

Intro to Deep Learning

Machine Learning

Ignacio Sarmiento-Barbieri

Universidad de La Plata

Agenda

- 1 Intro
- 2 Single Layer Neural Networks
- 3 Activation Functions
- 4 Output Functions
- 5 Training the network
- 6 Architecture Design
- 7 Multilayer Neural Networks
- 8 Network Tuning
- 9 Another Example: MNIST

Deep Learning: Intro

- ▶ Neural networks are simple models.
- ▶ Their strength lays in their simplicity
- ▶ The model has linear combinations of inputs that are passed through nonlinear activation functions called nodes (or, in reference to the human brain, neurons).

Deep Learning: Intro

- Let's start with a familiar and simple model, the linear model

$$y = f(X) + u \tag{1}$$

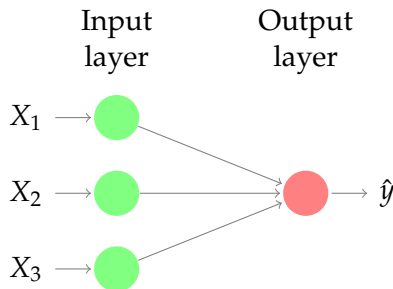
$$y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + u$$

Deep Learning: Intro

- Let's start with a familiar and simple model, the linear model

$$y = f(X) + u \quad (1)$$

$$y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + u$$



Agenda

- 1 Intro
- 2 Single Layer Neural Networks**
- 3 Activation Functions
- 4 Output Functions
- 5 Training the network
- 6 Architecture Design
- 7 Multilayer Neural Networks
- 8 Network Tuning
- 9 Another Example: MNIST

Single Layer Neural Networks

- ▶ Linear Models may miss the nonlinearities that best approximate $f^*(x)$
- ▶ We can overcome these limitations of linear models and handle a more general class of functions by incorporating hidden layers.
- ▶ Neural Networks are also called deep feedforward networks, feedforward neural networks, or multilayer perceptrons (MLPs), and are the quintessential deep learning models

Single Layer Neural Networks

- ▶ A neural network takes an input vector of p variables

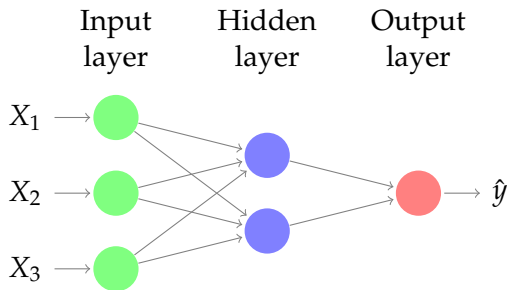
$$X = (X_1, X_2, \dots, X_p) \quad (2)$$

- ▶ and builds a nonlinear function $f(X)$ to predict the response y .

$$y = f(X) + u \quad (3)$$

- ▶ What distinguishes neural networks from previous methods is the particular structure of the model.

Single Layer Neural Networks



Single Layer Neural Networks

- ▶ The NN model has the form

$$f(X) = \beta_0 + \sum_{k=1}^K \beta_k h_k(X) \quad (4)$$

$$= \beta_0 + \sum_{k=1}^K \beta_k g \left(w_{k0} + \sum_{j=1}^p w_{kj} X_j \right) \quad (5)$$

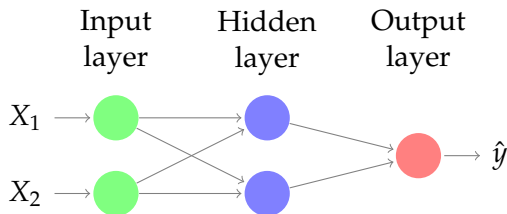
- ▶ where $g(\cdot)$ is a activation function specified in advance
- ▶ where the nonlinearity of $g(\cdot)$ is essential

Worked Example I: Single Layer Neural Networks

- ▶ $p = 2, X = (X_1, X_2)$
- ▶ $K = 2, h_1(X)$ and $h_2(X)$
- ▶ $g(z) = z^2$

Worked Example I: Single Layer Neural Networks

- ▶ $p = 2, X = (X_1, X_2)$
- ▶ $K = 2, h_1(X)$ and $h_2(X)$
- ▶ $g(z) = z^2$



Worked Example I: Single Layer Neural Networks

$$f(X) = \beta_0 + \sum_{k=1}^2 \beta_k g \left(w_{k0} + \sum_{j=1}^2 w_{kj} X_j \right) \quad (6)$$

