

ASSIGNMENT 7: FIRST NEURAL NETWORKS AND A GLIMPSE AT PYTORCH



Institute for Machine Learning

Copyright statement:

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.

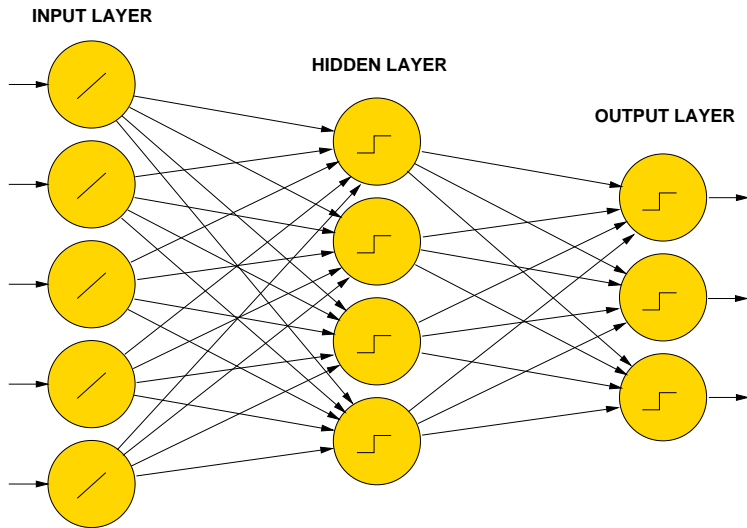
Perceptrons

- Introduced by Frank Rosenblatt in 1958
- A perceptron is a simple linear threshold unit:

$$g(\mathbf{x}_i; \mathbf{w}, \theta) = \begin{cases} 1 & \text{if } \mathbf{w}^T \mathbf{x}_i > \theta \\ 0 & \text{otherwise} \end{cases}$$

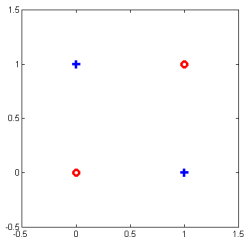
- In analogy to the biological model:
 - **inputs** $\mathbf{x}_i \Rightarrow$ charges received from connected cells
 - **weights** $\mathbf{w} \Rightarrow$ properties of the synaptic interface
 - **output** \Rightarrow impulse that is sent through the axon as soon as the charge exceeds the threshold θ
- Though it seems to be a (simplistic) model of a neuron, a **perceptron is nothing else but a simple linear classifier.**

Multi-layer perceptrons





Learning non-linear functions



x_1	x_2	y
0	0	0
0	1	1
1	0	1
1	1	0

This classification problem:

- ... can't be solved by one-layer network (single-layer perceptron, logistic regression)
- ... is trivial to solve using an MLP.



How powerful are neural networks?

- Useful measure of a machine learning algorithm: what are the **most complicated functions it can learn**?
- Measured by the **VC dimension**¹
- NNs can learn (much) more complex decision functions.
- **Neural Networks are “Universal Function Approximators”.**

¹**Vapnik-Chervonenkis** dimension – to learn more, see “Statistical Learning Theory” (Vladimir N. Vapnik)

NNs: Basic definitions and notation: Part 1



- a_i : **activity** of the i -th unit ($a_0 = 1$)
- w_{ij} : **weight** from unit j to unit i ($w_{i0} = b_i$ called **bias weight**)
- W : number of weights
- Q : number of units
- D : number of input units ($1 \leq i \leq D$) located in the first layer called **input layer**.
- K : number of output units ($Q - K + 1 \leq i \leq Q$) located in the last layer called **output layer**.
- H : number of hidden units ($D < i \leq Q - K$) located in the hidden layers.
- L : number of layers, where L_ν is the index set of the ν -th layer;
 $L_1 = \{1, \dots, D\}$ and $L_L = \{Q - K + 1, \dots, Q\}$.
- s_i : **network input or pre-activation** of the i -th unit ($i > D$):
$$s_i = \sum_{j=0}^Q w_{ij} a_j$$



NNs: Basic definitions and notation: Part 2

- f : **activation function** with $a_i = f(s_i)$ It is possible to define different activation functions f_i for different units.
- **Architecture** of a NN is given through:
 - ☐ number of layers
 - ☐ number of units in the layers
 - ☐ connections between units
 - ☐ also the activation function may be accounted to architecture
- A feed-forward NN has only connections from units in lower layers to units in higher layers:

$$j \in L_\nu \text{ and } i \in L_{\nu'} \text{ and } \nu' \leq \nu \Rightarrow w_{ij} = 0.$$

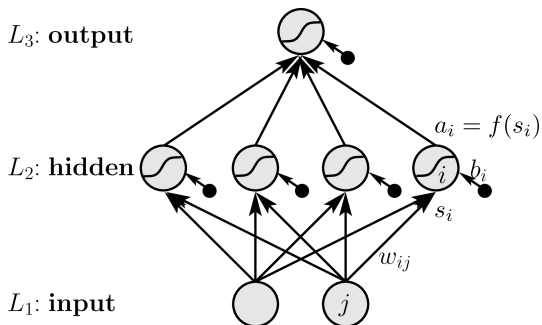
- For conventional NN: only connections or weights between consecutive layers. Other weights are fixed to zero.
- Connections between units in layers which are not adjacent are called **shortcut connections**.



