

Introduction to Neural Networks

Mathematical Tools for ITS (11MAI)

Mathematical tools, 2020

Jan Přikryl

11MAI, lecture 10

Monday, December 07, 2020

version: 2020-12-07 10:05

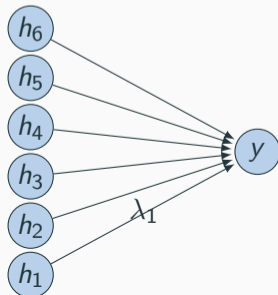
Department of Applied Mathematics, CTU FTS

- We used before weighted linear combination of feature values h_j and weights λ_j

$$y = \sigma(\lambda, d_i) = \sum_j \lambda_j h_j(d_i)$$

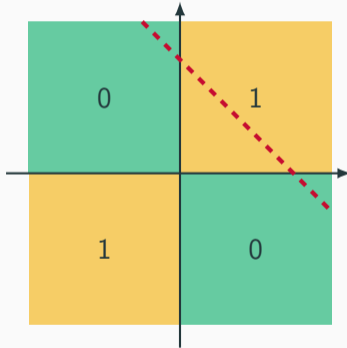
- Such models can be illustrated as a “network”

λ_6

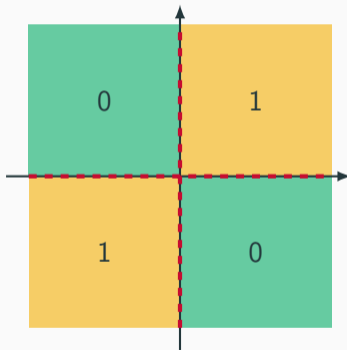


- We can give each feature a weight
- But not more complex value relationships, e.g.,
 - any value in the range $[0; 5]$ is equally good
 - values over 8 are bad
 - higher than 10 is not worse

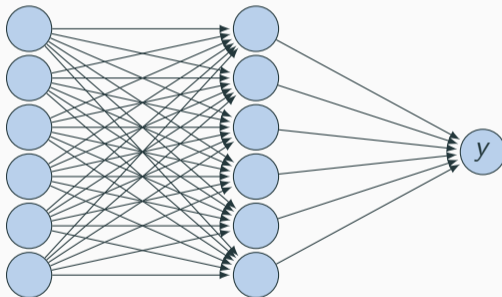
Linear models cannot model XOR-like behaviour:



Linear models cannot model XOR-like behaviour:



Solution: Add an intermediate (“hidden”) layer of processing (each arrow is a weight)



- Have we gained anything so far?

- Instead of computing a linear combination

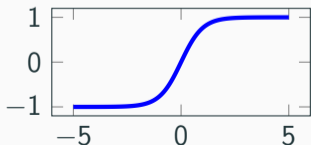
$$\sigma(\lambda, d_i) = \sum_j \lambda_j h_j(d_i)$$

- Add a non-linear function

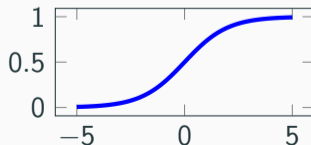
$$\sigma(\lambda, d_i) = f \left(\sum_j \lambda_j h_j(d_i) \right)$$

- Popular choices (“sigmoid” \equiv the logistic function)

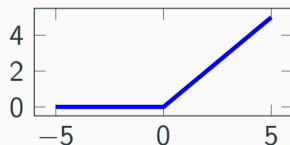
tanh(x)