► Suppose we get

$$\begin{array}{lll} \hat{\beta}_0 = 0 & \hat{\beta}_1 = \frac{1}{4} & \hat{\beta}_2 = -\frac{1}{4} \\ \hat{w}_{10} = 0 & \hat{w}_{11} = 1 & \hat{w}_{12} = 1 \\ \hat{w}_{20} = 0 & \hat{w}_{21} = 1 & \hat{w}_{22} = -1 \end{array}$$

Worked Example I: Single Layer Neural Networks

NN Minimalist Theory

- ▶ Why not a linear activation functions?
- ▶ Let's go back to our example
 - ▶ $p = 2, X = (X_1, X_2)$
 - ▶ $K = 2, h_1(X)$ and $h_2(X)$
 - ▶ Now $g(z) = z$
- ▶ Then

$$f(X) = \beta_0 + \sum_{k=1}^2 \beta_k A_k \quad (7)$$

$$= \beta_0 + \sum_{k=1}^2 \beta_k h_k(X) \quad (8)$$

$$= \beta_0 + \sum_{k=1}^2 \beta_k g \left(w_{k0} + \sum_{j=1}^p w_{kj} X_j \right) \quad (9)$$

NN Minimalist Theory

Why not a linear activation functions?

- ▶ Since $g(z) = z$ we get

$$f(X) = \beta_0 + \sum_{k=1}^2 \beta_k \left(w_{k0} + \sum_{j=1}^2 w_{kj} X_j \right) \quad (10)$$

- ▶ Replacing

$$f(X) = \beta_0 + \beta_1 (w_{10} + w_{11}X_1 + w_{12}X_2) + \beta_2 (w_{20} + w_{21}X_1 + w_{22}X_2) \quad (11)$$

- ▶ then

$$f(X) = \theta_0 + \theta_1 X_1 + \theta_2 X_2 \quad (12)$$

Worked Example II : The "Exclusive OR (XOR)" Function

- ▶ The exclusive disjunction of a pair of propositions, (p, q) , is supposed to mean that p is true or q is true, but not both
- ▶ It's truth table is:

q	p	$q \vee p$
0	0	0
0	1	1
1	0	1
1	1	0

- ▶ When exactly one of these binary values is equal to 1, the XOR function returns 1. Otherwise, it returns 0

Worked Example: The “Exclusive OR (XOR)” Function

- ▶ Let's use a linear model

$$y = \beta_0 + \beta_1 q + \beta_2 p + u \quad (13)$$

$$y = \begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \end{pmatrix} X = \begin{pmatrix} 0 & 0 \\ 0 & 1 \\ 1 & 0 \\ 1 & 1 \end{pmatrix} \mathbf{1} = \begin{pmatrix} 1 \\ 1 \\ 1 \\ 1 \end{pmatrix} \quad (14)$$

- ▶ Solution $\beta_0 = \frac{1}{2}, \beta_1 = 0, \beta_2 = 0$

- ▶ Prediction $\hat{y} = \begin{pmatrix} \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \\ \frac{1}{2} \end{pmatrix}$

Worked Example: The "Exclusive OR (XOR)" Function

- ▶ Let's use Single Layer NN containing two hidden units
- ▶ Activation Function: ReLU: $g(z) = \max\{0, z\}$
- ▶ NN

$$f(X) = \beta_0 + \sum_{k=1}^2 \beta_k g \left(w_{k0} + \sum_{j=1}^2 w_{kj} X_j \right) \quad (15)$$

Worked Example: The "Exclusive OR (XOR)" Function

- Suppose this is the solution to the XOR problem

$$f(x) = \max\{0, XW + W_0\} \beta + \beta_0$$

$$W = \begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}$$

$$W_0 = \begin{pmatrix} 0 & -1 \\ 0 & -1 \\ 0 & -1 \\ 0 & -1 \end{pmatrix}$$

$$\beta = \begin{pmatrix} 1 & -2 \end{pmatrix}$$

$$\beta_0 = 0$$

Worked Example: The "Exclusive OR (XOR)" Function

- Lets work out the example step by step

$$f(x) = \max\{0, XW + W_0\} \beta + \beta_0 \quad (16)$$

Worked Example: The “Exclusive OR (XOR)” Function

- ▶ In this example, we simply specified the solution, then showed that it obtained zero error.
- ▶ In a real situation, obviously we can't guess the solution
- ▶ What we do is to estimate the parameters via optimization methods
- ▶ All gain comes from using nonlinear activation function

Example: XOR



Agenda

- 1 Intro
- 2 Single Layer Neural Networks
- 3 Activation Functions**
- 4 Output Functions
- 5 Training the network
- 6 Architecture Design
- 7 Multilayer Neural Networks
- 8 Network Tuning
- 9 Another Example: MNIST