NNs: Forward pass

1. Provide input x
2. Input layer $\nu = 1$: for $i = 1$ to D : $a_i = x_i$
3. Further layers: for $\nu = 2$ to L
 - For i in L_ν :
 - $s_i = \sum_{j=0}^Q w_{ij} a_j$
 - $a_i = f(s_i)$
4. Output layer: for all $Q - K + 1 \leq i \leq Q$:
 - provide output $\hat{y}_i = g(x_i; w) = a_i$

NNs: Visualization



Back-propagation



Calculating gradients: Back-propagation

Part 1

- Similar as in Unit 4 and Unit 5 (logistic regression and gradient boosting): Expected Risk Minimization (ERM) by gradient descent (GD).
- Gradient of NN can be computed efficiently by **back-propagation**. Popular since mid 80ies.
- Empirical risk:

$$R_{\text{emp}}(\mathbf{w}, \mathbf{X}, \mathbf{Y}) = \frac{1}{N} \sum_{n=1}^N L(\mathbf{y}^n, \mathbf{g}(\mathbf{x}^n; \mathbf{w})).$$

\mathbf{w} now contains all adjustable weights of NN (i.e. all w_{ij} and b_i). Generalizes to multiple outputs \mathbf{y} . Superscript n labels entries from the data set of size N .

- GD update: $\mathbf{w}^{\text{new}} = \mathbf{w}^{\text{old}} - \eta \nabla_{\mathbf{w}} R_{\text{emp}}(\mathbf{w}, \mathbf{X}, \mathbf{Y})$.



Calculating gradients: Back-propagation

Part 2

- Select w_{ij} in arbitrary layer that connects from the unit j to unit i and calculate $\frac{\partial}{\partial w_{ij}} L(\mathbf{y}^n, \mathbf{g}(\mathbf{x}^n; \mathbf{w}))$ for all w_{ij} .
- Rewrite this as follows:

$$\begin{aligned}\frac{\partial}{\partial w_{ij}} L(\mathbf{y}^n, \mathbf{g}(\mathbf{x}^n; \mathbf{w})) &= \frac{\partial}{\partial s_i} L(\mathbf{y}^n, \mathbf{g}(\mathbf{x}^n; \mathbf{w})) \frac{\partial s_i}{\partial w_{ij}} \\ &= \frac{\partial}{\partial s_i} L(\mathbf{y}^n, \mathbf{g}(\mathbf{x}^n; \mathbf{w})) a_j.\end{aligned}$$

- Define δ -error at unit i as $\delta_i := \frac{\partial}{\partial s_i} L(\mathbf{y}^n, \mathbf{g}(\mathbf{x}^n; \mathbf{w}))$.
- Yields: $\frac{\partial}{\partial w_{ij}} L(\mathbf{y}^n, \mathbf{g}(\mathbf{x}^n; \mathbf{w})) = \delta_i a_j$.



Calculating gradients: Back-propagation

Part 3: Delta errors at the output layers.

- The δ -error at the output units, denoted as δ_k with preactivation s_k and activation $a_k = g(\mathbf{x}^n; \mathbf{w})$ for $Q - K + 1 \leq k \leq Q$ is $\delta_k = \frac{\partial}{\partial a_k} L(\mathbf{y}^n, g(\mathbf{x}^n; \mathbf{w})) f'(s_k)$
- f' : derivative of activation function f of output layer.
- f in output layer is usually different than in hidden layers.
- Expression typically results in a difference between network outputs, i.e. predicted labels, and true labels.

Loss function		output activation	delta-error
squared loss	$1/2(a - y)^2$	linear: $f(x) = x$	$\delta = a - y$
binary CE	$-y \log(a) - (1 - y) \log(1 - a)$	sigmoid: $f(x) = \frac{1}{1+e^{-x}}$	$\delta = a - y$
categorical CE	$-\sum_{k=1}^K y_k \log(a_k)$	softmax: $f(\mathbf{x}) = \text{softmax}(\mathbf{x})$	$\delta = \mathbf{a} - \mathbf{y}$



Calculating gradients: Back-propagation

Part 4: Delta errors at hidden layers.

