

Deep Learning

Chapter 6

Deep Feedforward Networks

Plan

- Introduction
- An Example: XOR
- Gradient Based Learning

Introduction

- Goal of the network is to approximate f^* , for some function f^* .
- Feedforward Network because no feedback connection from output.
- Network with multiple Layers.
- Depth is the number of layers.

Training Neural Network

- $f^*(x)$ - output for input
- We try to match $f(x)$ to $f^*(x)$
- Training Data consists of noisy, approximate examples of $f^*(x)$
- x accompanied with $y \approx f^*(x)$

Why Feedforward Network??

- Linear Models are appealing as they can be fit efficiently
- But model capacity is limited to linear functions
- To extend linear models to represent nonlinear functions of x , we can apply the linear model not to x itself but to a transformed input $\varphi(x)$, where φ is a nonlinear transformation
- We can think of φ as providing a set of features describing x

How to choose the mapping ϕ ??

- Option 1: use a very generic ϕ , such as the infinite-dimensional ϕ that is implicitly used by kernel machines based on the RBF kernel
- Problem: Generalization to the test set often remains poor

How to choose the mapping ϕ ??

- Option 2: manually engineer ϕ
- Problem: this approach requires decades of human effort for each separate task .

How to choose the mapping ϕ ??

- Option 2: Deep Learning
- In this case ϕ is represented in hidden layer.

An Example: XOR

- The XOR function (“exclusive or”) is an operation on two binary values, x_1 and x_2 . When exactly one of these binary values is equal to 1, the XOR function returns 1. Otherwise, it returns 0
- The XOR function provides the target function $y = f^*(\mathbf{x})$ that we want to learn. Our model provides a function $y = f(\mathbf{x};\theta)$ and our learning algorithm will adapt the parameters θ to make f as similar as possible to f^* .
- We want our network to perform correctly on the four points

$$\mathbb{X} = \{[0, 0]^\top, [0, 1]^\top, [1, 0]^\top, \text{ and } [1, 1]^\top\}$$

Learning XOR

Evaluated on our whole training set, the MSE loss function is

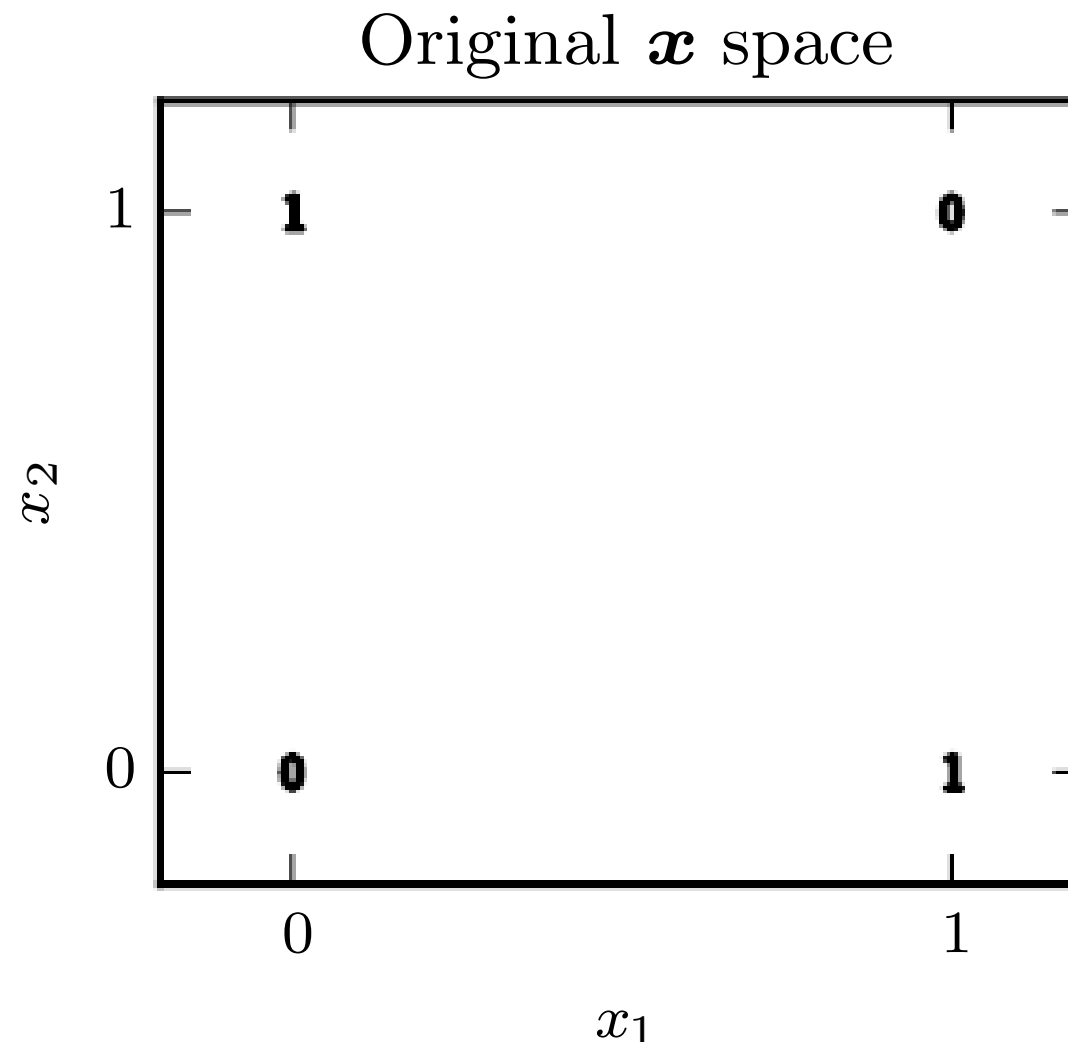
$$J(\boldsymbol{\theta}) = \frac{1}{4} \sum_{\mathbf{x} \in \mathbb{X}} (f^*(\mathbf{x}) - f(\mathbf{x}; \boldsymbol{\theta}))^2. \quad (6.1)$$

