data (*training set*):

$$\log p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \sigma) = -\frac{N}{2} \log(2\pi\sigma^2) - \frac{1}{\sigma^2} S(\mathbf{w})$$

$$S(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (y^n - \mathbf{w}^T \mathbf{x}^n)^2$$

$$\nabla_{\mathbf{w}} \log p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}, \sigma) = \frac{1}{\sigma^2} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}^n - y^n) \mathbf{x}^n$$

Linear regression: closed form solution

- Let us *learn* the prediction task, i.e., find the parameters that maximize the likelihood or minimize the quadratic loss:

$$0 = \frac{1}{\sigma^2} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}^n - y^n) \mathbf{x}^n$$

$$0 = \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y}$$

$$\mathbf{X}^T \mathbf{y} = \mathbf{X}^T \mathbf{X} \mathbf{w}$$

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

Linear NN: learning

- This is the solution to our optimization problem

$$\hat{\mathbf{w}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- It is a maximum-likelihood estimator
- It minimizes the least-squares problem
- The quantity

$$\mathbf{X}^+ = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

is called Moore-Penrose pseudo-inverse.

- Note that we have used the assumption that $\mathbf{X}^T \mathbf{X}$ has full rank.
- $\frac{1}{n} \mathbf{X}^T \mathbf{X}$ is an estimate of the covariance matrix (centered data)

Linear NN: gradient descent

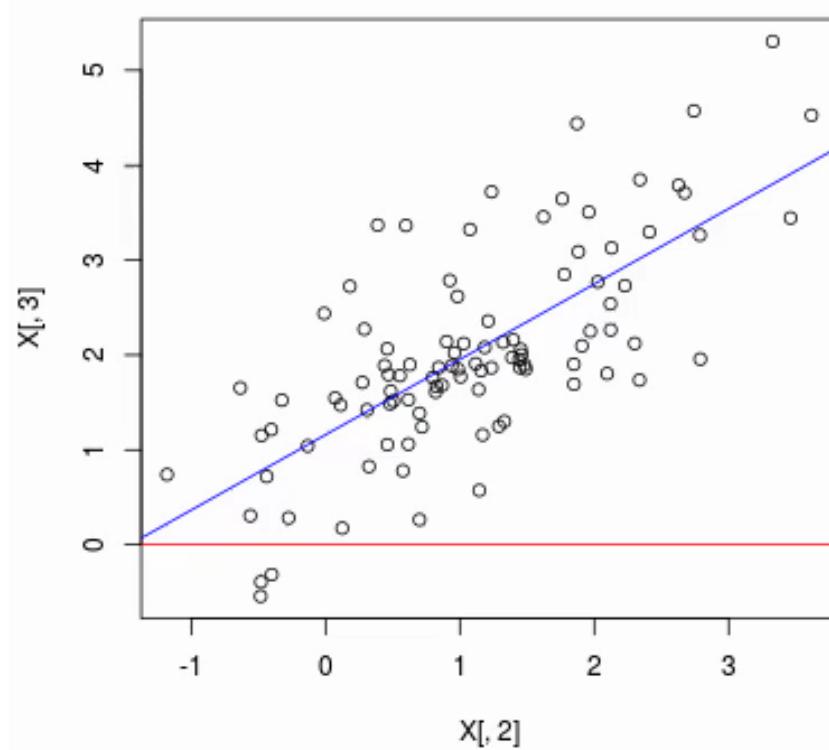
- We seek a $\hat{\mathbf{w}}$ that minimizes $S(\mathbf{w})$. To do so, we could also use a search algorithm that uses an initial guess for \mathbf{w} and gradually improves it:

$$\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} - \eta \nabla_{\mathbf{w}} S(\mathbf{w}) \mid_{\mathbf{w}^{\text{old}}}$$

- This algorithm is called *gradient descent*
 - The choice of the learning rate η is critical for convergence
 - Works well for the least squared problem because it is convex
- We already know the expression for the gradient:

$$\nabla_{\mathbf{w}} S(\mathbf{w}) = \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}^n - y^n) \mathbf{x}^n$$

Linear NN: gradient descent



Linear NN: stochastic grad. desc.

- We note that the *scaled gradient*

$$\nabla_{\mathbf{w}} S(\mathbf{w}) = \frac{1}{N} \sum_{n=1}^N (\mathbf{w}^T \mathbf{x}^n - y^n) \mathbf{x}^n$$

is a mean over data points.

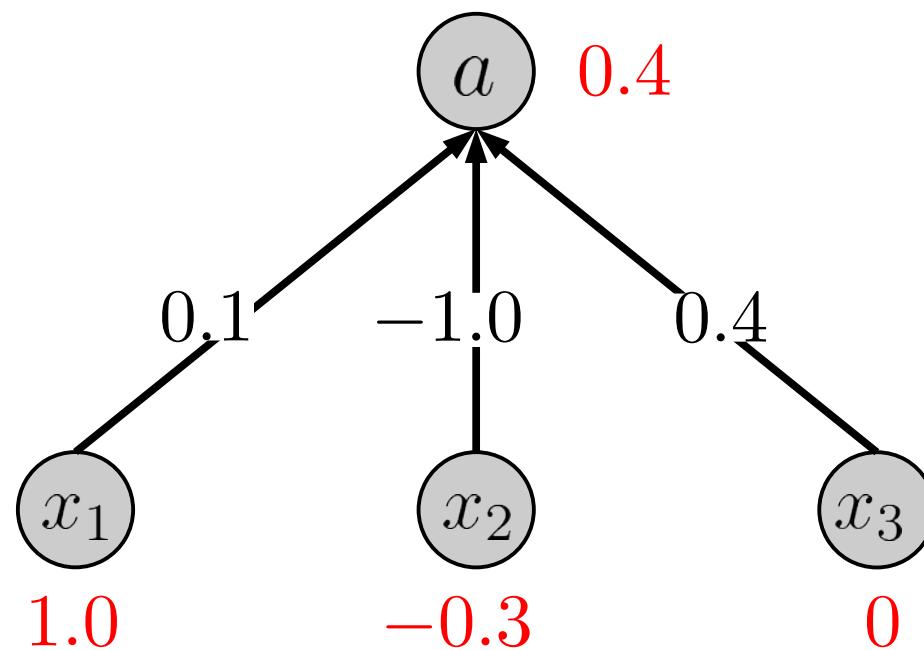
- Thus we can approximate the *scaled gradient* on a subset of datapoints, called *mini-batch*:

$$\nabla_{\mathbf{w}} S(\mathbf{w}) \approx \frac{1}{|\mathcal{B}|} \sum_{b=1}^{|\mathcal{B}|} (\mathbf{w}^T \mathbf{x}^{n_b} - y^{n_b}) \mathbf{x}^{n_b}$$

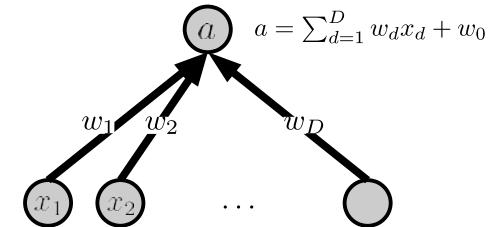
$$\mathcal{B} = \{\mathbf{x}_{n_1}, \dots, \mathbf{x}_{n_B}\} \subset \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$$