Activation Functions

ReLU Function

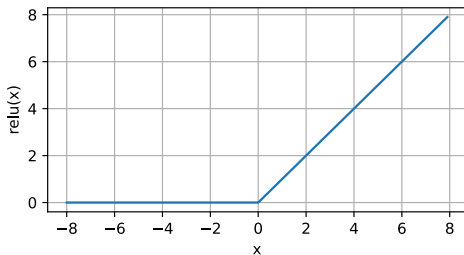
► ReLU Function

- The most popular choice, due to both simplicity of implementation and its good performance on a variety of predictive tasks, is the rectified linear unit (ReLU).
- ReLU provides a very simple nonlinear transformation. Given an element x , the function is defined as the maximum of that element and 0:

$$\text{ReLU}(x) = \max\{x, 0\}.$$

Activation Functions

- ▶ ReLU function retains only positive elements and discards all negative elements by setting them to 0.
- ▶ It is piecewise linear.



Activation Functions

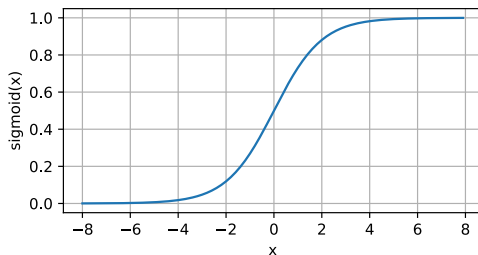
Sigmoid Function (Logit)

- ▶ The sigmoid function transforms its inputs, for which values lie in the domain \mathbb{R} , to outputs that lie on the interval $(0, 1)$.
- ▶ For that reason, the sigmoid is often called a squashing function: it squashes any input in the range $(-\infty, \infty)$ to some value in the range $(0, 1)$:

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

Activation Functions

Sigmoid Function (Logit)



- ▶ In the earliest neural networks, scientists were interested in modeling biological neurons which either fire or do not fire. Thus the pioneers of this field, going all the way back to McCulloch and Pitts, the inventors of the artificial neuron, focused on thresholding units.
- ▶ A thresholding activation takes value 0 when its input is below some threshold and value 1 when the input exceeds the threshold.
- ▶ When attention shifted to gradient based learning, the sigmoid function was a natural choice because it is a smooth, differentiable approximation to a thresholding unit.

Activation Functions

Tanh Function

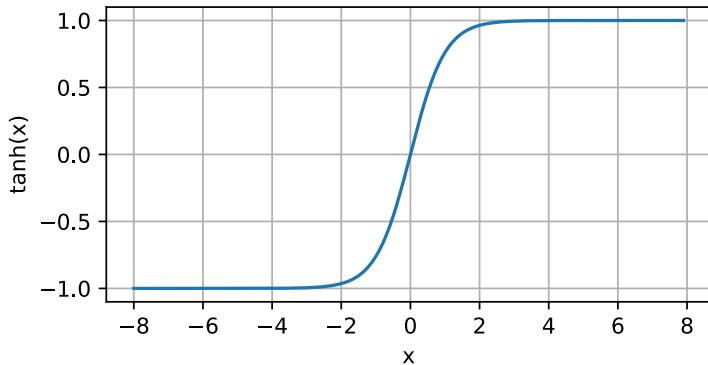
- ▶ Like the sigmoid function, the tanh (hyperbolic tangent) function also squashes its inputs, transforming them into elements on the interval between -1 and 1:

$$\tanh(x) = \frac{1 - \exp(-2x)}{1 + \exp(-2x)}.$$

- ▶ Although the shape of the function is similar to that of the sigmoid function, the *tanh* function exhibits point symmetry about the origin of the coordinate system.

Activation Functions

Tanh Function



Activation Functions

- ▶ Other Activation functions

- ▶ $h = \cos(Wx + b)$ Goodfellow et al. (2016) claim that on the MNIST dataset they obtained an error rate of less than 1 percent
- ▶ Radial basis function (RBF): $\exp\left(-\frac{1}{\sigma^2}||W - x||^2\right)$
- ▶ Softplus: $\log(1 + e^x)$
- ▶ Hard tanh: $\max(-1, \min(1, x))$

- ▶ Hidden unit design remains an active area of research, and many useful hidden unit types remain to be discovered

Agenda

- 1 Intro
- 2 Single Layer Neural Networks
- 3 Activation Functions
- 4 Output Functions**
- 5 Training the network
- 6 Architecture Design
- 7 Multilayer Neural Networks
- 8 Network Tuning
- 9 Another Example: MNIST