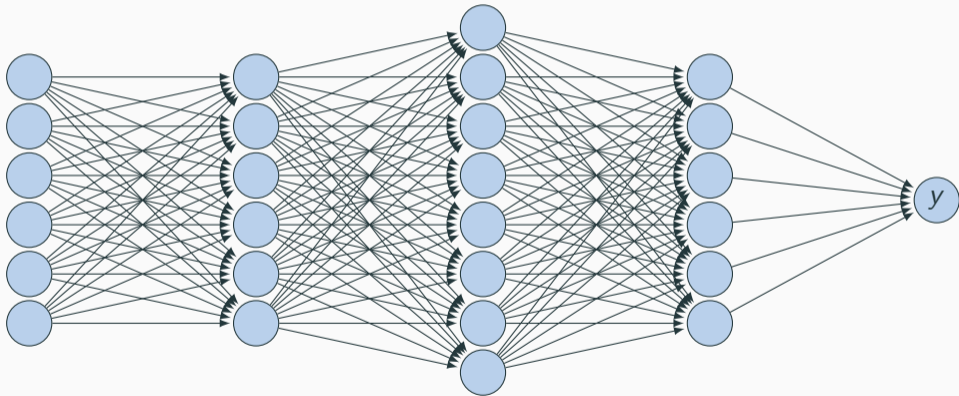
sgmd(x)



relu(x)



More layers = deep learning

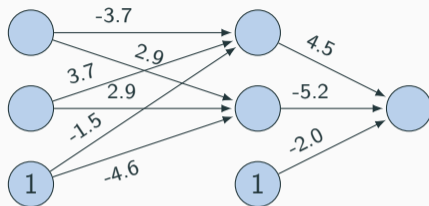


We can interpret the deep NN as follows:

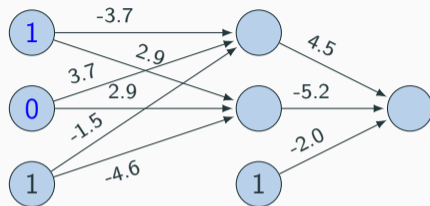
- Each layer is a processing step
- Having multiple processing steps allows complex functions
- Metaphor: NN and computing circuits
 - computer = sequence of Boolean gates
 - neural computer = sequence of layers
- Deep neural networks can implement complex functions e.g. sorting on input values

But in fact, a trained NN is just a clever **lookup table**.

Example



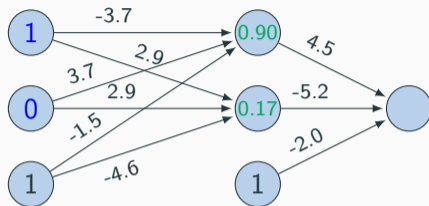
- One innovation: bias units (no inputs, always value 1)



- Try out two input values
- Hidden unit computation

$$\text{sgmd}(1.0 \times 3.7 + 0.0 \times 3.7 + 1 \times -1.5) = \text{sgmd}(2.2) = \frac{1}{1 + e^{-2.2}} = 0.90$$

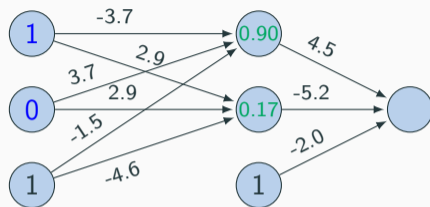
$$\text{sgmd}(1.0 \times 2.9 + 0.0 \times 2.9 + 1 \times -4.5) = \text{sgmd}(-1.6) = \frac{1}{1 + e^{-1.6}} = 0.17$$



- Try out two input values
- Hidden unit computation

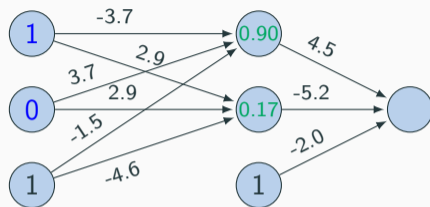
$$\text{sgmd}(1.0 \times 3.7 + 0.0 \times 3.7 + 1 \times -1.5) = \text{sgmd}(2.2) = \frac{1}{1 + e^{-2.2}} = 0.90$$

$$\text{sgmd}(1.0 \times 2.9 + 0.0 \times 2.9 + 1 \times -4.5) = \text{sgmd}(-1.6) = \frac{1}{1 + e^{-1.6}} = 0.17$$



- Output unit computation

$$\text{sgmd}(0.90 \times 4.5 + 0.17 \times (-5.2) + 1 \times (-2.0)) = \text{sgmd}(1.17) = \frac{1}{1 + e^{-1.17}} = 0.76$$