Linear neural network: Trained network

$$g(\mathbf{x}^{\text{new}}; \mathbf{w}) = a = \mathbf{w}^T \mathbf{x}^{\text{new}}$$



Summary



- Architecture of linear neural networks
- Objective function: sum of residuals, empirical error; motivation by maximum-likelihood
- Learning: minimization of objective function
 - Different ways to minimize
 - Properties of the solution
- Goal: Trained network

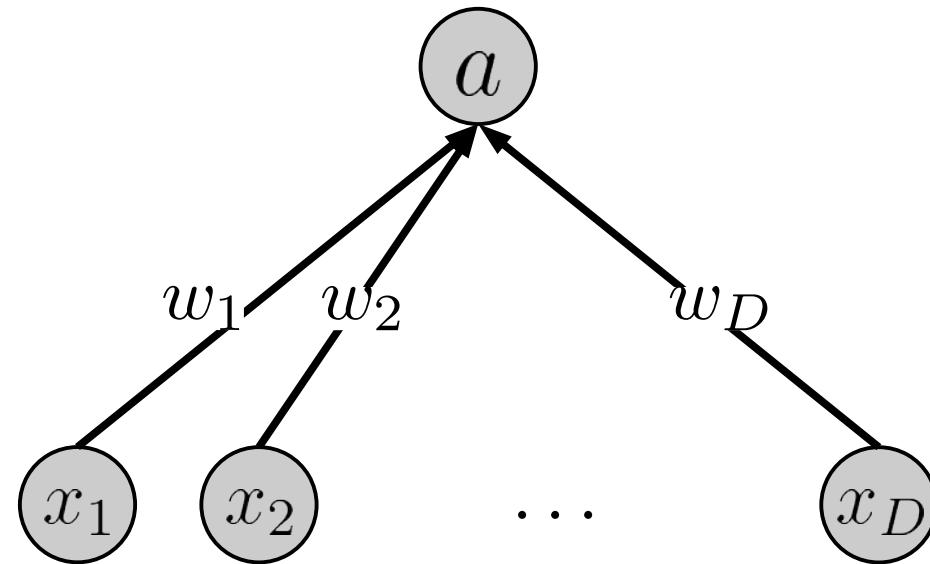
Perceptron: classification task

- Let us assume that the targets/labels only take the values -1 and 1 (classes)
- This is a supervised classification tasks
- Examples for classification tasks
 - Diseased vs healthy
 - Spam vs normal email
 - Win or loss
 - ...
- What can we do in this situation?
Should we just use LR as before?

X			y
x^1	0.13	0.03	1
x^2	0.15	0.14	1
	0.79	0.19	-1
	0.48	0.28	1
	0.93	0.43	-1
	0.43	0.41	-1
	0.62	0.63	1
	0.69	0.69	1
	0.19	0.79	1
	0.23	0.85	-1
x^N	0.11	0.99	-1

4.2 Perceptron

$$g(\mathbf{x}; \mathbf{w}) = a = \text{sign}(\mathbf{w}^T \mathbf{x})$$



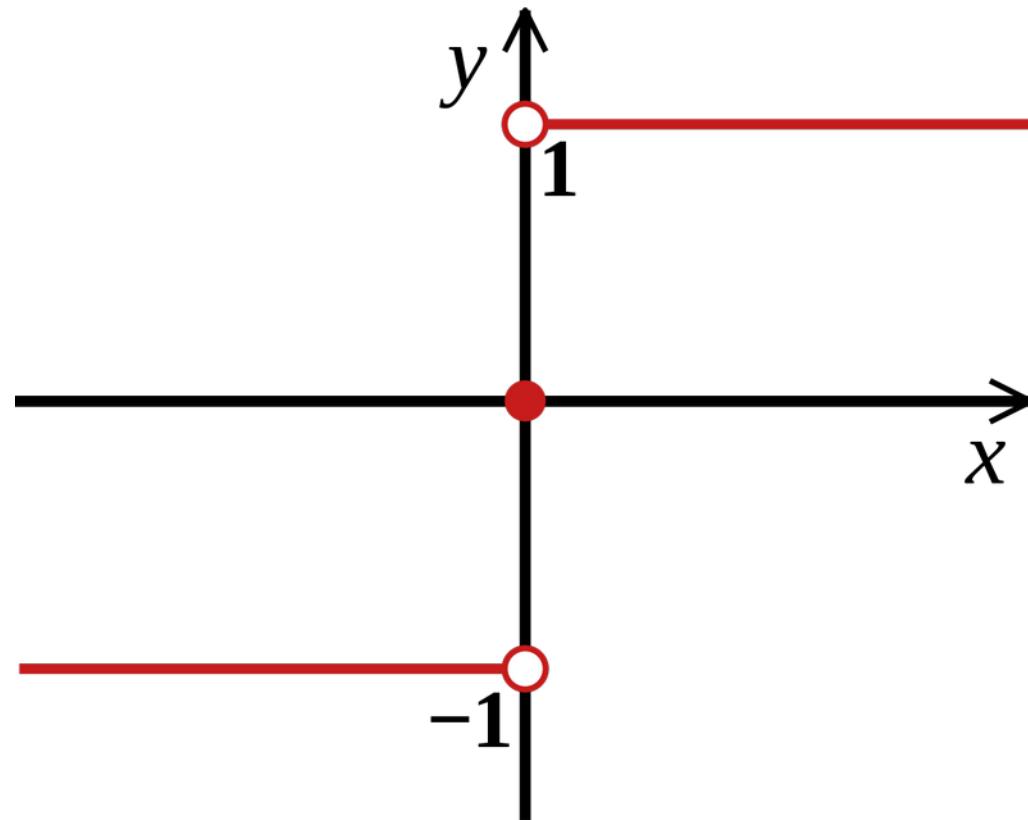
- Similar to linear neuron with linear activation

Reminder: Perceptron

- **1958:** Frank Rosenblatt
 - Perceptron
 - Electronic device with adaptable connections
 - 400 photocells as input neurons
 - Image recognition
 - Ability to learn

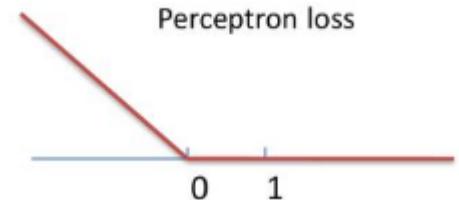


Perceptron: sign function



Perceptron: loss function

- Loss function: “perceptron loss”



$$L(y, \text{sign}(\mathbf{w}^T \mathbf{x})) = \begin{cases} 0, & \text{if } y \text{ sign}(\mathbf{w}^T \mathbf{x}) \geq 0 \\ y \mathbf{w}^T \mathbf{x} & \text{else} \end{cases}$$

- Empirical error

$$R_{\text{emp}}(\mathbf{y}, \mathbf{X}, \mathbf{w}) = - \sum_{n=1; a^n y^n = -1}^N y^n \mathbf{w}^T \mathbf{x}^n$$