Now we must choose the form of our model, $f(\mathbf{x}; \boldsymbol{\theta})$. Suppose that we choose a linear model, with $\boldsymbol{\theta}$ consisting of \mathbf{w} and b . Our model is defined to be

$$f(\mathbf{x}; \mathbf{w}, b) = \mathbf{x}^\top \mathbf{w} + b. \quad (6.2)$$

We can minimize $J(\boldsymbol{\theta})$ in closed form with respect to \mathbf{w} and b using the normal equations.

After solving the normal equations, we obtain $\mathbf{w} = 0$ and $b = 1/2$. The linear model simply outputs 0.5 everywhere. Why does this happen?



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 must apply a fixed coefficient w_2 to x_2 . The linear model therefore cannot use the value of x_1 to change the coefficient on x_2 and cannot solve this problem

Feedforward Network for Learning XOR

$$\mathbf{h} = \mathbf{g}(\mathbf{W}^\top \mathbf{x} + \mathbf{c}) \quad \text{where} \quad g(z) = \max\{0, z\}$$

and finally

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

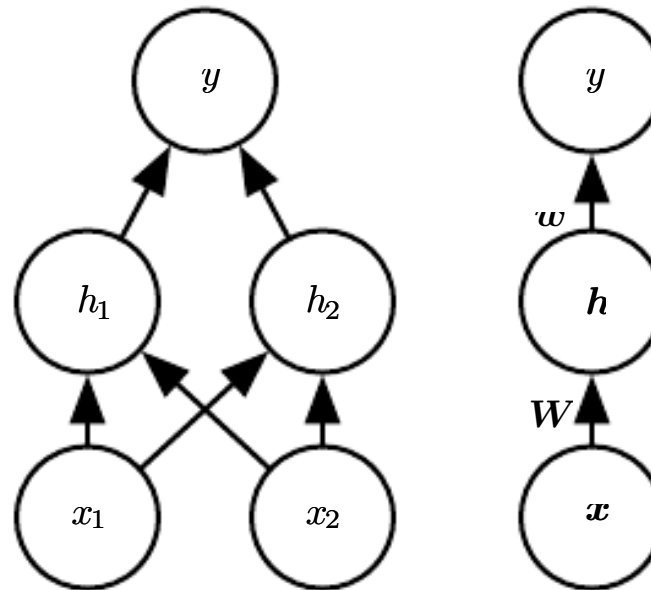


Figure 6.2: An example of a feedforward network, drawn in two different styles. Specifically, this is the feedforward network we use to solve the XOR example. It has a single hidden layer containing two units. *(Left)* In this style, we draw every unit as a node in the graph. This style is very explicit and unambiguous but for networks larger than this example it can consume too much space. *(Right)* In this style, we draw a node in the graph for each entire vector representing a layer's activations. This style is much more compact. Sometimes we annotate the edges in this graph with the name of the parameters that describe the relationship between two layers. Here, we indicate that a matrix W describes the mapping from x to h , and a vector w describes the mapping from h to y . We typically omit the intercept parameters associated with each layer when labeling this kind of drawing.

Solution!!

$$\mathbf{W} = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix},$$

$$\mathbf{c} = \begin{bmatrix} 0 \\ -1 \end{bmatrix},$$

$$\mathbf{w} = \begin{bmatrix} 1 \\ -2 \end{bmatrix},$$

$$b=0$$

$$\begin{bmatrix} 174 \\ \end{bmatrix}$$

Gradient Based Learning

- Neural network causes most interesting loss functions to become non-convex
- Neural networks are usually trained by using iterative, gradient-based optimizers
- Convex optimization algorithms with global convergence guarantees, whereas stochastic gradient descent applied to non-convex loss functions has no such convergence guarantee
- For Feedforward networks training algorithms are improvements and refinements on the ideas of gradient descent