- Output unit computation

$$\text{sgmd}(0.90 \times 4.5 + 0.17 \times (-5.2) + 1 \times (-2.0)) = \text{sgmd}(1.17) = \frac{1}{1 + e^{-1.17}} = 0.76$$

Input x_0	Input x_1	Hidden h_0	Hidden h_1	Output y_0
0	0	0.12	0.02	0.18 \rightarrow 0
0	1	0.88	0.27	0.74 \rightarrow 1
1	0	0.73	0.12	0.74 \rightarrow 1
1	1	0.99	0.73	0.33 \rightarrow 0

- Network implements XOR
 - hidden node h_0 is OR
 - hidden node h_1 is AND
 - final layer operation is $h_0 - (-h_1)$
- Power of deep neural networks: chaining of processing steps just as: more Boolean circuits \rightarrow more complex computations possible

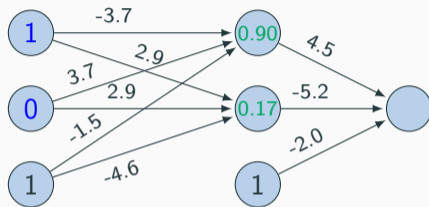
Why “neural” networks?

- The human brain is made up of about 100 billion neurons
 - Dendrite
 - Axon terminal
 - Soma
 - Axon Nucleus
- Neurons receive electric signals at the dendrites and send them to the axon

- The axon of the neuron is connected to the dendrites of many other neurons
Neurotransmitter
Synaptic vesicle Neurotransmitter transporter Axon terminal Voltage gated Ca^{++} channel
Receptor Postsynaptic density Synaptic cleft
Dendrite

- Similarities
 - Neurons, connections between neurons
 - Learning = change of connections, not change of neurons
 - Massive parallel processing
- But artificial neural networks are much simpler
 - computation within neuron vastly simplified
 - discrete time steps
 - typically some form of supervised learning with massive number of stimuli

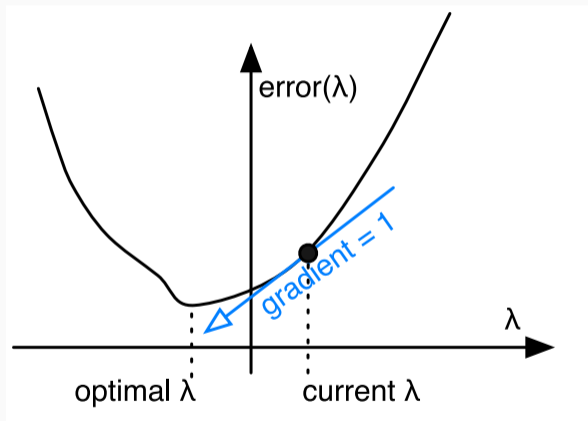
back-propagation training

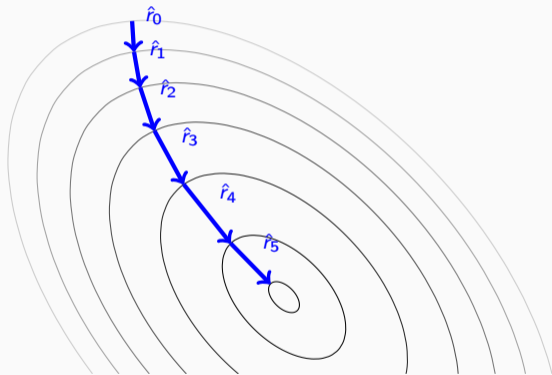


The output is not exact

- Computed output: $y = 0.76$
- Correct output: $t = 1.0$
- How do we adjust the weights?