- Note that the sum runs over all samples that are misclassified

Perceptron: learning rule

- Derivative with respect to weights:

$$\nabla_{\mathbf{w}} R_{\text{emp}}(\mathbf{y}, \mathbf{X}, \mathbf{w}) = \nabla_{\mathbf{w}} - \sum_{n=1; a^n y^n = -1}^N y^n \mathbf{w}^T \mathbf{x}^n = - \sum_{n=1; a^n y^n = -1}^N y^n \mathbf{x}^n$$

- Weight update rule:

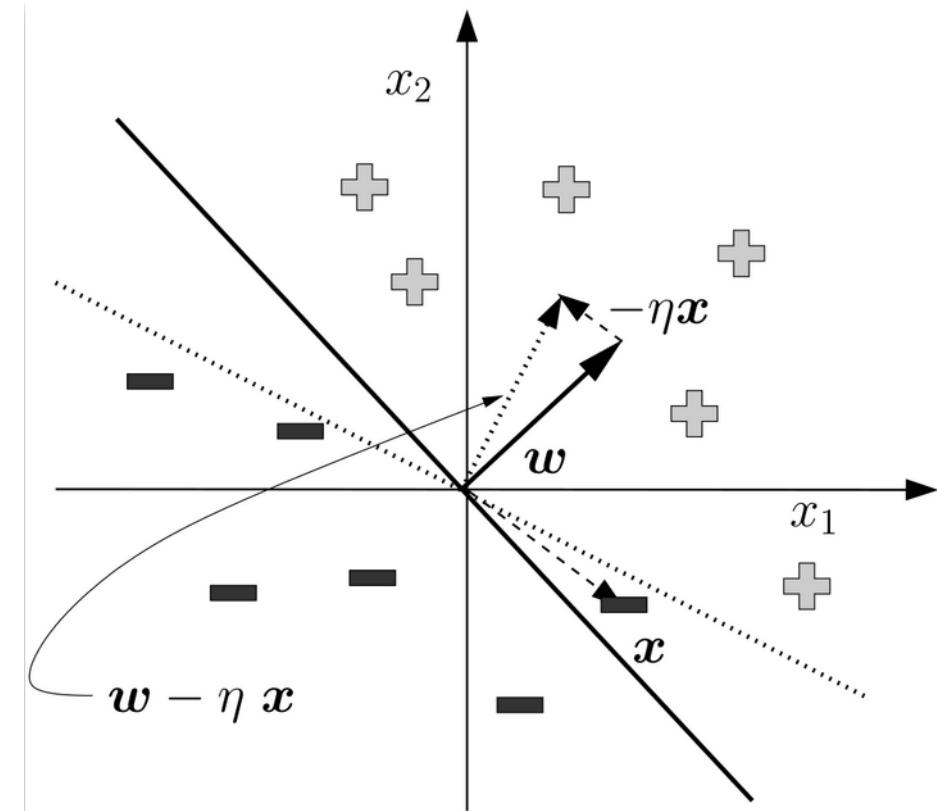
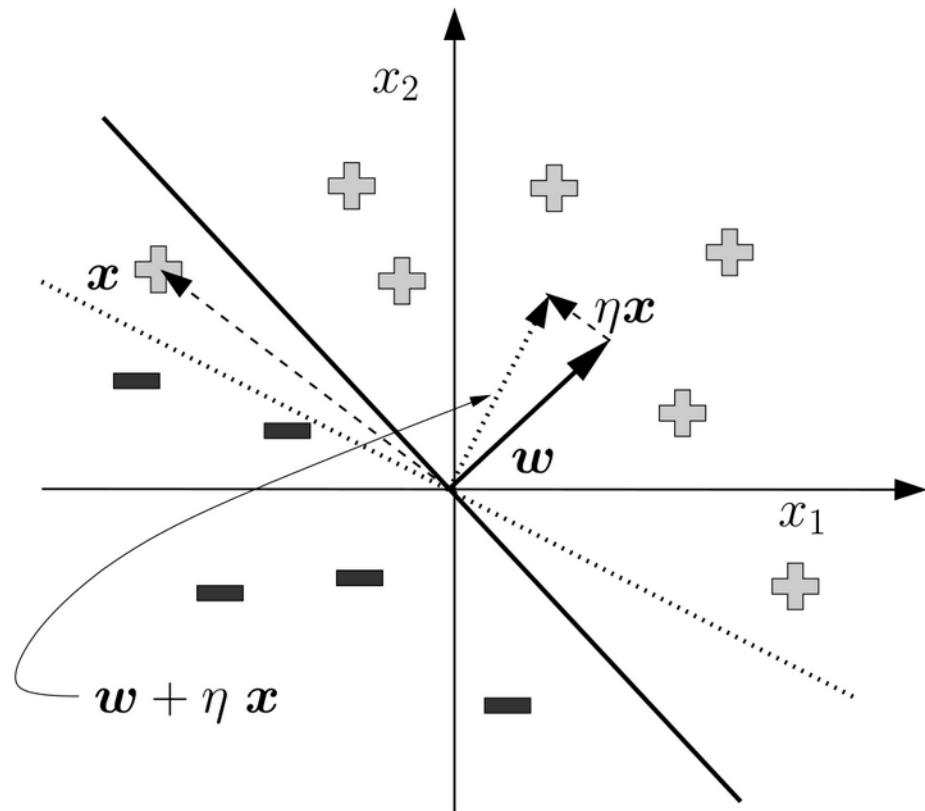
$$\Delta \mathbf{w} = \eta \sum_{n=1; a^n y^n = -1}^N y^n \mathbf{x}^n \quad \text{if } ay = -1 : \Delta \mathbf{w} = \eta y \mathbf{x}.$$

whole data set

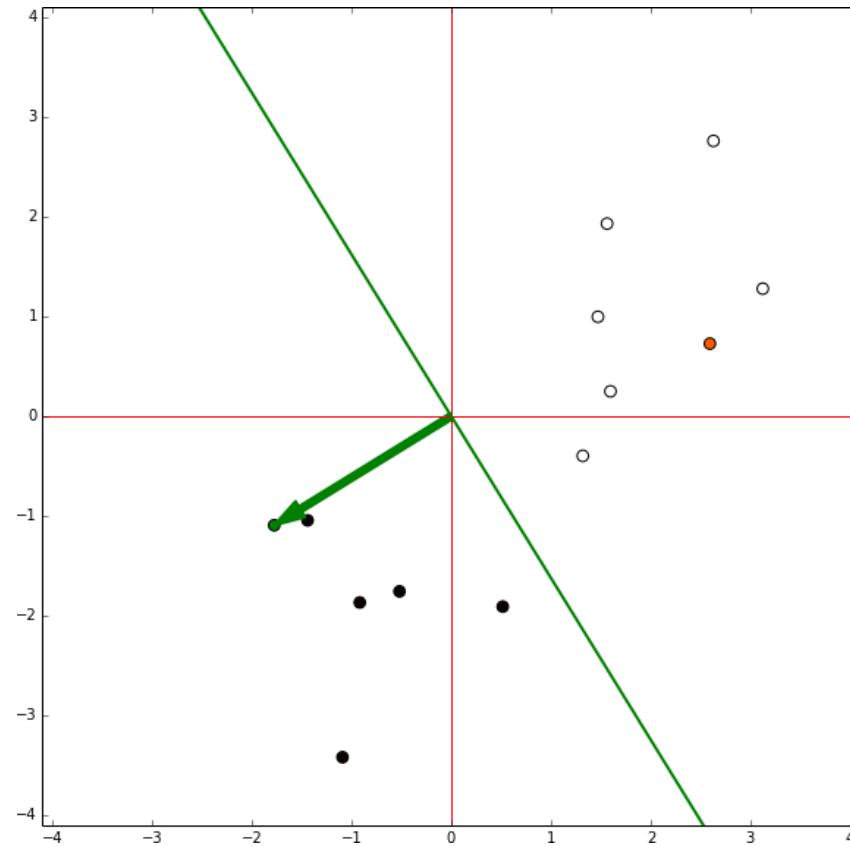
single sample

- Unique solution? Convergence?