Output Functions

- ▶ The choice of cost function is tightly coupled with the choice of output unit.
- ▶ Most of the time, we simply use the distance between the data distribution and the model distribution.
 - ▶ Linear $y = \beta_0 + \sum_{k=1}^K \beta_k h_k \rightarrow \mathbb{R}$
 - ▶ Sigmoid (Logistic) \rightarrow classification $\{0, 1\}$
 - ▶ Softmax \rightarrow classification multiple categories

Agenda

- 1 Intro
- 2 Single Layer Neural Networks
- 3 Activation Functions
- 4 Output Functions
- 5 Training the network**
- 6 Architecture Design
- 7 Multilayer Neural Networks
- 8 Network Tuning
- 9 Another Example: MNIST

Loss Functions

- We want to estimate

$$y_i = f(x) + u$$

- The first thing we need to specify is the loss function.
- For regression problems we used MSE.

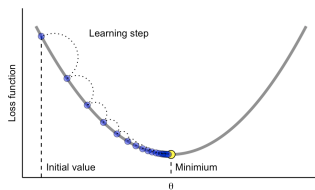
$$E(\theta) = \frac{1}{N} \sum (y - \hat{y}) \quad (17)$$

Loss Function Minimization: Gradient Descent

- ▶ Once we have defined a loss function, we need to minimize $E(\theta)$ with respect to θ
- ▶ Here, we employ an iterative method called gradient descent.
- ▶ Gradient Descent is a very generic optimization algorithm capable of finding optimal solutions to a wide range of problems.

Gradient-Based Descent

- ▶ The general idea of Gradient Descent is to tweak parameters iteratively in order to minimize a cost function.
- ▶ Intuitively, the method corresponds to “walking down the hill” in our many parameter landscape until we reach a (local) minimum.



Source: Boehmke, B., & Greenwell, B. (2019)

Batch Gradient Descent

- ▶ To implement Gradient Descent, you need to compute the gradient of the cost function with regards to each model parameter

$$\beta' = \beta - \epsilon \nabla_{\beta} f(\beta) \quad (18)$$

- ▶ In other words, you need to calculate how much the cost function will change if you change β just a little bit.
- ▶ You also need to define the learning step ϵ

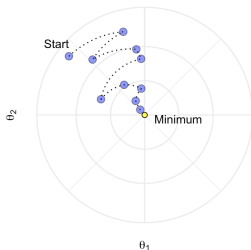
Batch Gradient Descent

- ▶ Notice that this formula involves calculations over the full data set X , at each Gradient Descent step!
- ▶ This is why the algorithm is called Batch Gradient Descent: it uses the whole batch of training data at every step.
- ▶ As a result it is terribly slow on very large data sets.
- ▶ However, Gradient Descent scales well with the number of variables; estimating a Linear Regression model when there are hundreds of thousands of features is much faster using Gradient Descent than using the Normal Equations or any decomposition.

Stochastic Gradient-Based Optimization

- ▶ Using the full data set can be terribly slow, specially in large data sets.
- ▶ At the opposite extreme, Stochastic Gradient Descent just picks a random observation at every step and computes the gradients based only on that single observation.

Stochastic Gradient-Based Optimization



Source: Boehmke, B., & Greenwell, B. (2019)

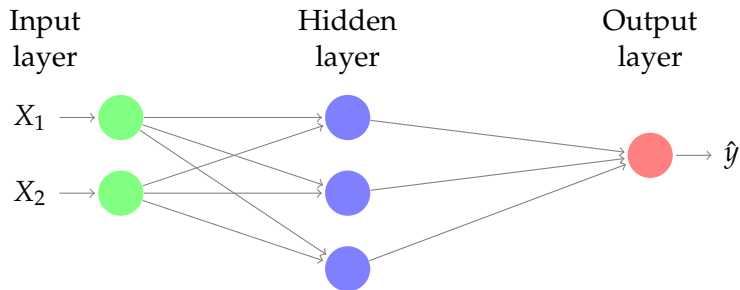
► This makes the algorithm faster but

- Adds some random nature in descending the risk function's gradient.
- Although this randomness does not allow the algorithm to find the absolute global minimum, it can actually help the algorithm jump out of local minima and off plateaus to get sufficiently near the global minimum.

Mini-batch Gradient Descent

- ▶ At each step, instead of computing the gradients based on the full dataset (as in Batch GD) or based on just one observation (as in Stochastic GD),
- ▶ Mini- batch GD computes the gradients on small random sets of observations called mini- batches.