- Gradient descent
 - error is a function of the weights
 - we want to reduce the error
 - gradient descent: move towards the error minimum
 - compute gradient \rightarrow get direction to the error minimum
 - adjust weights towards direction of lower error
- Back-propagation
 - first adjust last set of weights
 - propagate error back to each previous layer
 - adjust their weights





- Sigmoid function: $\text{sgmd}(x) = \frac{1}{1 + e^{-x}}$
- Derivative

$$\begin{aligned}\frac{d}{dx} \text{sgmd}(x) &= \frac{d}{dx} \frac{1}{1 + e^{-x}} \\ &= \frac{(1 - e^{-x}) \times 0 - 1 \times (-e^{-x})}{(1 + e^{-x})^2} \\ &= \frac{1}{1 + e^{-x}} \cdot \frac{e^{-x}}{1 + e^{-x}} \\ &= \frac{1}{1 + e^{-x}} \left(1 - \frac{1}{1 + e^{-x}} \right) \\ &= \text{sgmd}(x) (1 - \text{sgmd}(x))\end{aligned}$$

- Linear combination of weights $s = \sum_k w_k h_k$
- Activation function $y = \text{sgmd}(s)$
- Error (L2 norm) $E = (t - y)^2/2$
- Derivative of error with regard to one weight w_k

$$\frac{dE}{dw_k} = \frac{dE}{dy} \frac{dy}{ds} \frac{ds}{dw_k}$$

- Linear combination of weights $s = \sum_k w_k h_k$
- Activation function $y = \text{sgmd}(s)$
- Error (L2 norm) $E = (t - y)^2/2$
- Derivative of error with regard to one weight w_k

$$\frac{dE}{dw_k} = \frac{dE}{dy} \frac{dy}{ds} \frac{ds}{dw_k}$$

- Error E is defined with respect to y

$$\frac{dE}{dy} = \frac{d}{dy} \frac{1}{2} (t - y)^2 = -(t - y)$$

- Linear combination of weights $s = \sum_k w_k h_k$
- Activation function $y = \text{sgmd}(s)$
- Error (L2 norm) $E = (t - y)^2/2$
- Derivative of error with regard to one weight w_k

$$\frac{dE}{dw_k} = \frac{dE}{dy} \frac{dy}{ds} \frac{ds}{dw_k}$$

- y with respect to x is $\text{sgmd}(s)$:

$$\frac{dy}{ds} = \frac{d}{ds} \text{sgmd}(s) = \text{sgmd}(s) (1 - \text{sgmd}(s)) = y (1 - y)$$

- Linear combination of weights $s = \sum_k w_k h_k$
- Activation function $y = \text{sgmd}(s)$
- Error (L2 norm) $E = (t - y)^2/2$
- Derivative of error with regard to one weight w_k

$$\frac{dE}{dw_k} = \frac{dE}{dy} \frac{dy}{ds} \frac{ds}{dw_k}$$

- x is a weighted linear combination of hidden node values h_k

$$\frac{ds}{dw_k} = \frac{d}{dw_k} \sum_k w_k h_k = h_k$$

- Derivative of error with regard to one weight w_k

$$\frac{dE}{dw_k} = \frac{dE}{dy} \frac{dy}{ds} \frac{ds}{dw_k} = -(t - y)y(1 - y)h_k$$

- error
- derivative of sigmoid: y'
- Weight adjustment will be scaled by a fixed learning rate μ :

$$\Delta w_k = \mu (t - y) y' h_k$$

- Our example only had one output node
- Typically neural networks have multiple output nodes
- Error is computed over all j output nodes

$$E = \sum_j \frac{1}{2} (t_j - y_j)^2$$

- Weights $k \rightarrow j$ are adjusted according to the node they point to

$$\Delta w_{j \leftarrow k} = \mu (t_j - y_j) y_j' h_k$$

- In a hidden layer, we do not have a target output value
- But we can compute how much each node contributed to downstream error
- Definition of error term of each node

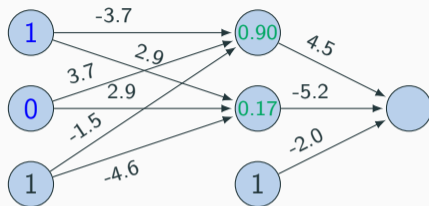
$$\delta_j = (t_j - y_j)y_j'$$

- Back-propagate the error term (why this way? there is math to back it up ...)