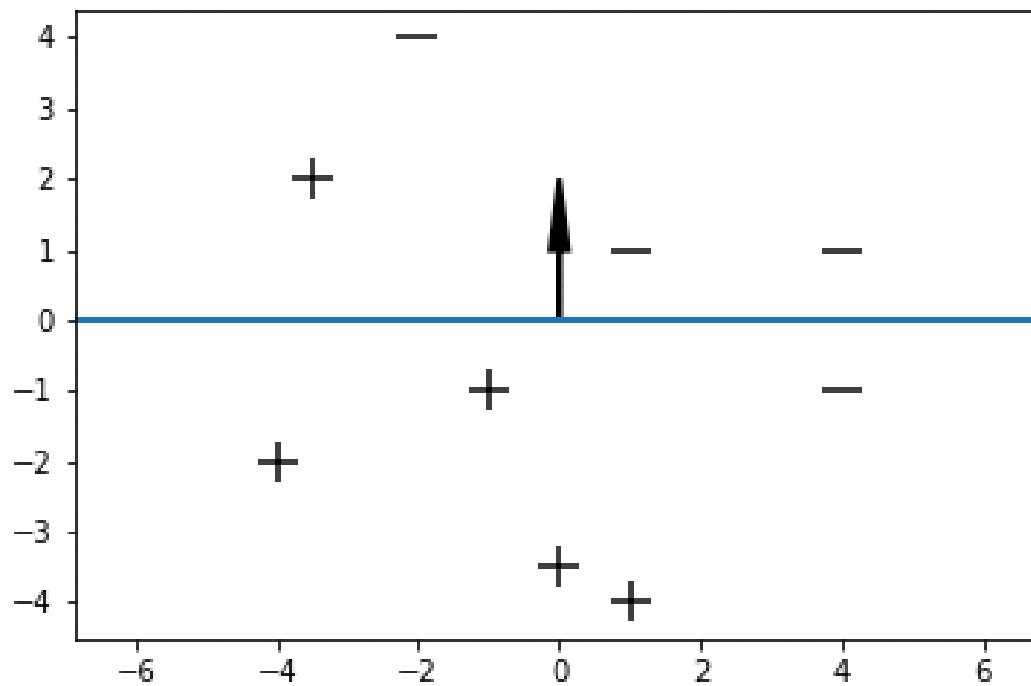
Perceptron: Interpretation of learning rule



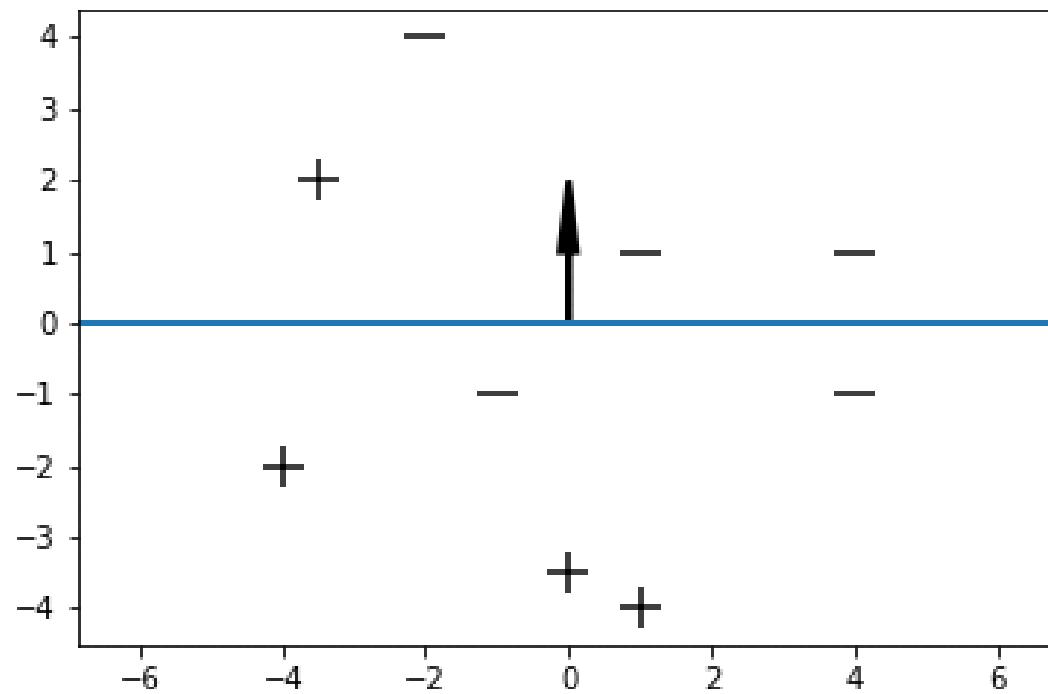
Perceptron: Interpretation of learning rule



Perceptron: separable data



Perceptron: non-separable data



Perceptron Convergence Theorem for linearly separable data (1/3)

- Since the data is linearly separable, there exists $\hat{\mathbf{w}}$:

$$\hat{\mathbf{w}}^T \mathbf{x}^n y^n > 0 \quad \forall n$$

- We assume that we start weights with the zero vector and update according to

$$\mathbf{w}^{\tau+1} = \mathbf{w}^\tau + \mathbf{x}^n y^n$$

if the data point is mis-classified.

- Then the weight vector is: $\mathbf{w} = \sum_n \tau^n \mathbf{x}^n y^n$

$\boxed{\tau^n}$

- Then we multiply:

$$\hat{\mathbf{w}}^T \mathbf{w} = \sum_n \tau^n \hat{\mathbf{w}}^T \mathbf{x}^n y^n \geq \tau \min_n (\hat{\mathbf{w}}^T \mathbf{x}^n y^n)$$

Perceptron Convergence Theorem for linearly separable data (2/3)

- We remember: $\mathbf{w}^{\tau+1} = \mathbf{w}^\tau + \mathbf{x}^n y^n$ and we have

$$\begin{aligned}\|\mathbf{w}^{(\tau+1)}\|^2 &= \|\mathbf{w}^{(\tau)}\|^2 + \|\mathbf{x}^n\|^2 (y^n)^2 + 2\mathbf{w}^{(\tau)}^T \mathbf{x}^n y^n \\ &\leq \|\mathbf{w}^{(\tau)}\|^2 + \|\mathbf{x}^n\|^2 (y^n)^2\end{aligned}$$

because the last term is negative.

