

An Introduction to Deep Learning

Clint P. George

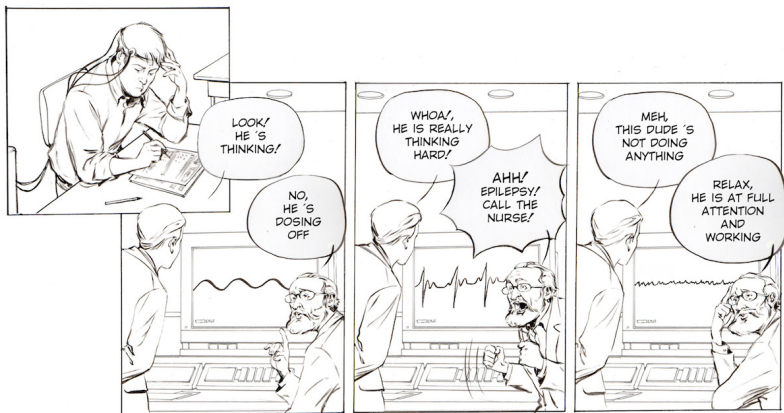
University of Florida Informatics Institute

March 14, 2017

Outline

- ① Introduction to feed-forward networks
- ② Relation to logistic regression
- ③ Notes on implementation
- ④ Illustration using synthetic and real data

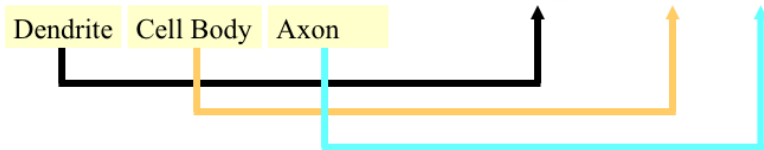
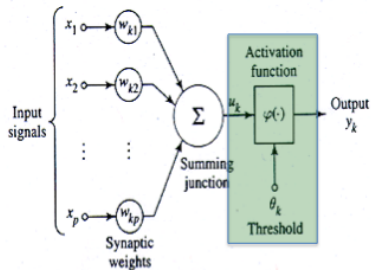
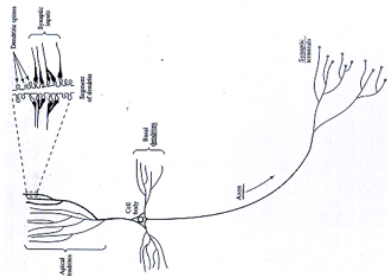
Motivation



The idea: Human intelligence may be due to a learning algorithm.
We aim to build algorithms that mimic the brain.¹

¹image: <https://backyardbrains.com/experiments/EEG>

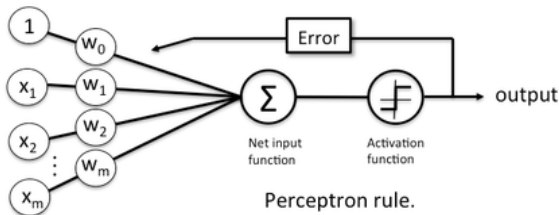
Imitate neurons in the brain: Artificial Neurons



Artificial Neuron (AN): input, weights, and output

Activation functions

The Perceptron (Rosenblatt et al. 1957 & 1962) computes a step function as an activation function.



$$\text{step}(z) = \begin{cases} 1 & z \geq t \\ 0 & z < t \end{cases}, \text{ where } t \text{ is a threshold}$$

As each input is applied to the perceptron its output is compared to the target. To keep the output closer to the target the **learning rule** adjusts the network parameters.

What can a single AN compute?

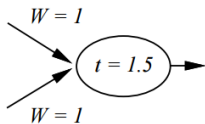
The Perceptron output is given by $y = \text{step}(b + \sum_{j=1}^p w_j x_j)$.

Perceptron can divide the input space into two regions.

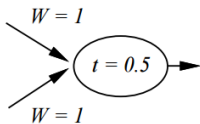
The decision boundary is given by: $b + \sum_{j=1}^p w_j x_j = 0$.

