

Deep learning

Zoltan Miklos

University of Rennes 1

2022

Plan

1 What is deep learning ?

What is deep learning ?

- layered representation learning
- recently very successfully applied for classification, regression tasks, as well as data generation
- a learning method that can approximate complicated decision functions

Chair ?



[http://www.dekoratifblog.com/
beautify-your-chair-room-with-a-variety-of-chairs/](http://www.dekoratifblog.com/beautify-your-chair-room-with-a-variety-of-chairs/)

Feature engineering

- A critical step in the machine learning workflow : which features of the data shall we use ?
- Deep learning completely removes this step
- However we lose on interpretability : a network can predict with high accuracy, for example, whether an image depicts a chair, but we have no information how he could achieve this. The large number of learned weights do not give any insights. (while in case of a decision tree classification, humans can interpret the learned structure)

Plan

2 Neural networks : basic concepts

Neural networks : roadmap for the topics to come

- Perceptron : basic neural network
- Multi-layer perceptron, feed-forward networks
- Training neural networks, backpropagation
- How to influence the training process ? (early stopping, regularization, dropout, etc.)
- Convolutional neural networks, basic image processing
- Recurrent neural networks, LSTM
- Word-embeddings, processing text or other sparse data
- Machine learning strategy
- Other topics

Perceptron

- Basic building block of artificial neural networks
- Forward propagation : the perceptron takes a weighted sum of the inputs and produces an output using a non-linear activation function

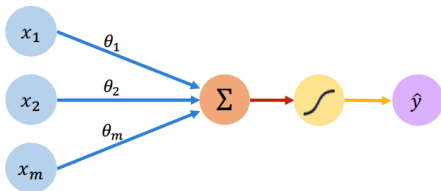
•

$$\hat{y} = g(\sum_{i=1}^n x_i \theta_i)$$

- in vector notation $\hat{y} = g(\Theta^\top X)$

Perceptron : forward propagation

The Perceptron: Forward Propagation



Linear combination of inputs

Output

$$\hat{y} = g \left(\sum_{i=1}^m x_i \theta_i \right)$$

Non-linear activation function

Inputs Weights Sum Non-Linearity Output

Bias term

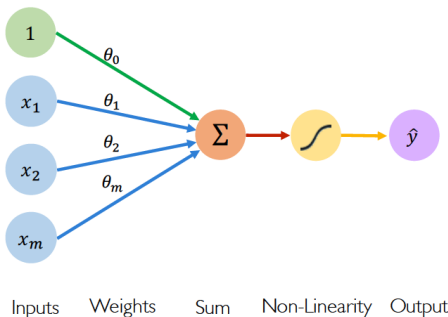
- Often we would like to add a bias term as well

$$\hat{y} = g(\theta_0 + \sum_{i=1}^n x_i \theta_i)$$

- we can consider that we have a constant 1 input (the 0th input)
- in vector notation $\hat{y} = g(\sum_{i=1}^n \Theta^\top X)$,
where $X = (1, x_1, \dots, x_n)^\top$ and $\Theta = (\theta_0, \theta_1, \dots, \theta_n)^\top$

Perceptron : the bias term

The Perceptron: Forward Propagation



Linear combination of inputs

$$\hat{y} = g \left(\theta_0 + \sum_{i=1}^m x_i \theta_i \right)$$

Output

Non-linear activation function

Bias

Activation functions

Non-linear functions :

- Sigmoid

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

- reLU (Rectified Linear Unit)

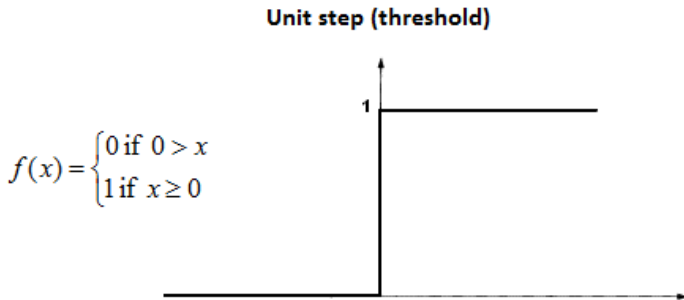
$$relu(x) = \max(0, x)$$

- tanh

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

Activation functions : Unit-step threshold

Not used in practice (as it is not differentiable), but helps to demonstrate the concepts in small toy examples



https:

[//www.saedsayad.com/artificial_neural_network.htm](http://www.saedsayad.com/artificial_neural_network.htm)

Logical gates

Exercise : Create a (single layer) perceptron that works as a logical gate and realizes the

- NOT gate
- AND gate
- OR gate

XOR : is it possible to realise with a single layer perceptron ?

