

— Training: algorithm —

• To operate a neural network, its parameters must first be determined; and to obtain the parameters, it must be trained: which training algorithm should be adopted? What are the main options regarding the network's hyperparameters and their configuration?

9a.1

1

▷ Training

• Parâmetros (recap)

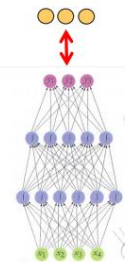
• weights of connections between neurons
• determined by training

• Training

• The goal is to find the function that maps the examples into respective desired labels
• that is, find the weights that support such a function

• we look for weights that minimize error or loss

• training is an optimization process



9a.2

2

Training algorithm

• 1. Forward pass

• make prediction

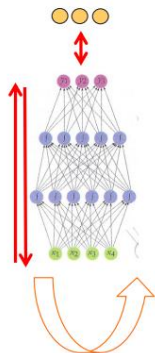
• 2. Estimate error/loss

• compare with desired value

• 3. Backward pass

• adjust the weights

• Repeat the above steps until the parameter values stabilize (without over-adjustment: see below)



9a.3

3

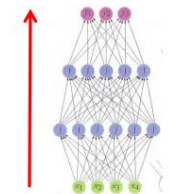
▷ Step 1: Forward pass para a frente ("forward pass")

• Only in the initial step:

Assign initial weights to connections between neurons (possibly randomly)

• Compute the prediction over the network

Go from the input layer to the output layer, taking an example training (or a batch of them) as input



9a.4

4

Network Computing (recap) 1/2 1/2

q Entrance

\vec{y} vector $\mathbf{x} = [x_1, \dots, x_n]$

q Weights

\vec{y} matrix \mathbf{W} with n neurons of the source layer in the rows and m neurons of the target layer in the columns

layer origin

[w11, ..., w1m]

[wn1, ..., wnm]

layer target

\vec{y} $W[i,j]$ is the connection weight from source neuron i to target neuron j

9a.5

5

Network Computing (recap) 2/2 2/2

q $\mathbf{h} = \mathbf{g}(\mathbf{x} \mathbf{W} + \mathbf{b})$ (fully connected layer)

\vec{y} \mathbf{h} is output vector

each component with the output of each neuron

\vec{y} \mathbf{g} is nonlinear activation function

e.g. ReLU – Rectifier Linear Unit, $g(x) = \max(0, x)$

\vec{y} \mathbf{x} is input vector and \mathbf{W} is transition matrix

\vec{y} $\mathbf{x} \mathbf{W} =$

$[x_1, \dots, x_n][w_{11}, \dots, w_{1m}] = [x_1.w_{11} + \dots + x_n.w_{n1}, \dots, x_1.w_{1m} + \dots + x_n.w_{nm}]$

[wn1, ..., wnm]

\vec{y} \mathbf{b} is bias term

9a.6

6

X is an input vector and W is a transition matrix

Step 2: Estimate the loss ("loss")

q Loss function $L(\mathbf{y}^{\wedge}, \mathbf{y})$

\vec{y} error/loss occurs when the prediction \mathbf{y}^{\wedge} is different from the correct output \mathbf{y}

\vec{y} function maps the two vectors to a scalar

Which two vectors?

\vec{y} The goal of training is to minimize loss: less is better

9a.7

7

Loss

q Examples of commonly used functions

\vec{y} For binary classification problems

- Loss of binary articulation
- Binary cross-entropy
- ...

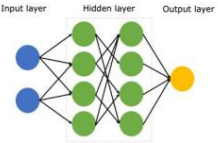
\vec{y} For multi-class classification problems

- Multi-class articulation loss
- Categorical cross-entropy
- ...

9a.8

8

Binary classification 1/2



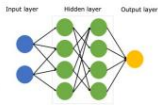
q Loss of binary articulation ("hinge")

- also known as marginal loss
- binary classification: prediction y^{\wedge} is a scalar and **correct outputs are -1 or +1**
- correct classification: y^{\wedge} and y have the same sign
- zero loss when prediction y^{\wedge} and output y have the same sign and $|y^{\wedge}| \geq 1$
- $L(y^{\wedge}, y) = \max(0, 1 - yy^{\wedge})$

9a.9

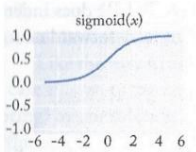
9

Binary classification 2/2



q Binary cross-entropy

- also known as logistic loss
- correct outputs are 0 or 1
- prediction y^{\wedge} in $[0, 1]$ after y^{\wedge} is passed through the sigmoid function
- correct classification: $y^{\wedge} < 0.5$ and $y = 0$, or $y^{\wedge} \geq 0.5$ and $y = 1$



• $L(y^{\wedge}, y) = -\log y^{\wedge} = -\log(1 - y^{\wedge})$

- if $y = 1$
- if $y = 0$

9a.10

10

Intermezzo: softmax function (recap)


q $\text{softmax}(x[i]) = \frac{\exp(x[i])}{\sum_j \exp(x[j])}$

- after application to the components of the vector x :
- each component with value in $[0, 1]$
- the sum of all components results in 1
- component values interpretable as a probability distribution
- normalizing function