Perceptron can learn to classify any linearly separable set of inputs—convergence theorem (Rosenblatt 1962)

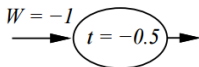
Boolean functions



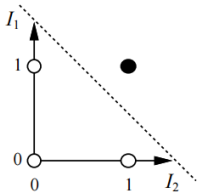
AND



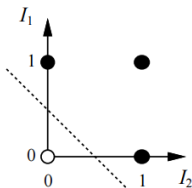
OR



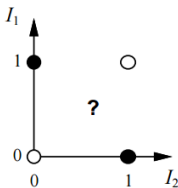
NOT



(a) I_1 and I_2



(b) I_1 or I_2



(c) I_1 xor I_2

Examples² of linearly separable and non separable problems

²Veloso, 2001

Learning the XOR function

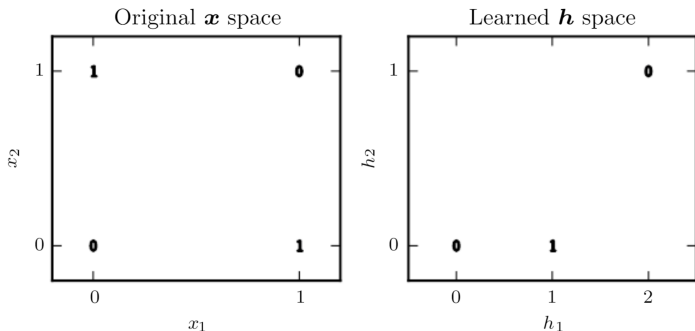
Consider this as a regression problem and use the MSE loss function:

$$\text{MSE}(\theta) = \frac{1}{4} \sum_{\mathbf{x} \in \mathcal{X}} (y - f_{\theta}(\mathbf{x}))^2$$

where $\theta = (\mathbf{w}, b)$ and $f_{\theta}(\mathbf{x}) = b + \sum_{j=1}^p w_j x_j$ —a linear model.

Using the normal equations we can minimize $\text{MSE}(\theta)$ w.r.t. \mathbf{w} and b in closed form. It gives $\mathbf{w} = 0$ and $b = .5$ —This gives the model output .5 everywhere.

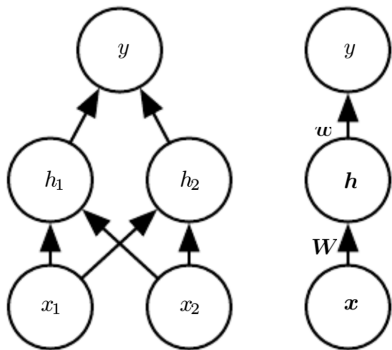
Learning the XOR function



- When $x_1 = 0$, the model's output must increase as x_2 increases
- When $x_1 = 1$, the model's output must decrease as x_2 increases
- A linear model applies a fixed coefficient w_2 to x_2 . It cannot use the value of x_1 to change the coefficient w_2 on x_2 and cannot solve this problem.

Intuition behind multilayer neural network

One way to solve the XOR problem is to transform the input by introducing a feed-forward network:



The complete model will then be, in a function form:

$$f(\mathbf{x}; \mathbf{W}, \mathbf{w}, b, c) = f^{(2)} \left(f^{(1)}(\mathbf{x}; \mathbf{W}, c); \mathbf{w}, b \right)$$

What function should $f^{(1)}$ be?

We consider $f^{(2)}$ as a linear function.

We can write the hidden layer output as

$$\mathbf{h} = f^{(1)}(\mathbf{x}; \mathbf{W}, \mathbf{c}) = g(\mathbf{W}^\top \mathbf{x} + \mathbf{c})$$

What function should $f^{(1)}$ be?

We consider $f^{(2)}$ as a linear function.

We can write the hidden layer output as

$$\mathbf{h} = f^{(1)}(\mathbf{x}; \mathbf{W}, \mathbf{c}) = g(\mathbf{W}^T \mathbf{x} + \mathbf{c})$$

If $f^{(1)}$ is also linear, then the network as a whole would remain a linear function of its input.

What function should $f^{(1)}$ be?

We consider $f^{(2)}$ as a linear function.

We can write the hidden layer output as

$$\mathbf{h} = f^{(1)}(\mathbf{x}; \mathbf{W}, \mathbf{c}) = g(\mathbf{W}^T \mathbf{x} + \mathbf{c})$$

So we must use a nonlinear activation function for g . A popular nonlinear activation function g is rectified linear unit (ReLU)

$$g(z) = \max\{0, z\}$$

A feed-forward network solution to XOR

We wish the network to perform well on the four cases

$$\mathbf{X} = \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{bmatrix}, \mathbf{y} = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

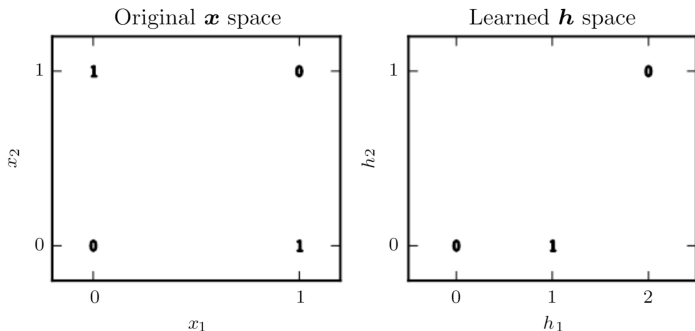
We can specify a two-layer network solution to XOR as

$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \mathbf{c} = \begin{bmatrix} 0 \\ -1 \end{bmatrix}, \mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix}, b = 0$$

The complete network is given by

$$f(\mathbf{x}; \mathbf{W}, \mathbf{w}, b, \mathbf{c}) = \mathbf{w}^T \max\{0, \mathbf{W}^T \mathbf{x} + \mathbf{c}\} + b$$

A feed-forward network solution to XOR



In the proposed network, the nonlinear hidden layer has mapped both $\mathbf{x} = [1, 0]$ and $\mathbf{x} = [0, 1]$ to a single point in feature space, $\mathbf{h} = [1, 0]$.

A linear model can now describe the function as increasing in h_1 and decreasing in h_2 .

Activation functions

Modern ANs use a variety of activation functions that are smoother than the step function.

Linear - no input squashing

$$y = x$$

Logistic sigmoid - squash input into $[0, 1]$

$$y = \text{sigmoid}(x) = \frac{1}{1 + \exp(-x)}$$

Hyperbolic tangent - squash input into $[-1, 1]$

$$y = \tanh(x) = \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)}$$

Rectified linear unit

$$y = \max\{0, x\}$$

Relation to logistic regression

Let $p(y = 1 \mid X = \mathbf{x}) = p(\mathbf{x}; \mathbf{w})$ be the conditional probability that a particular sample belongs to class 1 given its predictors \mathbf{x} . We write the logistic regression model as

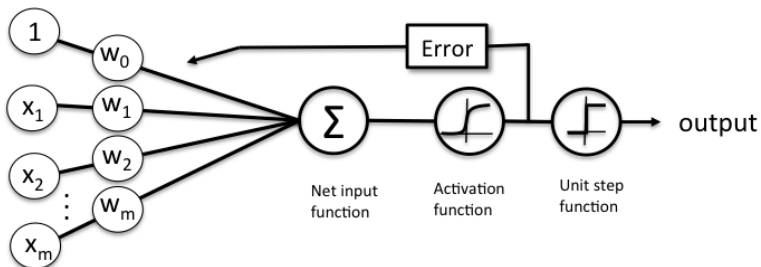
$$\text{logit}(p(\mathbf{x}; \mathbf{w})) = w_0 + \mathbf{w}^\top \mathbf{x}$$

Solving for $p(\mathbf{x}; \mathbf{w})$ gives

$$p(\mathbf{x}; \mathbf{w}) = \text{sigmoid}(w_0 + \mathbf{w}^\top \mathbf{x})$$

To minimize misclassification rate, one should predict $y = 1$ when $p \geq .5$, and vice versa.—i.e. guess 1 whenever $w_0 + \mathbf{w}^\top \mathbf{x} \geq 0$ and 0 otherwise. So logistic regression gives a linear classifier.

Relation to logistic regression



Schematic of a logistic regression classifier.

3

Learning parameters w and b : maximum likelihood estimation

- Perceptron algorithm: online and error-driven.
- Logistic regression: batch algorithms—e.g. gradient descent, limited-memory BFGS, or online algorithms—stochastic gradient descent.

³http://rasbt.github.io/mlxtend/user_guide/classifier/LogisticRegression/

Gradient-based learning

The nonlinearity of a neural network causes most interesting loss functions to become non-convex

- Optimization procedure – using iterative, gradient-based optimizers that drive the cost function to a very low value
- Cost function, $C(\theta)$ – they are more or less the same as those for other parametric models, such as linear models

Gradient descent (GD) algorithm⁴

Let $C(\theta) = \int \mathcal{L}(f_\theta, z)P(z)dz$, where in supervised learning $z = (x, y)$ and $f_\theta(x)$ is the predictive function for y given θ .

Gradient descent: find a θ that minimizes the cost

- By solving $\frac{\partial C(\theta)}{\partial \theta} = 0$ we can find the minima, maxima, and saddle points.
- In general we cannot find the solutions of this equation, hence we seek numerical optimization methods
- *local descent*: iteratively modify θ so as to decrease $C(\theta)$, until we reach a local minima

$$\theta^{(t+1)} = \theta^{(t)} - \epsilon \frac{\partial C(\theta^{(t)})}{\partial \theta^{(t)}},$$

ϵ is the learning rate

⁴<http://www.iro.umontreal.ca/~pift6266/H10/notes/gradient.html>

Stochastic gradient descent (SGD) algorithm

Let $C(\theta) = \int \mathcal{L}(f_\theta, z)P(z)dz$, where in supervised learning $z = (x, y)$ and $f_\theta(x)$ is the predictive function for y given θ .