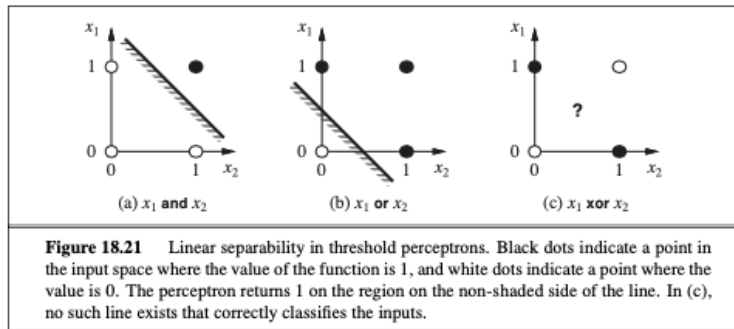
x_1	x_2	$x_1 \text{ XOR } x_2$
0	0	0
0	1	1
1	0	1
1	1	0

Table – XOR

Decision surface and the XOR gate

Decision boundary : a hypersurface that partitions the underlying vector space into two sets, for each class (0/1).

One can identify a separating hyperplane in the case of AND and OR functions, but there is no separating hyperplane for XOR :



XOR : a historical side-note

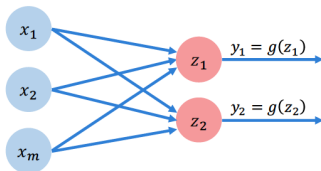
- Frank Rosenblatt (1957)
- Marvin Minsky and Seymour Papert's book on Perceptrons (1969) showed that they (single-layer version) cannot compute XOR
- While Minsky and Papert understood that multi-layer version can do it, but incorrect citations contributed to a decline of neural network research, until the 1980's

Deep neural network

- Use multiple hidden layers
- use non-linear activation functions
- IMPORTANT : one can consider the deep neural network as an ensemble of perceptions (where we apply logistic regression several times). However the high number of neurons (i.e. weights) changes the behaviour of this learning system qualitatively. This phenomenon is poorly understood theoretically.

Multi Output Perceptron

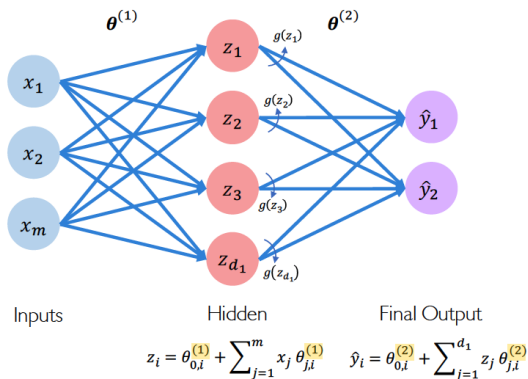
Multi Output Perceptron



$$z_i = \theta_{0,i} + \sum_{j=1}^m x_j \theta_{j,i}$$

<http://introductiontodeeplearning.com> $\theta_{j,i}$ the weight of the edge from input x_j to node z_i

Neural Network with a hidden layer



<http://introductiontodeeplearning.com>

Forward calculations in the neural network

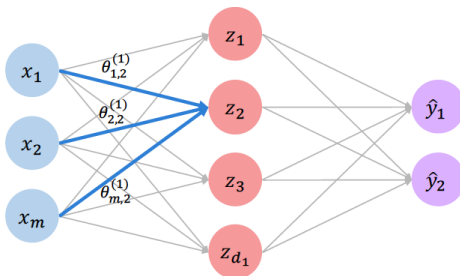
- Input layer, Hidden layer, Output layer



$$z_i = \theta_{0,i}^{(1)} + \sum_{j=1}^n x_j \theta_{j,i}^{(1)}$$

$$\hat{y}_i = \theta_{0,i}^{(2)} + \sum_{j=1}^{d_1} z_j \theta_{j,i}^{(2)}$$

A visual help for understanding the formulas



$$\begin{aligned} z_2 &= \theta_{0,2}^{(1)} + \sum_{j=1}^m x_j \theta_{j,2}^{(1)} \\ &= \theta_{0,2}^{(1)} + x_1 \theta_{1,2}^{(1)} + x_2 \theta_{2,2}^{(1)} + x_m \theta_{m,2}^{(1)} \end{aligned}$$