Forward and Back Propagation



Model Accuracy

- ▶ With the training completed, we want to understand how well the final model performs.
- ▶ For that, we can introduce another metric

$$E(\theta) = \frac{1}{N} \sum |y - \hat{y}|$$

Example: Regression



Agenda

- 1 Intro
- 2 Single Layer Neural Networks
- 3 Activation Functions
- 4 Output Functions
- 5 Training the network
- 6 Architecture Design**
- 7 Multilayer Neural Networks
- 8 Network Tuning
- 9 Another Example: MNIST

Architecture Design

- ▶ Another key design consideration for neural networks is determining the architecture.
- ▶ The word architecture refers to the overall structure of the network: how many units it should have and how these units should be connected to each other.
- ▶ The universal approximation theorem guarantees that regardless of what function we are trying to learn, a sufficiently large MLP will be able to represent this function.

Architecture Design

- ▶ The universal approximation theorem (Hornik et al., 1989; Cybenko, 1989) states that:
 - ▶ A feedforward network with a linear output layer and at least one hidden layer with any “squashing” activation function (such as the logistic sigmoid activation function) can approximate any Borel measurable function from one finite-dimensional space to another with any desired nonzero amount of error, provided that the network is given enough hidden units.

Architecture Design

- ▶ We are not guaranteed, however, that the training algorithm will be able to learn that function.
- ▶ Even if the network is able to represent the function, learning can fail for two different reasons.
 - 1 The optimization algorithm used for training may not be able to find the value of the parameters that corresponds to the desired function.
 - 2 The training algorithm might choose the wrong function as a result of overfitting

Architecture Design

- ▶ A feedforward network with a single layer is sufficient to represent any function, but the layer may be infeasibly large and may fail to learn and generalize correctly.
- ▶ In many circumstances, using deeper models can reduce the number of units required to represent the desired function and can reduce the amount of generalization error.
- ▶ The ideal network architecture for a task must be found via experimentation guided by monitoring the validation set error

Agenda

- 1 Intro
- 2 Single Layer Neural Networks
- 3 Activation Functions
- 4 Output Functions
- 5 Training the network
- 6 Architecture Design
- 7 Multilayer Neural Networks**
- 8 Network Tuning
- 9 Another Example: MNIST

Multilayer Neural Networks

- ▶ Modern neural networks typically have more than one hidden layer, and often many units per layer.
- ▶ In theory a single hidden layer with a large number of units has the ability to approximate most functions.
- ▶ However, the learning task of discovering a good solution is made much easier with multiple layers each of modest size.

Multilayer Neural Networks: MNIST Digits



Agenda

- 1 Intro
- 2 Single Layer Neural Networks
- 3 Activation Functions
- 4 Output Functions
- 5 Training the network
- 6 Architecture Design
- 7 Multilayer Neural Networks
- 8 Network Tuning**
- 9 Another Example: MNIST

Network Tuning

- ▶ Training networks requires a number of choices that all have an effect on the performance:
 - ▶ The number of hidden layers,
 - ▶ The number of units per layer
 - ▶ Regularization tuning parameters
 - ▶ Details of stochastic gradient descent.
- ▶ This is an active research area that involves a lot of trial and error, and overfitting is a latent danger at each step.

Agenda

- 1 Intro
- 2 Single Layer Neural Networks
- 3 Activation Functions
- 4 Output Functions
- 5 Training the network
- 6 Architecture Design
- 7 Multilayer Neural Networks
- 8 Network Tuning
- 9 Another Example: MNIST

Example: MNIST



Example: MNIST

Loss Functions

- ▶ We want to estimate

$$y_i = f(x) + u$$

- ▶ The first thing we need to specify is the loss function.
- ▶ For regression problems we used MSE, MAE.
- ▶ Cross-entropy is a commonly used loss function for (multi-class) classification tasks.

$$E(\theta) = - \sum_{i=1}^K I(Y = k) \log_2 (\hat{Pr}(Y = k))$$

- ▶ where p is the empirical distribution and \hat{p} is the estimated distribution
- ▶ for 2 classes it becomes

$$l(\theta) = - [y \log_2 (\hat{Pr}(Y = 1)) + (1 - y) \log_2 (1 - \hat{Pr}(Y = 1))] \quad (19)$$

Example: MNIST

Model Accuracy

- ▶ With the training completed, we want to understand how well the final model performs in recognizing handwritten digits.
- ▶ For that, we introduce the accuracy defined by

$$\text{accuracy} = \frac{\text{correct predictions}}{\text{total predictions}}.$$

Example: MNIST