- The δ -error at units not in the output layer, i.e. hidden layers, reads:

$$\begin{aligned}\delta_j &= \frac{\partial}{\partial s_j} L(\mathbf{y}^n, \mathbf{g}(\mathbf{x}^n; \mathbf{w})) = \sum_i \frac{\partial}{\partial s_i} L(\mathbf{y}^n, \mathbf{g}(\mathbf{x}^n; \mathbf{w})) \frac{\partial s_i}{\partial s_j} \\ &= f'(s_j) \sum_i \delta_i w_{ij},\end{aligned}$$

- This is a recursive formula (lower layer unit error as function of above layer unit errors), allowing us to calculate all δ -errors of the network.
- \sum_i goes over all i for which w_{ij} exists.
- Typically, i goes over all units in the layer *above* the layer where unit j is located.



Calculating gradients: Back-propagation

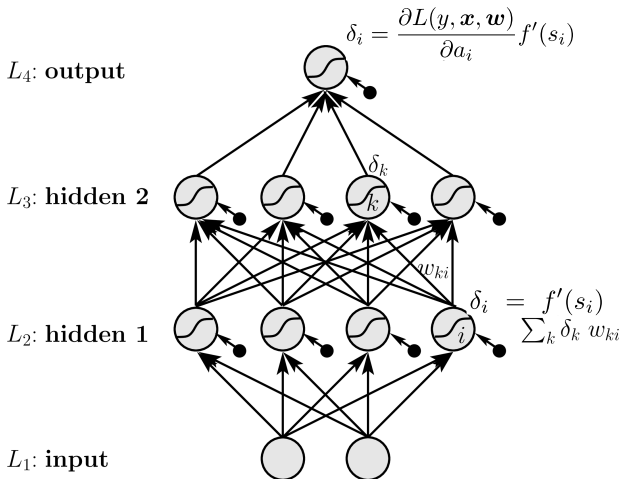
Part 5: Recap of the algorithm.

- Select a sample x randomly from the training set and perform a forward-pass. Memorize activations and preactivations of all neurons.
- Calculate delta-errors at output units according to Part 3.
- Calculate delta-errors for hidden layers starting from the hidden layers closest to the output layer, thus backpropagating the error signal according to Part 4.
- Calculate gradients for weights as done in Part 2.
- Update the weights by performing a gradient descent step using the weight changes: $w_{ij}^{\text{new}} = w_{ij}^{\text{old}} - \eta \delta_i a_j$. The learning rate η must be chosen appropriately for the algorithm to converge. (Typical learning rates are $\eta = 0.1$ or $\eta = 0.01$).



Calculating gradients: Back-propagation

Part 6: Visualization of 4-layer NN.

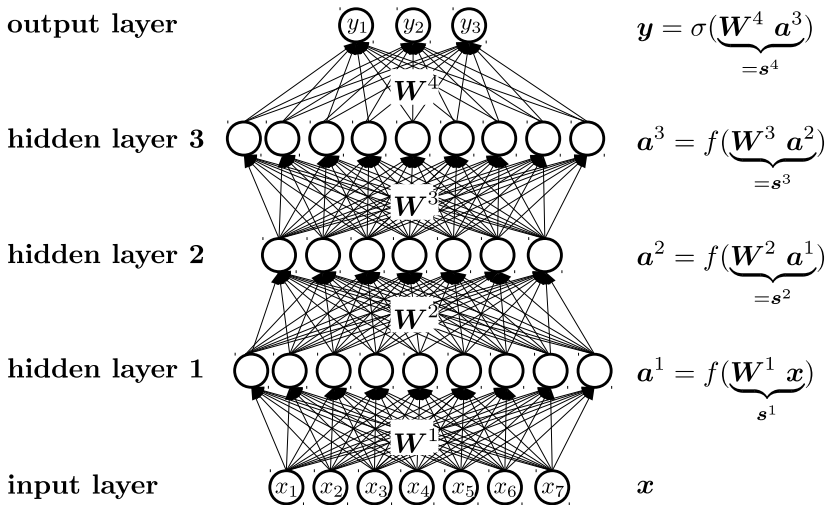




Nowadays used notation for NNs: Matrix-vector notation

- First introduce the following expressions:
 - x : input for layer 0; equivalent to $a^{[0]}$
 - $W^{[l]}$: weight matrix connecting layer $l - 1$ and layer l
 - $s^{[l]}$: pre-activations of layer l
 - $a^{[l]}$: activations of layer l
 - f : activation function applied element-wise to vector.
- Given inputs $x^{[l]}$ to a layer l , i.e. the activations $a^{[l-1]}$ of the lower layer $l - 1$, a NN computes activations of next layer as follows: $a^{[l]} = f(s^{[l]}) = f(W^{[l]}a^{[l-1]})$
- Function from the input layer x to output layer: $\hat{y} = g(x; W^{[1]}, \dots, W^{[L]}) = \sigma(W^{[L]}(\dots f(W^{[2]}f(W^{[1]}x)) \dots))$,
- \hat{y} are activations of output layer and x is input of network. Outermost activation function σ typically different from activation function f in hidden layers.