<http://introductiontodeeplearning.com>

Plan

3 Training neural networks

Training neural networks

We would like to use the neural networks for classification/prediction

- Supervised learning : we use a (large) set of labelled objects
- Training data : $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$
- Training data : set of data items \mathbf{x}_i , for which the correct label is y_i (ideally hundreds of millions of data items)
- We define the topology of the network
- The model we would like to learn is the set of weights $\Theta_{i,j}^{(k)}$ (k -th layer, weight between the input i and output j)

Loss function

- Characterize the differences between the predicted and real output
- Cross entropy loss (for outputs between 0 and 1) (discrete case) $-\sum_i y_i \log_2(h_w(x_i)) + (1 - y_i) \log_2(h_w(x_i))$
- Mean-squared error loss
- (Besides the widely used loss functions you can also a specific one, more adapted to your particular problem.)

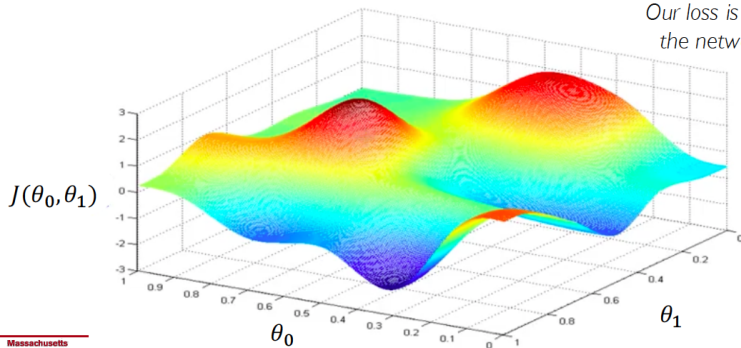
Training neural networks

- We would like to find the weights that result the lowest loss
- We can consider this as an optimization : minimize the loss function (that is a function of all weights and biases)
- Optimization : use the partial derivatives to find the optimum
- Analytical solution : calculate all the partial derivatives, set them equal to 0 and solve the equations (however, this is only possible for very small networks, already obtaining the derivatives is too costly)

Loss Optimization

$$\theta^* = \underset{\theta}{\operatorname{argmin}} J(\theta)$$

Remember:
Our loss is a function of
the network weights!