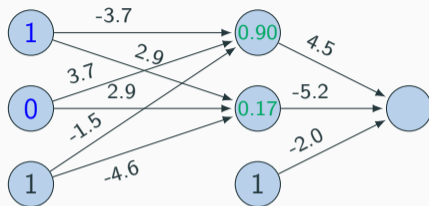
$$\delta_i = \left(\sum_j w_{j \leftarrow i} \delta_j \right) y_i'$$

- Universal update formula

$$\Delta w_{j \leftarrow k} = \mu \delta_j h_k$$



- Computed output: $y = 0.76$
- Correct output: $t = 1.0$
- Final layer weight updates (learning rate $\mu = 10$)
 - $\delta_G = (t - y)y' = (1 - 0.76) \times 0.181 = 0.0434$
 - $\Delta w_{GD} = \mu \delta_G h_D = 10 \times 0.0434 \times 0.90 = 0.391$
 - $\Delta w_{GE} = \mu \delta_G h_E = 10 \times 0.0434 \times 0.17 = 0.074$
 - $\Delta w_{GF} = \mu \delta_G h_F = 10 \times 0.0434 \times 1 = 0.434$



- Hidden node **D**

- $\delta_D = \left(\sum_j w_{j \leftarrow i} \delta_j \right) y'_D = w_{GD} \delta_G y'_D = 4.5 \times .0434 \times .0898 = .0175$
- $\Delta w_{DA} = \mu \delta_D h_A = 10 \times 0.0175 \times 1.0 = 0.175$
- $\Delta w_{DB} = \mu \delta_D h_B = 10 \times 0.0175 \times 0.0 = 0$
- $\Delta w_{DC} = \mu \delta_D h_C = 10 \times 0.0175 \times 1 = 0.175$

- Hidden node **E**

- $\delta_E = \left(\sum_j w_{j \leftarrow i} \delta_j \right) y'_E = w_{GE} \delta_G y'_E = -5.2 \times 0.0434 \times 0.2055 = -0.0464$
- $\Delta w_{EA} = \mu \delta_E h_A = 10 \times -0.0464 \times 1.0 = -0.464$

Some additional aspects

- Weights are initialized randomly, e.g. uniformly from interval $[-0.01, 0.01]$
- Glorot and Bengio (2010) suggest

- for shallow neural networks

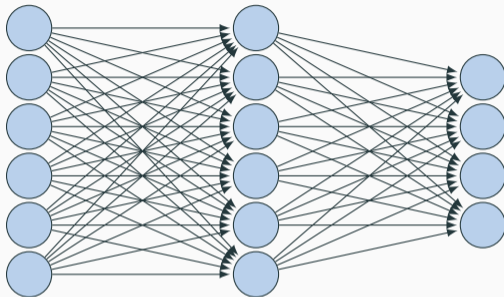
$$\left[-\frac{1}{\sqrt{n}}, \frac{1}{\sqrt{n}} \right]$$

where n is the size of the previous layer

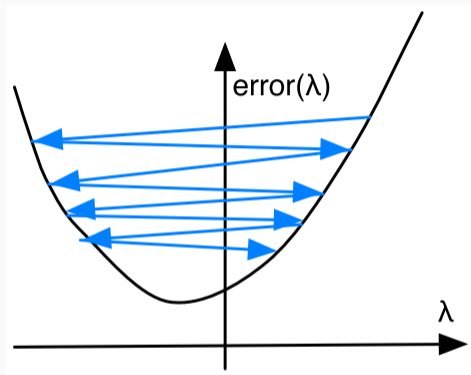
- for deep neural networks

$$\left[-\frac{\sqrt{6}}{\sqrt{n_j + n_{j+1}}}, \frac{\sqrt{6}}{\sqrt{n_j + n_{j+1}}} \right]$$

n_j is the size of the previous layer, n_{j+1} the size of the next layer

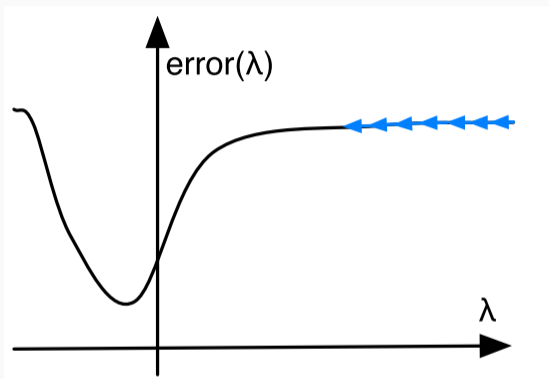


- Predict class: one output node per class
- Training data output: “One-hot vector”, e.g. $\mathbf{y} = (0, 0, 1)^T$
- Prediction
 - predicted class is output node y_i with highest value
 - obtain posterior probability distribution by soft-max, $\text{softmax}(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}}$