9a.11

11

Multi-class classification 1/2



q Loss of multi-class hinge

- correct output will be a one-hot y vector : **components a 0 except one**
- correct class: $t = \text{argmax}_i y[i]$
- network output is a vector $y^{\wedge} = y^{\wedge}[1], \dots, y^{\wedge}[m]$
- predicted class: $k = \text{argmax}_i y^{\wedge}[i]$
- $L(y^{\wedge}, y) = \max (0, 1 - (y^{\wedge}[t] - y^{\wedge}[k]))$ com $t \neq k$

9a.12

12

Multi-class classification 2/2



Categorical cross-entropy

• also known as negative log likelihood
• correct output will be a vector $y = y[1], \dots, y[m]$ representing the multinomial distribution over the labels

• item may belong to different classes with some degree of certainty

• network output is a vector $y^{\wedge} = y^{\wedge}[1], \dots, y^{\wedge}[m]$ after being transformed by the softmax function (each value in $[0,1]$ and all added together results in 1)

$$L(y^{\wedge}, y) = - \sum_i y[i] \log(y^{\wedge}[i])$$

9a.13

9a.13

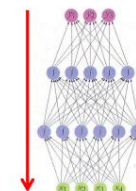
13

Step 3: Backward pass ("backward pass")

Back-propagation

• compute the gradient of each parameter relative to the loss that was estimated

• move the parameters in the opposite direction to the sign of the gradient



9a.14

9a.14

14

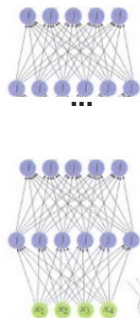
Interlude: networks calculate composite functions (recap)

Multi-layer neural network

• is a composite function that involves the successive
• multiplication of vectors by matrices,
• passing vectors through non-linear functions,
• and sum with bias vectors:

$$y^{\wedge} = g(\dots g(g(xW_1 + b_1)W_2 + b_2)W_3 + b_3) \dots) W_n + b_n$$
$$y^{\wedge} = f(x)$$

• which, during training, is composed with the function of loss: $L(f(x), c(x))$
where $c(x) = y$ the correct classification for each x



9a.15

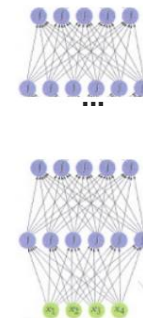
9a.15

15

Optimization based on 1/4 gradient

Stochastic gradient descent ("stochastic gradient descent + SGD")

1. compute the gradient g of each parameter with respect to the loss



9a.16

9a.16

16

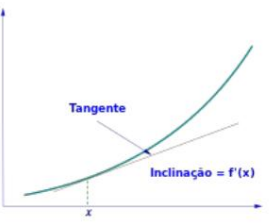
Interlude Gradient

q Derivative of a function

• the derivative function f' of a function f measures the sensitivity of your output to change your entry

• the derivative function of a unary function in your entry is the slope of the tangent line to the graph at that point

• the sign of the derivative function, negative or positive, indicates whether the object function is decreasing or increasing.



9a.17

17

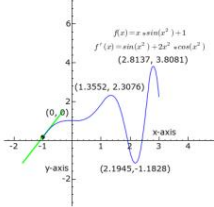
Gradient-based optimization 2/4

q Objective: minimize the loss function

• the loss function measures the error, i.e. the distance between the prediction and the correct output

• the goal of training is to find the weights that minimize the loss, i.e. that make the distance between the prediction and the correct output zero

• that is, the objective is to adjust the weights of such so that the output of the (loss) function calculated by the network is zero



9a.18

18

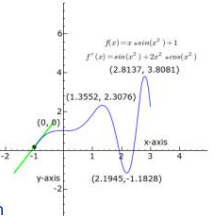
Gradient-based optimization 3/4

q Objective: "counteract the sign of the gradient"

• that is, the objective is to adjust the network weights in the opposite direction to the (de)growth of the function of loss

• that is, in terms of gradient, the goal is change each weight in the opposite direction to the sign of its gradient

• that is, if the gradient is negative the (loss) function is decreasing, an increase in weight will contribute to increasing the prediction value and thus making the loss continue to decrease and therefore approaching zero (and vice versa)

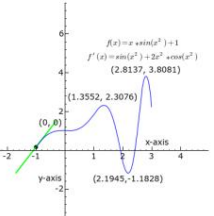


9a.19

19

q N.B.

q This graph was used in the previous slides only to help form intuitions about the relevance of the derivative to the optimization process; does not represent "the" loss function curve, nor any curve especially relevant to deep learning



9a.20

20

Optimization based on 4/4 gradient

q Stochastic gradient descent

("stochastic gradient descent = SGD")

• 1. compute the gradient g of each parameter with respect to the loss

• 2. move each parameter p in the opposite direction to its gradient:
replace p with what results from its adjustment by the gradient
after this is weighted by a learning rate t :

$$p \leftarrow p - t \cdot g$$

q Variants [advanced topics]

• Gradients are accumulated: SGD+Momentum, Nesterov Momentum

• Adaptive learning rate: AdaGrad, AdaDelta, RMSProp, Adam

21

- Conclusion -

q Index

- Training algorithm
- Forward pass computation
- Loss functions
- Backward pass computation

What hyperparameters should I configure a model with (to train it)?

22