Cost Functions

- Most popular Cost Function comes from MLE
- Alternatively Cost Function is the cross-entropy between the training data and the model distribution

- $J(\boldsymbol{\theta}) = -\mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} \log p_{\text{model}}(\mathbf{y} \mid \mathbf{x})$

- If $p_{\text{model}}(\mathbf{y} \mid \mathbf{x}) = \mathcal{N}(\mathbf{y}; f(\mathbf{x}; \boldsymbol{\theta}), \mathbf{I})$ then,

$$J(\theta) = \frac{1}{2} \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \hat{p}_{\text{data}}} \|\mathbf{y} - f(\mathbf{x}; \boldsymbol{\theta})\|^2 + \text{const.}$$

- Advantage of MLE: Specifying a model $p(\mathbf{y} \mid \mathbf{x})$ automatically determines a cost function $\log p(\mathbf{y} \mid \mathbf{x})$
- The negative log-likelihood helps to avoid the saturation problem for many models

Learning Conditional Statistics

- Instead of learning a full probability distribution $p(y | x; \theta)$ we often want to learn just one conditional statistic of y given x

- Neural Networks to represent any function f from a wide class of functions. Hence cost function is a functional

- Hence $f^* = \arg \min_f \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}} ||\mathbf{y} - f(\mathbf{x})||^2$ yields,

$$f^*(\mathbf{x}) = \mathbb{E}_{\mathbf{y} \sim p_{\text{data}}(\mathbf{y}|\mathbf{x})} [\mathbf{y}].$$

- Also $f^* = \arg \min_f \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim p_{\text{data}}} ||\mathbf{y} - f(\mathbf{x})||_1$ yields the median

Linear Output Units for Gaussian output Distributions

- Given features \mathbf{h} , a layer of linear output units produces a vector $\hat{\mathbf{y}} = \mathbf{W}^\top \mathbf{h} + \mathbf{b}$.
- Linear output layers are often used to produce the mean of a conditional Gaussian distribution:

$$p(\mathbf{y} \mid \mathbf{x}) = \mathcal{N}(\mathbf{y}; \hat{\mathbf{y}}, \mathbf{I})$$

- MLE framework makes it straightforward to learn the covariance of the Gaussian
- The covariance must be constrained to be a positive definite matrix for all inputs which is difficult with linear output layer

Sigmoid Output Units for Bernoulli Output Distributions

- Suppose we were to use a linear unit, and threshold its value to obtain a valid probability:

$$P(y = 1 \mid \mathbf{x}) = \max \left\{ 0, \min \left\{ 1, \mathbf{w}^\top \mathbf{h} + b \right\} \right\}$$

- Problem: Any time that $\mathbf{w}^\top \mathbf{h} + b$ strayed outside the unit interval, the gradient of the output of the model with respect to its parameters would be 0
- Instead use sigmoid unit: $\sigma(\mathbf{w}^\top \mathbf{h} + b)$

Sigmoid to Normalize

$$\log \tilde{P}(y) = yz$$

$$\tilde{P}(y) = \exp(yz)$$

$$P(y) = \frac{\exp(yz)}{\sum_{y'=0}^1 \exp(y' z)}$$

$$P(y) = \sigma \left((2y - 1)z \right) .$$

Importance of Sigmoid

- The cost function used with maximum likelihood is $-\log P(y \mid \mathbf{x})$, the log in the cost function undoes the exp of the sigmoid
- Without this effect, the saturation of the sigmoid could prevent gradient-based learning from making good progress
- Loss Function:

$$\begin{aligned} J(\boldsymbol{\theta}) &= -\log P(y \mid \mathbf{x}) \\ &= -\log \sigma((2y - 1)z) \\ &= \zeta((1 - 2y)z). \end{aligned}$$

- Saturation thus occurs only when the model already has the right answer—when $y = 1$ and z is very positive, or $y = 0$ and z is very negative
- When we use other loss functions, such as mean squared error, the loss can saturate anytime $\sigma(z)$ saturates

Other Output Units

- Softmax Units for Multinoulli Output Distributions
- We may wish to learn the variance of a conditional Gaussian for y , given x . there is a closed form expression
- Alternative approach is to simply include the variance as one of the properties of the distribution $p(y \mid x)$ that is controlled by $\omega = f(x; \theta)$.
- In the simple case where the standard deviation does not depend on the input, we can make a new parameter in the network that is copied directly into ω

Other Output Units

- We often want to perform multimodal regression, that is, to predict real values that come from a conditional distribution $p(y \mid x)$ that can have several different peaks in y space for the same value of x
- In this case, a Gaussian mixture is a natural representation for the output
- A Gaussian mixture output with n components is defined by the conditional probability distribution:

$$p(\mathbf{y} \mid \mathbf{x}) = \sum_{i=1}^n p(c = i \mid \mathbf{x}) \mathcal{N}(\mathbf{y}; \boldsymbol{\mu}^{(i)}(\mathbf{x}), \boldsymbol{\Sigma}^{(i)}(\mathbf{x}))$$