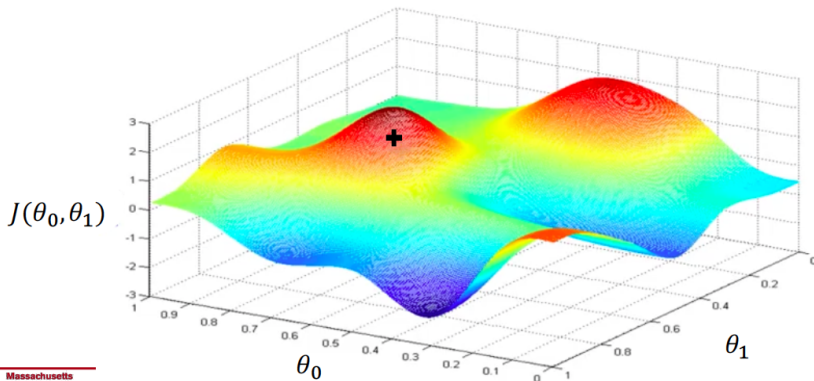


Gradient descent

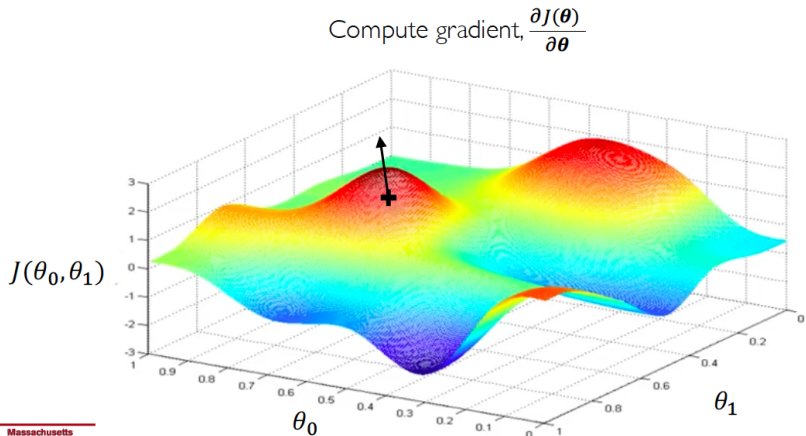
- Initialize the weights randomly
- Loop until convergence
 - Compute the components of the gradient $\nabla J = (\frac{\partial J}{\partial \theta_1}, \dots, \frac{\partial J}{\partial \theta_n})$
 - Update the weights $\theta^{(i+1)} \leftarrow \theta^{(i)} - \gamma \nabla J$ where $\gamma \in [0, 1]$ is the learning rate (parameter).
 - Attention : θ here is a vector (of n dimension) that contains all network weights (including the bias terms)
- Return the weights

Gradient Descent

Randomly pick an initial (θ_0, θ_1)

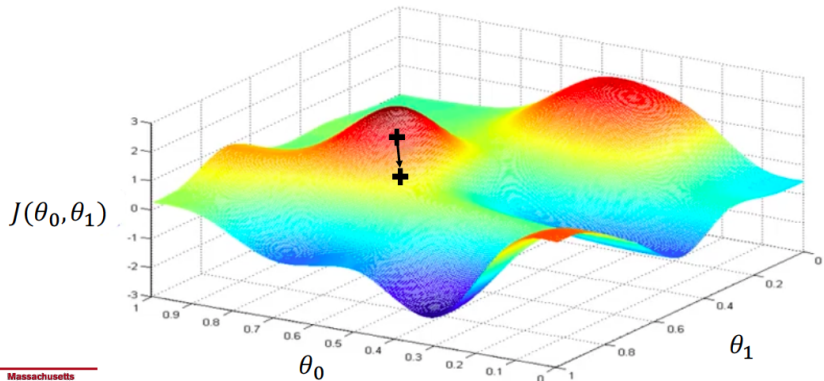


Gradient Descent



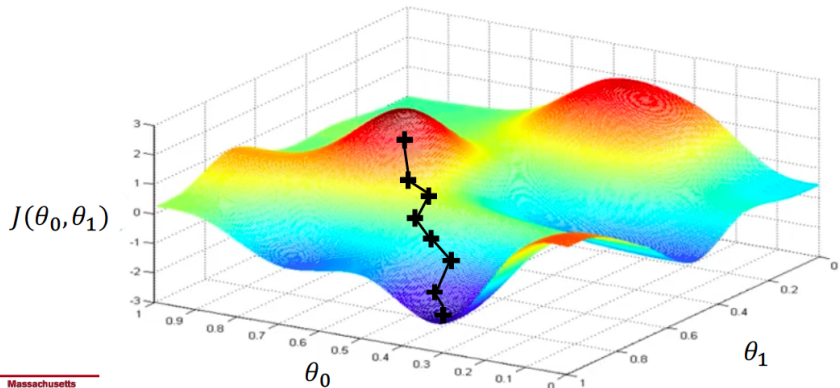
Gradient Descent

Take small step in opposite direction of gradient



Gradient Descent

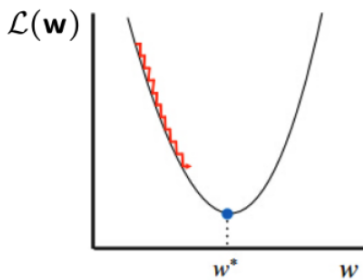
Repeat until convergence



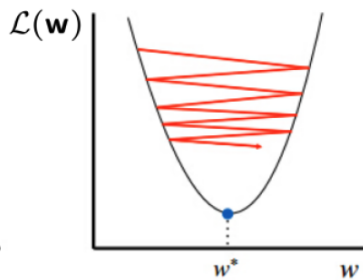
Learning rate

- Setting the learning rate is not easy
 - Small learning rate can lead to slow convergence and the algorithm could get stuck in a local minimum
 - Large learning rate : the algorithm might even diverge and we could miss the minimum
- It is important to monitor the learning process (e.g. by measuring accuracy of the trained models)