Schematic view of a deep feed-forward neural network in matrix-vector notation





Backpropagation in Matrix-vector notation

- Now we can write:

$$\frac{\partial}{\partial w_{ij}^{[1]}} L(\mathbf{y}, g(\mathbf{x}; \mathbf{W}^{[1]}, \dots, \mathbf{W}^{[L]})) = \frac{\partial L(\mathbf{y}, \hat{\mathbf{y}})}{\partial \hat{\mathbf{y}}} \frac{\partial \hat{\mathbf{y}}}{\partial \mathbf{s}^{[L]}} \cdots \frac{\partial \mathbf{s}^{[l]}}{\partial \mathbf{s}^{[l-1]}} \cdots \frac{\partial \mathbf{s}^{[1]}}{\partial w_{ij}^{[1]}}.$$

- Introduce delta errors as follows: $\delta^{[l]} = \left(\frac{\partial L(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{s}^{[l]}} \right)^T$.
- Note that we typically use these deltas $\delta^{[l]T}$ as **row vectors**.
- Similar arguments as before lead to recursive formula for delta errors:

$$\delta^{[l-1]T} = \frac{\partial L(\mathbf{y}, \hat{\mathbf{y}})}{\partial \mathbf{s}^{[l]}} \frac{\partial \mathbf{s}^{[l]}}{\partial \mathbf{s}^{[l-1]}} = \delta^{[l]T} \underbrace{\mathbf{W}^{[l]} \text{diag} \left(f'(\mathbf{s}^{[l-1]}) \right)}_{\mathbf{J}^{[l]}}.$$



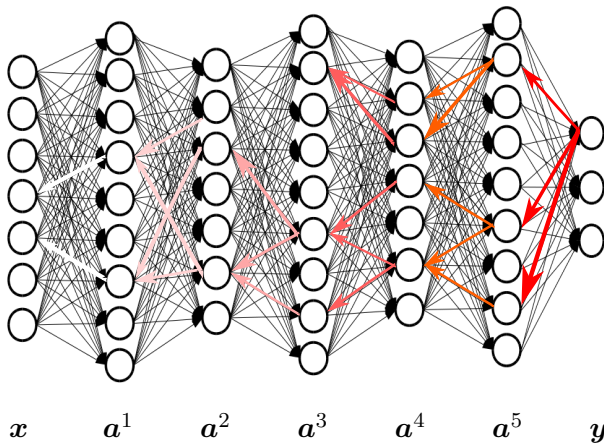
Vanishing and exploding gradients

- $\delta^{[l-1]T} = \delta^{[l]T} \mathbf{W}^{[l]} \text{diag} \left(f'(\mathbf{s}^{[l-1]}) \right)$:
- Important property that influences the ability of a NN to learn: size (norm) of the delta-errors that backpropagate through the network.
- Since for each layer delta errors are multiplied by Jacobian \rightarrow exponential behavior (growth or shrinkage):
 - $\|\delta^{[l-1]}\| < \|\delta^{[l]}\|$: **Vanishing gradients**: Has been typical case since the derivatives of the sigmoid activation function is at most $|f'(x)| = 0.25 \rightarrow \|\text{diag} \left(f'(\mathbf{s}^{[l-1]}) \right)\| \leq 0.25 \rightarrow$ norm of delta-errors becomes smaller through each layer (if not compensated by the norm of $\mathbf{W}^{[l]}$).
 - $\|\delta^{[l-1]}\| \approx \|\delta^{[l]}\|$: **Stable gradients**: Ideal case that the delta errors have a similar norm in each layer; many recent algorithmic improvements aim at keeping this quantity close to one (e.g. initialization strategies).
 - $\|\delta^{[l-1]}\| > \|\delta^{[l]}\|$: **Exploding gradients**: If norm of weight matrix is large \rightarrow norm of delta-errors grows through each layer \rightarrow unstable learning or even numeric overflows.



Vanishing Gradient: Visualization

- Gradient signal gets lost in the noise:





Example: ReLUs and ELUs

■ Rectified linear unit, by Nair and Hinton in 2010:

- Idea: N sigmoids with shared weights but different biases:

$$\sum_i^N \sigma(h - i + 0.5) \approx \log(1 + e^h) \approx \max(0, h)$$

where $h = \mathbf{w} \cdot \mathbf{x} + b$

- Gradient is either 0 or 1 (helps against vanishing gradients!)
- ReLU nets learn a piecewise linear function
- Problem: Dying ReLUs

■ Exponential linear unit, by Clevert, Unterthiner, Hochreiter in 2015:

