Multi-Layer Perceptrons

Johanni Brea

Introduction to Machine Learning



Feature Engineering is great, but couldn't it be automatized?

With smart features basically any problem can be solved by a linear method.

How should we find the smart features?

Idea: Let us take a more flexible function family and find "features" and "regression coefficients" at the same time with gradient descent.

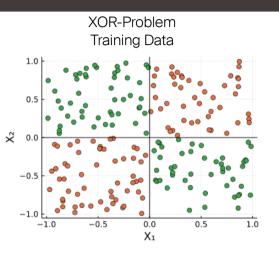


Table of Contents

- 1. Solving the XOR Problem Without Feature Engineering
- 2. Artificial Neurons
- 3. Multilayer Perceptrons
- 4. Regression with Multilayer Perceptrons
- 5. Classification with Multilayer Perceptrons



Recap: Vector-Features

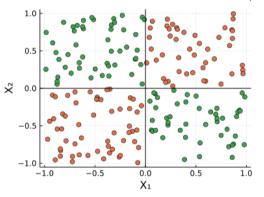


Logistic Regression fails:
There is no linear decision boundary.



Recap: Vector-Features

Project data to a higher dimensional space by computing the scalar products between feature vectors w_1, \ldots, w_q and input vectors x_i and thresholding.



For example $h_{21} = \max(0, w_1^T x_2)$.

Logistic Regression on the features works.

Solving the XOR Problem without Feature Engineering

Logistic Regression:
$$P(Y = 1|x, \beta) = \sigma(\beta_0 + \beta_1 x_1 + \cdots + \beta_p x_p)$$

Logistic Regression on features:

$$P(Y = 1|x, \beta) = \sigma(\beta_0 + \beta_1 \underbrace{g(w_1^T x)}_{h_1} + \cdots + \beta_q \underbrace{g(w_q^T x)}_{h_q})$$

with hand-picked feature vectors w_1, \ldots, w_q and activation function $g(x) = \text{relu}(x) = \max(0, x)$.

Idea

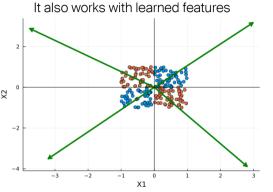
Why don't we learn the features with gradient descent?

$$P(Y = 1 | x, \beta, w_1, \dots, w_q) = \sigma(\beta_0 + \beta_1 g(w_1^T x) + \dots + \beta_q g(w_q^T x))$$

$$\Rightarrow \hat{\beta}, \hat{w} = \arg \min_{\beta, w_1, \dots, w_q} - \sum_{i=1}^n \log P(y_i | x_i, \beta, w_1, \dots, w_q)$$



Solving the XOR Problem without Feature Engineering



- ▶ We just fitted our first neural network ^②.
- ► The loss function has local minima; gradient descent does not find for all initial guesses a good solution.
- With more than 4 feature vectors, gradient descent finds good solutions for most initial guesses.

Table of Contents

- 1. Solving the XOR Problem Without Feature Engineering
- 2. Artificial Neurons
- 3. Multilayer Perceptrons
- 4. Regression with Multilayer Perceptrons
- 5. Classification with Multilayer Perceptrons



Artificial Neurons

Artificial neurons take a *d*-dimensional input $x = (x_1, \dots, x_d)^T$ and output a scalar

$$a = g(w_0 + w_1x_1 + w_2x_2 + \cdots + w_dx_d)$$

with parameters (or weights) w_0 , w_1 , ..., w_d and activation function g. w_0 is also called **bias** (instead of intercept).



Popular Activation Functions

rectified linear unit
$$\operatorname{relu}(x) = \max(0, x) = \begin{cases} x & x \ge 0, \\ 0 & x < 0 \end{cases}$$

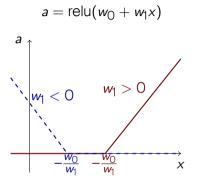
sigmoid
$$\sigma(x) = \frac{1}{1+e^{-x}}$$

tangent hyperbolic
$$tanh(x)$$

softplus softplus(
$$x$$
) = log(exp(x) + 1)

heaviside
$$H(x) = \begin{cases} 1 & x \ge 0, \\ 0 & x < 0 \end{cases}$$

Artificial relu-Neurons



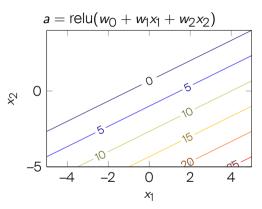




Table of Contents

- 1. Solving the XOR Problem Without Feature Engineering
- 2. Artificial Neurons
- 3. Multilayer Perceptrons
- 4. Regression with Multilayer Perceptrons
- 5. Classification with Multilayer Perceptrons



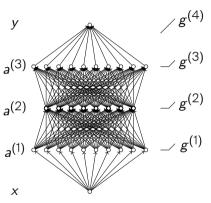
Multilayer Perceptrons

Multilayer Perceptrons (MLP) consist of multiple neurons organized in layers 1, 2, ..., L. Each layer has $d^{(I)}$ neurons and activation $g^{(I)}$.

output $a_k^{(I)}$ of k-th neuron in I-th layer

$$a_k^{(I)} = g^{(I)} \Big(w_{k0}^{(I)} + w_{k1}^{(I)} a_1^{(I-1)} + \dots + w_{kd^{(I-1)}}^{(I)} a_{d^{(I-1)}}^{(I-1)} \Big)$$

input layer $a_k^{(0)} = x_k$.



one input neuron x
one linear output neuron y
3 hidden layers with relu-neurons

Multilayer Perceptrons: Matrix Notation

$$a_k^{(I)} = g^{(I)} \left(w_{k0}^{(I)} + w_{k1}^{(I)} a_1^{(I-1)} + \dots + w_{kd^{(I-1)}}^{(I)} a_{d^{(I-1)}}^{(I-1)} \right)$$
matrix notation $a^{(I)} = g^{(I)} \left(b^{(I)} + w^{(I)} a^{(I-1)} \right)$
with $b_k^{(I)} = w_{k0}^{(I)}$.

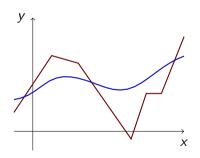
For example the network on the previous slide can be written as

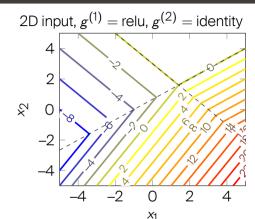
$$y = b^{(4)} + w^{(4)} \operatorname{relu}\left(b^{(3)} + w^{(3)} \operatorname{relu}\left(b^{(2)} + w^{(2)} \operatorname{relu}\left(b^{(1)} + w^{(1)}x\right)\right)\right)$$



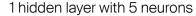
Multilayer Perceptrons

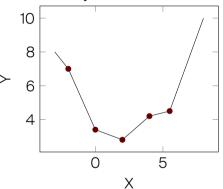
$$g^{(1)} = \text{relu}, g^{(1)} = \text{tanh}, g^{(2)} = \text{identity}$$





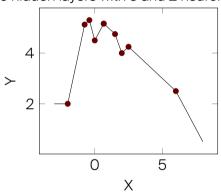
Depth versus Width





$$10 + 6 = 16$$
 parameters

two hidden layers with 3 and 2 neurons



6+8+3=17 parameters

Quiz

- ► The number of input units in a neural network is equal to the number of training samples.
- The activation of an artificial neuron with inputs $x_1 = 1$, $x_2 = 3$, $x_3 = 0$, weights $w_1 = 1$, $w_2 = -1$, $w_3 = 10$, bias $w_0 = 1$ and relu activation function is **A** -1 **B** 0 **C** 1
- For a network with 3-dimensional input, 2 hidden layers of each 10 neurons and one output neuron the number of free parameters (weights and biases) is

 A 24

 B 131

 C 161
- ► The XOR problem could also be solved with $g^{(1)}(x) = x$ and $g^{(2)}(x) = \sigma(x)$ (instead of $g^{(1)}(x) = \text{relu}(x)$ and $g^{(2)}(x) = \sigma(x)$).



Table of Contents

- 1. Solving the XOR Problem Without Feature Engineering
- 2. Artificial Neurons
- 3. Multilayer Perceptrons
- 4. Regression with Multilayer Perceptrons
- 5. Classification with Multilayer Perceptrons



Regression with Multilayer Perceptrons

The output of the neural network is used to parametrize the conditional density.

The parameters are fitted with gradient descent

on the negative log-loglikelihood loss

$$-\log \ell(\theta) = -\sum_{i=1}^{n} \log p(y_i|x_i,\theta).$$

For example: Assume the wind speed in Luzern is distributed normally around some mean that correlates with the measurements done 5 hours earlier.

We take a neural network with as many input neurons as predictors, some hidden neurons and one output neuron. Gradient descent finds the parameters such that the output activity approaches the mean of the conditional density.

$$p(y|x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2} \left(y - g^{(2)}(b^{(2)} + w^{(2)}g^{(1)}(b^{(1)} + w^{(1)}x))\right)^2\right)$$



Regression with Multilayer Perceptrons

A neural network with more outputs can be used to predict more complex densities.

Example: Assume the wind speed in Luzern is log-normally distributed with mean and standard deviation correlating with the pressure in Luzern.

We take a neural network with one input neuron, some hidden neurons and two output neurons. Gradient descent finds the parameters such that the output neurons code for the mean and the variance¹ of the log-normal density.

$$p(y|x) = \frac{1}{\sqrt{2\pi f(a_2^{(2)})}} \exp\left(-\frac{1}{2f(a_2^{(2)})} \left(\log(y) - a_1^{(2)}\right)^2\right)$$



¹The output of the second neuron is additionally transformed with function f to be positive.

Initialization and Data Preprocessing

The initialization of neural networks and the choice the hyper-parameters (number of hidden neurons, non-linearities, learning rate, etc.) is an art.

Common choices are:

biases = 0 weights
$$w_{ij}^{(I)}$$
 sampled uniformly from $[-x, x]$ with $x = \sqrt{\frac{6}{d^{(I-1)} + d^{(I)}}}$ learning rate between 10^{-4} and 10^{-1} .

These choices work typically well for input data between 0 and 1 or standardized input with each predictor having mean 0 and variance 1. For regression it is advisable to also scale or standardize the output.



Regularization

For the weights and biases in each layer one can apply L1 or L2 regularization.

Early stopping in gradient descent can be used.

Using fewer hidden neurons also reduces the flexibility of the neural network.

Another popular and effective regularization method is **Dropout**: During training, for each training example a randomly selected fraction of p neurons is dropped out (inactivated).

This prevents neurons from becoming over-specialized. All neurons are active when testing, but their weights are scaled by 1 - p.



Flexibility of a Neural Network and Its Number of Parameters

More layers or more neurons ⇒ more parameters.

More parameters ⇒ more flexibility?

Not necessarily! Regularization has a strong effect on the flexibility. Even without explicit regularization (L1, L2, Dropout) and without explicitly monitored early stopping one does stop gradient descent usually after some number of iterations and before perfect convergence; therefore one regularizes by implicit early stopping.

Wisely regularized large neural networks often work better than small ones, because they tend not to get stuck at sub-optimal losses.



Why Multilayer Perceptrons?

Flexibility by Composition of Simple Elements.

- Individual neurons should not be simpler: the composition of linear functions is a linear function.
- Individual neurons do not need to be more complex: complexity is achieved by using multiple neurons.
- ▶ With sufficiently many neurons one can approximate any function.



Table of Contents

- 1. Solving the XOR Problem Without Feature Engineering
- 2. Artificial Neurons
- 3. Multilayer Perceptrons
- 4. Regression with Multilayer Perceptrons
- 5. Classification with Multilayer Perceptrons



Classification with Multilayer Perceptrons

With K output neurons we can use neural networks to parametrize categorical distributions suitable for classification problems.

Example: We take $28 \times 28 = 784$ input neurons, some hidden neurons and 10 output neurons to classify MNIST images. The softmax of the 10 output neurons is the predicted probability of the different class labels.

$$P(C_i|x) = s\left(g^{(2)}(b^{(2)} + w^{(2)}g^{(1)}(b^{(1)} + w^{(1)}x))\right)_i$$

where s is the softmax function (see slides "Supervised Learning").



Quiz

- ▶ With L1 or L2 regularization applied to the weights of a neural network, the final weights at the end of gradient descent tend to be smaller than without regularization.
- ▶ With early stopping gradient descent usually stops in a local minimum of the training loss.
- ► A neural network with no hidden layers, sigmoid activation function and negative log-likelihood loss is equivalent to logistic regression.
- Gradient descent in neural networks always finds the global minimum of the loss function of the training set.
- ▶ Which activation function should be chosen in the output layer to predict the mean in a regression setting?

A relu B tanh C identity

Suggested Reading

- ▶ 10.1. Single Layer Neural Networks
- ▶ 10.2. Multilayer Neural Networks