- The change of weights satisfies:

$$\Delta \|\mathbf{w}\|^2 = \|\mathbf{w}^{(\tau+1)}\|^2 - \|\mathbf{w}^{(\tau)}\|^2 \leq \|\mathbf{x}\|_{\max}^2$$

- After τ updates we have: $\|\mathbf{w}\|^2 \leq \tau \|\mathbf{x}\|_{\max}^2$

Perceptron Convergence Theorem for linearly separable data (3/3)

- Thus the function $\hat{w}^T w$ is bounded below by a function that grows linearly with tau

$$\hat{w}^T w \geq \tau \min_n (\hat{w}^T x^n y^n)$$

- The length of the weight vector increases no faster than with $\tau^{1/2}$

$$\|w\|^2 \leq \tau \|x\|_{\max}^2$$

- Therefore, since \hat{w} is fixed, for sufficiently large time, the results would incompatible. Hence the algorithm must converge in a finite number of time steps. *qed*

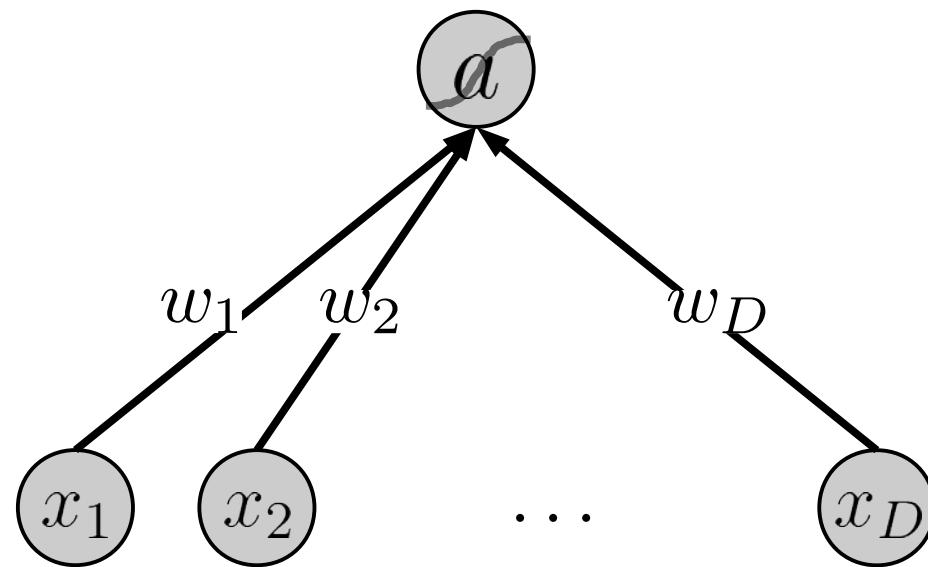
Logistic neuron: classification

- Let us assume that the targets/labels only take the values 0 and 1 (classes)
- This is a supervised classification tasks
- Examples for classification tasks
 - Diseased vs healthy
 - Spam vs normal email
 - Win or loss
 - ...
- What can we do in this situation?
Should we just use LR as before?

X			y
x^1	0.13	0.03	1
x^2	0.15	0.14	1
	0.79	0.19	0
	0.48	0.28	1
	0.93	0.43	0
	0.43	0.41	0
	0.62	0.63	1
	0.69	0.69	1
	0.19	0.79	1
	0.23	0.85	0
x^n	0.11	0.99	0

4.3 Logistic neuron: model

$$g(\mathbf{x}; \mathbf{w}) = a = \sigma(\mathbf{w}^T \mathbf{x})$$



- Corresponds to a linear model as used for logistic regression

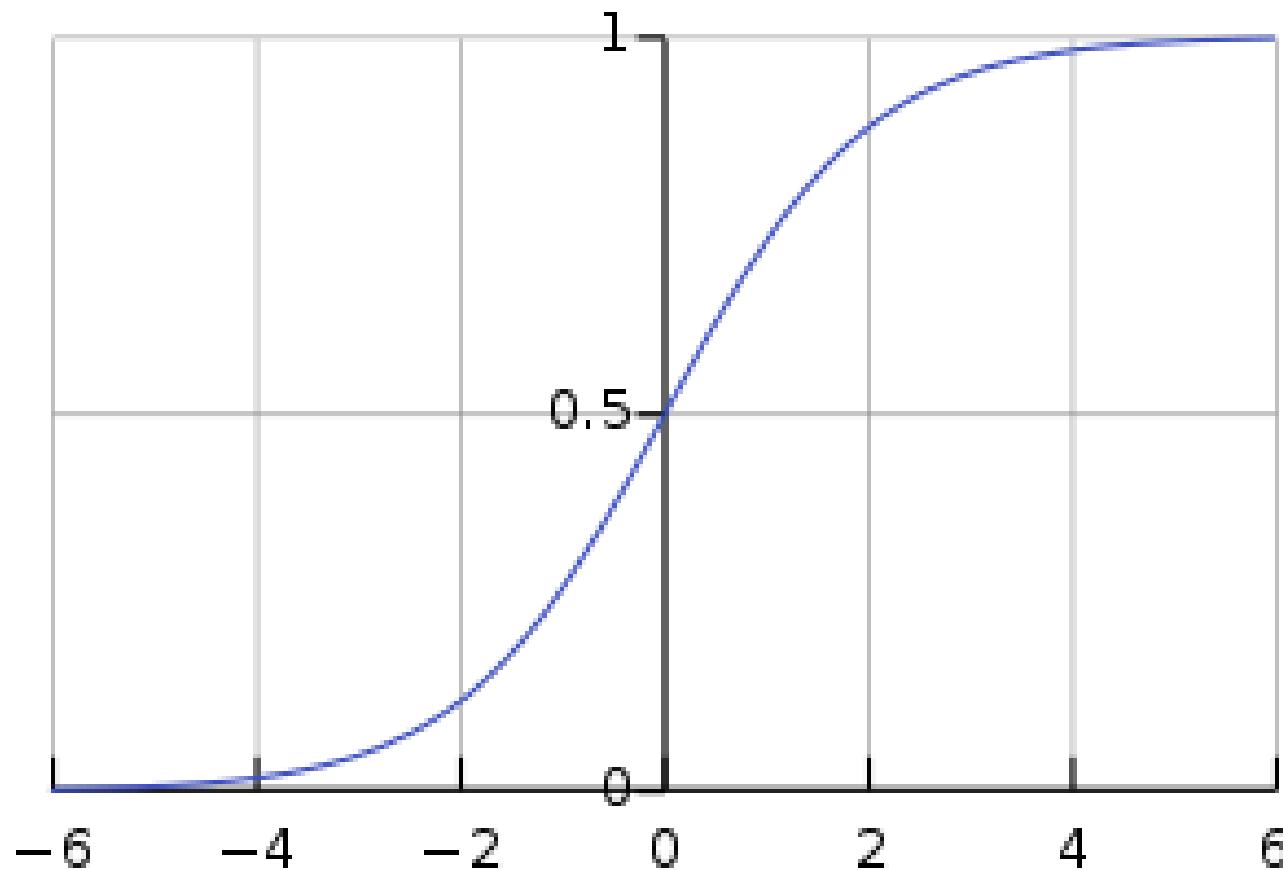
Linear neuron with sigmoid activation:

- The main idea is to fit a model of the following form

$$g(\mathbf{x}; \mathbf{w}) = \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} = \sigma(\mathbf{w}^T \mathbf{x}),$$

where σ is the logistic sigmoid function.

Logistic regression: The sigmoid function



Logistic regression

- Property of the sigmoid function:

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

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

Logistic regression: likelihood

- We aim at maximizing the likelihood again

$$p(y = 1 \mid \mathbf{x}; \mathbf{w}) = g(\mathbf{x}; \mathbf{w})$$

$$p(y = 0 \mid \mathbf{x}; \mathbf{w}) = 1 - g(\mathbf{x}; \mathbf{w})$$

- This can be written more compactly as:

$$p(y \mid \mathbf{x}; \mathbf{w}) = g(\mathbf{x}; \mathbf{w})^y (1 - g(\mathbf{x}; \mathbf{w}))^{(1-y)}$$

- Assuming that the data points are drawn independently,

the likelihood is

$$p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}) = \prod_{n=1}^N g(\mathbf{x}^n; \mathbf{w})^{y^n} (1 - g(\mathbf{x}^n; \mathbf{w}))^{(1-y^n)}$$

Logistic regr: loss and error function

- The likelihood is

$$p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}) = \prod_{n=1}^N g(\mathbf{x}^n; \mathbf{w})^{y^n} (1 - g(\mathbf{x}^n; \mathbf{w}))^{(1-y^n)}$$

- The negative log likelihood is

$$-\log p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}) = -\sum_{n=1}^N y^n \log g(\mathbf{x}^n; \mathbf{w}) + (1 - y^n) \log(1 - g(\mathbf{x}^n; \mathbf{w}))$$

- The function

$$L(y, g(\mathbf{x}; \mathbf{w})) = -(y \log g(\mathbf{x}; \mathbf{w}) + (1 - y) \log(1 - g(\mathbf{x}; \mathbf{w})))$$

is called *binary cross-entropy*.

Logistic regression: learning

- One way to learn the model is again *gradient descent*

$$\begin{aligned}\nabla_{\mathbf{w}} - \log p(\mathbf{y} \mid \mathbf{X}, \mathbf{w}) &= \nabla_{\mathbf{w}} R_{\text{emp}}(\mathbf{y}, \mathbf{X}, \mathbf{w}) = \\ &= \nabla_{\mathbf{w}} - \sum_{n=1}^N y^n \log \sigma(\mathbf{w}^T \mathbf{x}^n) + (1 - y^n) \log(1 - \sigma(\mathbf{w}^T \mathbf{x}^n)) = \\ &= - \sum_{n=1}^N \left(y^n \frac{1}{\sigma(\mathbf{w}^T \mathbf{x}^n)} - (1 - y^n) \frac{1}{(1 - \sigma(\mathbf{w}^T \mathbf{x}^n))} \right) \nabla_{\mathbf{w}} \sigma(\mathbf{w}^T \mathbf{x}^n) = \\ &= - \sum_{n=1}^N \left(y^n \frac{1}{\sigma(\mathbf{w}^T \mathbf{x}^n)} - (1 - y^n) \frac{1}{(1 - \sigma(\mathbf{w}^T \mathbf{x}^n))} \right) \sigma(\mathbf{w}^T \mathbf{x}^n)(1 - \sigma(\mathbf{w}^T \mathbf{x}^n)) \mathbf{x}^n = \\ &= \sum_{n=1}^N (\sigma(\mathbf{w}^T \mathbf{x}^n) - y^n) \mathbf{x}^n\end{aligned}$$

- Does that look familiar? Is it exactly the same as for linear regression?

Logistic regression: second derivative

- Gradient and Hessian for a single sample

$$\nabla_{\mathbf{w}} R_{\text{emp}}(y, \mathbf{x}, \mathbf{w}) = (\sigma(\mathbf{w}^T \mathbf{x}) - y) \mathbf{x}$$

$$\nabla_{\mathbf{w}}^2 R_{\text{emp}}(y, \mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x})(1 - \sigma(\mathbf{w}^T \mathbf{x})) \mathbf{x} \mathbf{x}^T$$

Note: $a = g(\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x})$

- Gradient and Hessian for multiple samples

$$\nabla_{\mathbf{w}} R_{\text{emp}}(\mathbf{y}, \mathbf{X}, \mathbf{w}) = \sum_{n=1}^N (\sigma(\mathbf{w}^T \mathbf{x}^n) - y^n) \mathbf{x}^n = \mathbf{X}^T (\mathbf{a} - \mathbf{y})$$

$$\nabla_{\mathbf{w}}^2 R_{\text{emp}}(\mathbf{y}, \mathbf{X}, \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x}^n)(1 - \sigma(\mathbf{w}^T \mathbf{x}^n)) \mathbf{x}^n \mathbf{x}^{nT} = \mathbf{X}^T \mathbf{R} \mathbf{X}$$

where \mathbf{R} is a diagonal matrix with $\sigma(\mathbf{w}^T \mathbf{x}^n)(1 - \sigma(\mathbf{w}^T \mathbf{x}^n))$ as entries

Logistic regression: learning

- Again we can view the negative log likelihood as data driven error and use *gradient descent*:

$$\boldsymbol{w}^{\text{new}} = \boldsymbol{w}^{\text{old}} - \eta \nabla_{\boldsymbol{w}} R_{\text{emp}}(\boldsymbol{y}, \boldsymbol{X}, \boldsymbol{w})|_{\boldsymbol{w}^{\text{old}}}$$

- Also *stochastic gradient descent* could be used:

$$\boldsymbol{w}^{\text{new}} = \boldsymbol{w}^{\text{old}} - \eta \tilde{\nabla}_{\boldsymbol{w}} R_{\text{emp}}(\boldsymbol{y}, \boldsymbol{X}, \boldsymbol{w})|_{\boldsymbol{w}^{\text{old}}}$$

- Typically, *Iteratively Reweighted Least-Squares* is used:

$$\boldsymbol{w}^{\text{new}} = \boldsymbol{w}^{\text{old}} - \boldsymbol{H}^{-1} \nabla_{\boldsymbol{w}} R_{\text{emp}}(\boldsymbol{y}, \boldsymbol{X}, \boldsymbol{w})|_{\boldsymbol{w}^{\text{old}}}$$

where \boldsymbol{H} is the Hessian matrix, i.e. the matrix of second derivatives

Intermezzo: Optimization

- Generic optimization problem:

Let $x \in \mathbb{R}^n$, $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$, $h : \mathbb{R}^n \rightarrow \mathbb{R}$. Find

$$\min_x f(x)$$

This would be our empirical error that we aim to minimize

$$\text{s.t. } g(x) \leq 0, h(x) = 0$$

- Notation above is general and not the specific ML notation of this course!