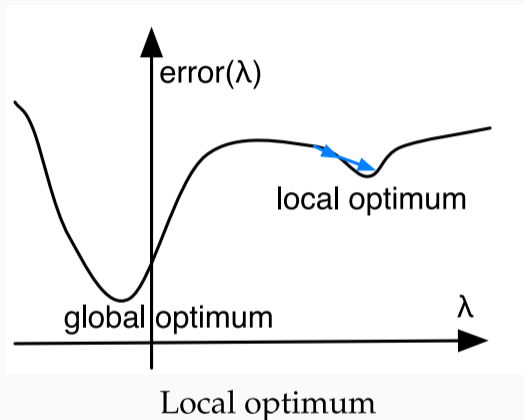
Too high learning rate

Too high learning rate



Bad initialization

Bad initialization



Local optimum

- Updates may move a weight slowly in one direction
- To speed this up, we can keep a memory of prior updates ...

$$\Delta w_{j \leftarrow k}(n-1)$$

- ...and add these to any new updates (with decay factor ρ)

$$\Delta w_{j \leftarrow k}(n) = \mu \delta_j h_k + \rho \Delta w_{j \leftarrow k}(n-1)$$

- Typically reduce the learning rate μ over time
 - at the beginning, things have to change a lot
 - later, just fine-tuning
- Adapting learning rate per parameter
- Adagrad update: based on error E with respect to the weight w at time $t = g_t = \frac{dE}{dw}$

$$\Delta w_t = \frac{\mu}{\sqrt{\sum_{\tau=1}^t g_{\tau}^2}} g_t$$

- A general problem of machine learning: overfitting to training data (very good on train, bad on unseen test)
- Solution: *regularization*, e.g., keeping weights from having extreme values
- Dropout: randomly remove some hidden units during training
 - mask: set of hidden units dropped
 - randomly generate, say, 10–20 masks
 - alternate between the masks during training
- Why does that work? → bagging, ensemble, ...

- Each training example yields a set of weight updates Δw_i .
- Batch up several training examples
 - sum up their updates
 - apply sum to model
- Mostly done for speed reasons

46

computational aspects

- Forward computation: $\mathbf{s} = \mathbf{W}\mathbf{h}$
- Activation function: $\mathbf{y} = \text{sgmd}(\mathbf{h})$
- Error term: $\boldsymbol{\delta} = (\mathbf{t} - \mathbf{y}) \cdot \text{sgmd}(\mathbf{s})'$
- Propagation of error term: $\boldsymbol{\delta}_i = \mathbf{W}\boldsymbol{\delta}_{i+1} \cdot \text{sgmd}(\mathbf{s})'$
- Weight updates: $\Delta\mathbf{W} = \mu\boldsymbol{\delta}\mathbf{h}^\top$

- Neural network layers may have, say, 200 nodes
- Computations such as $\mathbf{s} = \mathbf{Wh}$ require $200 \times 200 = 40\,000$ multiplications
- Graphics Processing Units (GPU) are designed for such computations
 - Real-time graphics (projections, shading) requires fast vector and matrix operations
 - GPU has massive number of multi-core but lean processing units
 - *Example:* NVIDIA Tesla K20c GPU provides 2496 thread processors, NVIDIA Tesla V100 GPU provides 5120 of them + 640 tensor cores operating on 4×4 matrices
- Extensions to C to support programming of GPUs, such as CUDA
- MATLAB is able to offload computations to GPU if parallel toolbox is installed

- Theano
- Tensorflow (Google) — <https://playground.tensorflow.org/>
- PyTorch (Facebook)
- MXNet (Amazon)
- DyNet

MATLAB: Deep Learning Toolbox