Learning rate



Too small: converge
very slowly



Too big: overshoot and
even diverge

Backpropagation

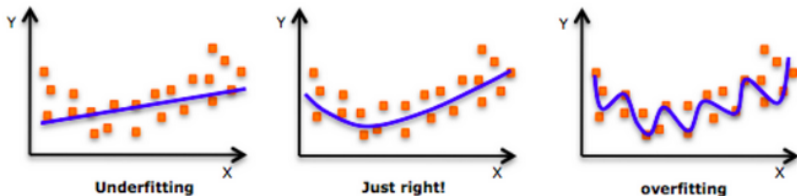
Learning algorithm of the neural networks

- Rumelhart, Hinton, Williams in 1986
- Use forward propagation to compute output, evaluate the loss function and then
- propagate back, that is update the weights (from the output layer back towards the input layer)
- updates correspond to the chain rule (of derivation)
- In this way the network computes implicitly the derivatives of the cost function
- epoch : complete set of computations where you see all your training data

Stochastic gradient descent (SGD)

- Gradient descent can be slow if the training set is very large
- Gradient descent : $\theta_i \leftarrow \theta_i - \eta \frac{\partial J}{\partial \theta_i}$
- Stochastic gradient descent :
 - Randomly shuffle the training data
 - The gradient is estimated using a small set of samples (instead of the entire training set)
- Other optimization techniques (available in tensorflow or elsewhere)
 - Momentum
 - Adagrad
 - Adam
 - RMSprop

Overfitting



<https://www.analyticsvidhya.com/blog/2018/04/fundamentals-deep-learning-regularization-techniques/>

Regularization

- Techniques to constraint the optimization problem (to avoid too complex models)
- Can help to avoid overfitting
- $L1$ (Lasso) :

$$Cost = Loss + \frac{\lambda}{2m} \sum |\theta|$$

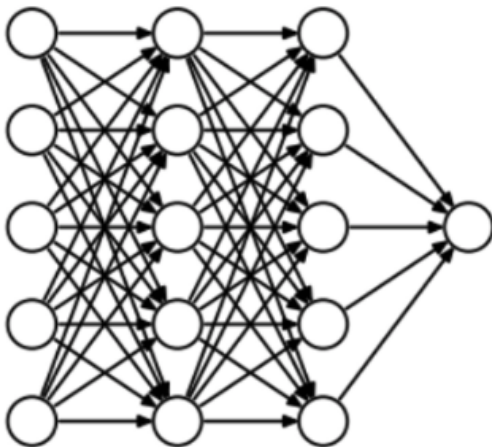
- $L2$ (Ridge) :

$$Cost = Loss + \frac{\lambda}{2m} \sum |\theta|^2$$

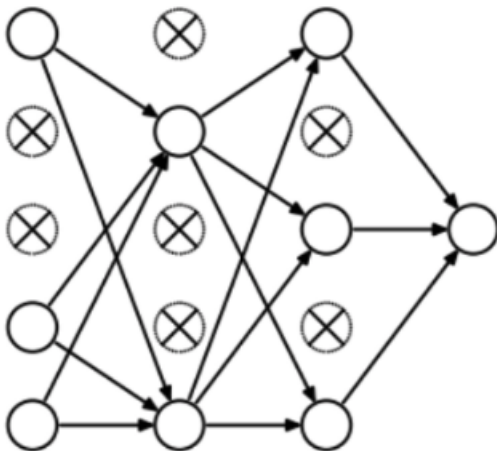
Dropout

- Set some activations to 0 (randomly, with some probability that is a parameter)
- Avoids that a network relies on some particular nodes

Dropout : example (complete network)



Dropout : example (connections after the dropout operation)

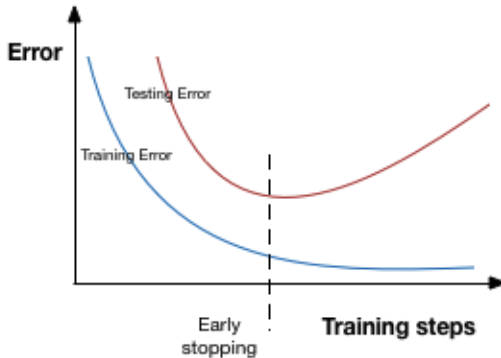


Early stopping

Early stopping is a technique to prevent overfitting

- monitor the error of the trained models
- stop the training of the network early (before convergence)
- In this way, we will not reach the maximum possible accuracy on the training set, but this could lead to better performance on the test set (if we manage to avoid overfitting)

Early stopping



<https://www.analyticsvidhya.com/blog/2018/04/>