Stochastic gradient descent: find a θ that minimizes the cost

- Similar to GD, but exploits the fact that $C(\theta)$ is an average, generally over i.i.d. samples
- Make updates much often, the most common case

$$\theta^{(t+1)} = \theta^{(t)} - \epsilon \frac{\partial \mathcal{L}(\theta^{(t)}, z)}{\partial \theta^{(t)}},$$

z is an the next sample from the training set.—It can be implemented online

- In SGD, the update direction is a random variable whose expectation is the true gradient of interest.

Learning Conditional Distributions with Maximum Likelihood

Popular choice: the neural network defines a distribution $p(\mathbf{y} | \mathbf{x}, \theta)$ and uses the principle of maximum likelihood

Cost function: the negative log-likelihood—the cross-entropy between the training data and the model distribution

$$C(\theta) = -\mathbb{E}_{\mathbf{x}, \mathbf{y}} \log p_{\text{model}}(\mathbf{y} | \mathbf{x}, \theta)$$

An example: if $p_{\text{model}}(\mathbf{y} | \mathbf{x}, \theta) = \mathcal{N}(f_{\theta}(\mathbf{x}), \mathbf{I})$, then we have

$$C(\theta) = \frac{1}{2} \mathbb{E}_{\mathbf{x}, \mathbf{y}} \|\mathbf{y} - f_{\theta}(\mathbf{x})\|^2 + \text{const.}$$

Output Units

The output layer provides additional transformation from the hidden features to complete the network's intended task.

Linear units for Gaussian output distributions: the output units based on an affine transformation with no nonlinearity

- Given hidden features \mathbf{h} , we define outputs $\hat{\mathbf{y}} = \mathbf{w}^T \mathbf{h} + b$
- Used to produce the mean of $p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\hat{\mathbf{y}}, \mathbf{I})$
- MLE is equivalent to minimizing the MSE
- The linear units do not saturate, hence ideal for gradient-based optimization algorithms

Output Units

Sigmoid units for Bernoulli output distributions: e.g. for binary classification

Given hidden features \mathbf{h} , we define output units

$$\hat{y} = \text{sigmoid}(\mathbf{z}) = \text{sigmoid}(\mathbf{w}^T \mathbf{h} + b)$$

The sigmoid can be motivated by constructing an unnormalized probability distribution that doesn't sum to 1. We then normalize to yield a Bernoulli distribution

$$p(y | \mathbf{x}) = \text{sigmoid}((2y - 1)\mathbf{z})$$

The cost function for MLE is defined by the negative log-likelihood

Hidden Units

ReLU - the default choice, which is similar to linear

- $g(z) = \max\{0, z\}$, not differentiable at $z = 0$; but one can safely disregard this problem: *"... neural network training algorithms do not usually arrive at a local minimum of the cost function, but instead merely reduce its value significantly"*
- Drawback: they cannot learn via gradient-based methods on examples for which their activation is zero

Hidden Units

Sigmoid and Hyperbolic Tangent

- Sigmoidal units saturate across most of their domain—they saturate to a high value when z is very positive
- It can make gradient-based learning very difficult. Hyperbolic Tangent typically performs better than the logistic sigmoid and can be used

An overview of back-propagation algorithm

Computing an analytical expression for the gradient is straightforward, but numerical evaluation of such an expression can be expensive.—**backprop** gives an inexpensive procedure to evaluate this gradient.

backprop is a generic algorithm which can be applied to any problem where we need to evaluate $\nabla_{\theta}C(\theta)$

backprop uses the chain rule of calculus: Suppose $\mathbf{x} \in \mathbb{R}^m$, $\mathbf{y} \in \mathbb{R}^n$, $g : \mathbb{R}^m \rightarrow \mathbb{R}^n$, $f : \mathbb{R}^n \rightarrow \mathbb{R}$. If $\mathbf{y} = g(\mathbf{x})$, $z = f(\mathbf{y})$, then

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \cdot \frac{\partial y_j}{\partial x_i}.$$

We can apply this to vectors and tensors in the multilayer network

Notes on architectural considerations

In practice, the overall structure of the network is important: how many units it should have and how these units should be connected to each other.

Other examples

- Convolutional Neural Networks — imitates human memory
- Auto Encoders — unsupervised learning and dimensionality reduction

Hands on experiments

Datasets

- A synthetic spiral dataset with multiple classes
- The MNIST data set with 0-9 handwritten characters

Algorithm: A basic back-propagation algorithm implementation with one hidden layer

Link to R scripts: <http://bit.ly/deep-learning-stats>

Links to serious implementations

Deep learning frameworks, which uses CPU and GPU

- TensorFlow with R
<https://rstudio.github.io/tensorflow/index.html>
- Theano with python
<http://deeplearning.net/software/theano>

Available R packages

- neuralnet
- deepnet
- h2o

Questions?