- Optimization types:

- Blackbox optimization: we only have access to $f(x)$
- Gradient-based optimization: we have access to $f(x), \nabla f(x)$
- Second order optimization: we have access to $f(x), \nabla f(x), \nabla^2 f(x)$

CONSTRAINTS

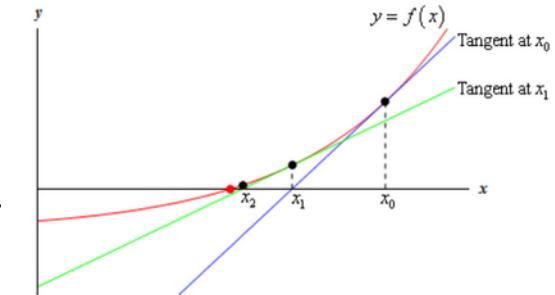
Typically, we do not have constraints

Intermezzo: Optimization

Newton's method

- Finding roots (zeros) of $f(x)$:

$$x^{\text{new}} = x^{\text{old}} - \frac{f(x^{\text{old}})}{f'(x^{\text{old}})}$$



- Finding optima of $f(x)$ in 1D:

$$x^{\text{new}} = x^{\text{old}} - \frac{f'(x^{\text{old}})}{f''(x^{\text{old}})}$$

- Finding optima of $f(x)$ in n-dim:

$$x^{\text{new}} = x^{\text{old}} - \nabla^2 f(x^{\text{old}})^{-1} \nabla f(x^{\text{old}})$$

Intermezzo: Optimization

Newton's method

- Taylor approximation:

$$f(x) \approx f(x_0) + (x - x_0)f'(x_0)$$

- Setting to zero and solving:

$$f(x_0) + (x - x_0)f'(x_0) = 0 \Rightarrow$$

$$x = \frac{-f(x_0) + x_0 f'(x_0)}{f'(x_0)} = x_0 - \frac{f(x_0)}{f'(x_0)}$$

- Analogously for first derivative and multiple dimensions

Second order methods for linear and logistic regression

$$\boldsymbol{w}^{\text{new}} = \boldsymbol{w}^{\text{old}} - (\nabla_{\boldsymbol{w}}^2 R(\boldsymbol{w}^{\text{old}}))^{-1} \nabla_{\boldsymbol{w}} R(\boldsymbol{w}^{\text{old}})$$

- Linear regression:

$$\nabla_{\boldsymbol{w}} R(\boldsymbol{w}) = \mathbf{X}^T \mathbf{X} \boldsymbol{w} - \mathbf{X}^T \mathbf{y} \quad \nabla_{\boldsymbol{w}}^2 R(\boldsymbol{w}) = \mathbf{X}^T \mathbf{X}$$

$$\boldsymbol{w}^{\text{new}} = \boldsymbol{w}^{\text{old}} - (\mathbf{X}^T \mathbf{X})^{-1} (\mathbf{X}^T \mathbf{X} \boldsymbol{w}^{\text{old}} - \mathbf{X}^T \mathbf{y})$$

$$\boldsymbol{w}^{\text{new}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

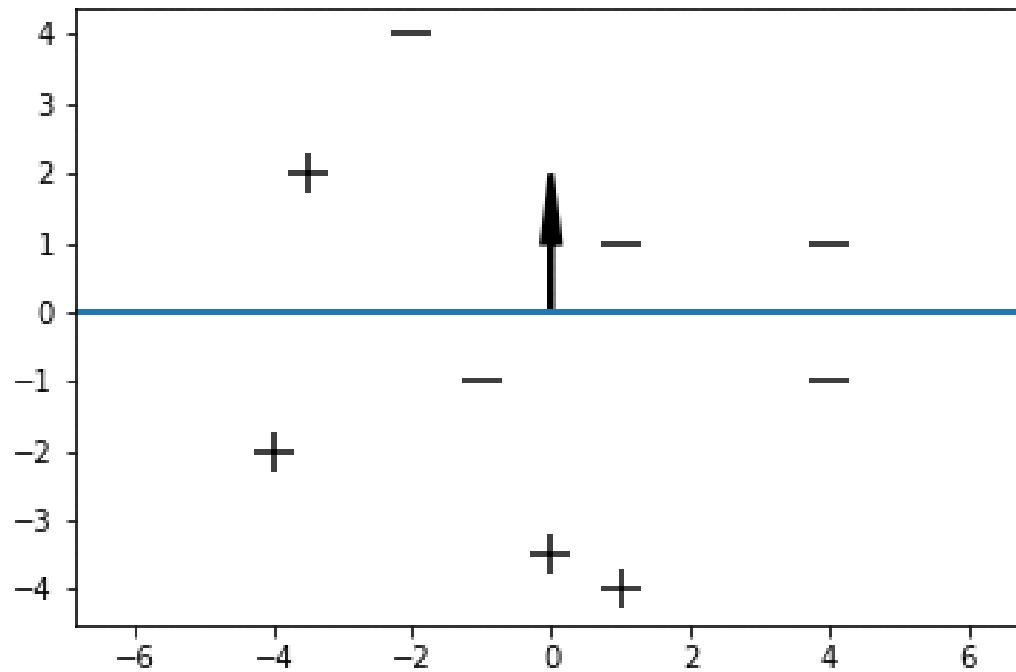
- Logistic regression:

$$\nabla_{\boldsymbol{w}} R(\boldsymbol{w}) = \mathbf{X}^T (\mathbf{a} - \mathbf{y}) \quad \nabla_{\boldsymbol{w}}^2 R(\boldsymbol{w}) = -\mathbf{X}^T \mathbf{R} \mathbf{X}$$

$$\boldsymbol{w}^{\text{new}} = \boldsymbol{w}^{\text{old}} - (\mathbf{X}^T \mathbf{R} \mathbf{X})^{-1} (\mathbf{X}^T (\mathbf{a} - \mathbf{y}))$$

This iterative procedure is called *iterative reweighted least squares* because the matrix \mathbf{R} gives new “weights” in each iteration.

Logistic regr: non-separable data



Logistic regression: example

- In 2019, the Austrian job center announced that a algorithm would be used to assess job-seekers
- The model is a logistic regression model using several person-centric variables
- The algorithm was show to be unfair, e.g. w.r.t. women

Anhang 1

Das vorliegende Modell ist ein logistisches Regressionsmodell, das die kurzfristige Integrationschance (d.h. die Wahrscheinlichkeit, 90 Beschäftigungstage in 7 Monaten zu erreichen) zu Geschäftsfallbeginn in Abhängigkeit von 12 erklärenden Variablen¹ und Interaktionen zwischen diesen schätzt.

Das Modell enthält alle paarweisen Interaktionsterme, alle Interaktionsterme höherer Ordnung zwischen den Merkmalen GESCHLECHT, ALTER, BERUFSGRUPPE, AUSBILDUNG und STAATENGRUPPE sowie alle Interaktionsterme höherer Ordnung zwischen den Merkmalen GESCHLECHT, ALTER, BEEINTRÄCHТИGUNG, BESCHAFTIGUNGSVERLAUF, und GESCHÄFTSFALL-FREQUENZ.

Overview

- 4. Basic neural network architectures
 - 4.1 Linear model: Neuron
 - 4.2 Perceptron
 - 4.3 Linear neuron with sigmoid activation: logistic regression
 - 4.4 Linear neuron with softmax activation: softmax regression

4.4 Softmax regression

Softmax: classification task with multiple classes

- Let us assume that the targets/labels can have 1 of K classes
- This is a supervised classification tasks
- Examples for classification tasks
 - Image recognition: cars, dogs, cats, airplanes, ...
- What can we do in this situation?
Should we just use LR as before?

X			y
x^1	0.13	0.03	A
x^2	0.15	0.14	A
	0.79	0.19	C
	0.48	0.28	D
	0.93	0.43	C
	0.43	0.41	B
	0.62	0.63	B
	0.69	0.69	A
	0.19	0.79	B
	0.23	0.85	C
x^n	0.11	0.99	B

One-hot encoding

X					y		
x^1	0.13	0.03	A	1	0	0	0
x^2	0.15	0.14	A	1	0	0	0
	0.79	0.19	C	0	0	1	0
	0.48	0.28	D	0	0	0	1
	0.93	0.43	C	0	0	1	0
	0.43	0.41	B	0	1	0	0
	0.62	0.63	B	0	0	0	0
	0.69	0.69	A	1	0	0	0
	0.19	0.79	B	0	1	0	0
	0.23	0.85	C	0	0	1	0
x^n	0.11	0.99	B	0	1	0	0

Softmax regression: model

$$g(\mathbf{x}, \mathbf{W}) = \mathbf{a} = \text{softmax}(\mathbf{W}\mathbf{x})$$

- We have a weight matrix instead of a weight vector
- The softmax function is:

$$\text{softmax}(\mathbf{s}) = \left(\frac{e^{s_1}}{\sum_{j=1}^K e^{s_j}}, \dots, \frac{e^{s_k}}{\sum_{j=1}^K e^{s_j}}, \dots, \frac{e^{s_K}}{\sum_{j=1}^K e^{s_j}} \right)^T$$

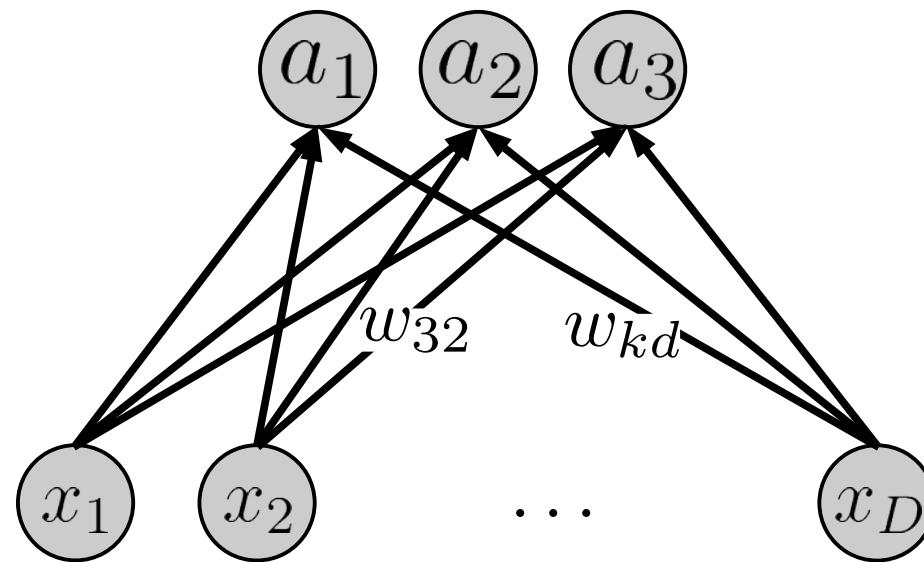
- The resulting vector is a probability vector:

$$\mathbf{a} = \text{softmax}(\mathbf{s})$$

with positive entries that sum to one.

$$g(\mathbf{x}, \mathbf{W}) = \mathbf{a} = \text{softmax}(\mathbf{W}\mathbf{x})$$

Softmax regression: net



$$g(\mathbf{x}, \mathbf{W}) = \mathbf{a} = \text{softmax}(\mathbf{W}\mathbf{x})$$

Softmax regr: loss and error

- Likelihood of label vector:

$$p(\mathbf{y}|\mathbf{a}) = \prod_{k=1}^K a_k^{y_k}$$

- Categorical cross-entropy loss:

$$-\log p(\mathbf{y}|\mathbf{a}) = -\sum_{k=1}^K y_k \log(a_k)$$

- Empirical error:

$$R_{\text{emp}}(\mathbf{W}, \mathbf{X}, \mathbf{Y}) = -\sum_{n=1}^N \sum_{k=1}^K y_{nk} \log(a_{nk}) = -\sum_{n=1}^N \sum_{k=1}^K (\mathbf{y}^n)_k \log((\mathbf{a}^n)_k)$$

$$g(\mathbf{x}, \mathbf{W}) = \mathbf{a} = \text{softmax}(\mathbf{W}\mathbf{x})$$

Softmax regression: learning

- Derivative w.r.t. weight:

$$\begin{aligned} \frac{\partial}{\partial w_{\xi d}} \left(- \sum_{k=1}^K y_k \log(a_k) \right) &= - \frac{\partial}{\partial w_{\xi d}} \sum_{k=1}^K y_k \log \left(\frac{e^{\mathbf{w}_k^T \mathbf{x}}}{\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}}} \right) \\ &= - \frac{\partial}{\partial s_k} \frac{\partial s_k}{\partial w_{\xi d}} \sum_{k=1}^K y_k \log \left(\frac{e^{s_k}}{\sum_{j=1}^K e^{s_j}} \right) \\ &= (a_\xi - y_\xi) x_d \end{aligned}$$

- Derivative w.r.t. to weight matrix:

$$\nabla_{\mathbf{W}} R_{\text{emp}}(\mathbf{W}, \mathbf{x}, \mathbf{y}) = (\mathbf{a} - \mathbf{y}) \mathbf{x}^T$$

Softmax regression learning

- Again, we can use gradient descent for learning
- Type of optimization problem?
- Convergence?

Comparing gradients

- Linear regression:

$$\nabla_{\mathbf{w}} R_{\text{emp}}(y, \mathbf{x}, \mathbf{w}) = (\mathbf{w}^T \mathbf{x} - y) \mathbf{x}$$

- Logistic regression:

$$\nabla_{\mathbf{w}} R_{\text{emp}}(y, \mathbf{x}, \mathbf{w}) = (\sigma(\mathbf{w}^T \mathbf{x}) - y) \mathbf{x}$$

- Softmax regression

$$\nabla_{\mathbf{W}} R_{\text{emp}}(\mathbf{y}, \mathbf{x}, \mathbf{W}) = (\text{softmax}(\mathbf{W} \mathbf{x}) - \mathbf{y}) \mathbf{x}^T$$

Summary

- We have discussed single-layer networks
 - Linear neuron, linear regression
 - Perceptron
 - Linear neuron with sigmoid activation, log regr
 - Softmax regression
- Learning, convergence and solutions
- For many applications, these networks are sufficient