Neural Networks

CMPUT 261: Introduction to Artificial Intelligence

P §3.1-3.6

Lecture Outline

- Recap
- 2. Nonlinear models
- 3. Feedforward neural networks

After this lecture, you should be able to:

- define an activation function

- describe the basic procedure for training a neural network
- identify the parameters of a feedforward neural network

define a rectified linear activation and give an expression for its value describe how the units in a feedforward neural network are connected give an expression in matrix notation for a layer of a feedforward network explain (high level) what the Universal Approximation Theorem guarantees

Recap: Supervised Learning

• Supervised learning task:

predict the values of target features Y based on input features X

- Formally: Choose a hypothesis $h: \mathcal{X} \to \mathcal{Y}$ from a hypothesis space \mathcal{H}
- We use the value of a loss function L applied to a set of training examples $S = \{(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)})\}$ to choose the hypothesis
 - **Regularization** penalty biases optimization toward simpler functions:

$$\hat{h} = \arg\min_{h \in \mathcal{H}} L(h) = \sum_{i=1}^{n} \ell(h(\mathbf{x}^{(i)}), y^{(i)}) + \lambda \operatorname{penalty}(h)$$

- Simpler functions are more likely to generalize
- Generalization performance is evaluated on the test set \bullet
- Another way to reduce overfitting: Learn **distribution** over hypotheses (Bayesian)
 - Many regularization approaches amount to a MAP estimate with a particular prior

(Generalized) Linear Models

- Linear model
 - Linear classification / regression \bullet
 - Logistic regression \bullet
- lacksquare
- **Disadvantages:** Can be really **limited**



Advantages: Efficient to fit (closed form sometimes!)

Example: XOR

- The function $h(x_1, x_2) = (x_1 \text{ XOR } x_2)$ is not linearly separable
 - There is no way to draw a straight line with all of the 1's on one side and all of the 0's on the other
 - This means that no linear model can represent XOR exactly; there will always be some errors
- **Question:** What else could we do?





(Image: Goodfellow 2017)

Nonlinear Features

 $h(\mathbf{x}; \mathbf{w}) = g(\mathbf{w})$

One option: Learn a linear model on richer inputs

- 2. Learn a linear model of the **features** instead of the **inputs**

$$h(\mathbf{x}; \mathbf{w}) = g(\mathbf{w}^{\top} \boldsymbol{\phi}(\mathbf{x})) = g\left(\sum_{j=1}^{d} w_{j}[\boldsymbol{\phi}(\mathbf{x})]_{j}\right)$$

$$(\mathbf{x}^{\mathsf{T}}\mathbf{x}) = g\left(\sum_{j=1}^{d} w_j x_j\right)$$

1. Define a feature mapping $\phi(\mathbf{x})$ that returns functions of the original inputs

Nonlinear Features for XOR

Question: What additional features would help?

• The product of x_1 and $x_2!$

•
$$\phi(x_1, x_2) = [1, x_1, x_2, x_1 x_2]^{\mathsf{T}}$$

- $\mathbf{w} = [-0.2, 0.5, 0.5, -2]^{\top}$
- $h(\mathbf{x}; \mathbf{w}) = \mathbf{w}^{\mathsf{T}} \phi(\mathbf{x}) > 0$ for (0,1) and (1,0) $h(\mathbf{x}; \mathbf{w}) = \mathbf{w}^{\mathsf{T}} \phi(\mathbf{x}) < 0$ for (1,1) and (0,0)



Learning Nonlinear Features

- Manually constructing good features is hard
- Manually constructed features are not transferrable between domains
 - e.g., SIFT features were a revolution in computer vision, but are **only** for computer vision
- Deep learning aims to learn ϕ automatically from the data

Neural Units

- Deep learning learns ϕ by composing little functions
- These function are called **units**



Question: How is this different from a generalized linear model?

- A neural network is many units composed together
- Feedforward neural network: Units arranged into layers
 - Each layer takes outputs of **previous layer** as its **inputs**



Feedforward Neural Network

$$h_{1} = g \left(w_{1,1}^{(1)} x_{1} + w_{1,2}^{(1)} x_{2} + b_{1}^{(1)} \right)$$

$$h_{1} w_{1}^{(2)}$$

$$w_{2}^{(2)} y y = g \left(w_{1}^{(2)} h_{1} + w_{2}^{(2)} h_{2} + b^{(2)} \right)$$

$$h_{2} = g \left(w_{2,1}^{(1)} x_{1} + w_{2,2}^{(1)} x_{2} + b_{2}^{(1)} \right)$$

Example: XOR network



- Activation: $g(z) = \max\{0, z\}$ ("rectified linear unit")
- Offsets: 0 lacksquare
- Weights:
 - [+1, -1] for h_1 ; [-1, +1] for h_2
 - [+1, +1] for y

Question:

When does $h_1 = 1$?

Matrix Representation of Layers • You can think of the **outputs** of each layer as a vector h into a matrix $\mathbf{W} \in \mathbb{R}^{k \times d}_{\uparrow}$ X_2

- The weights from all the outputs of a previous layer to each of the units of the layer can be collected
- The offset term for each unit can be collected into a vector $\mathbf{b} \in \mathbb{R}^k$:

$$\mathbf{h} = g\left(\mathbf{W}\mathbf{x} + \mathbf{b}\right)$$



Design decisions:

- **Depth:** number of layers 1.
- 2. Width: number of nodes in each layer
- 3. Fully connected?

Universal Approximation Theorem

- Theorem: (Hornik et al. 1989; Cybenko 1989; Leshno et al. 1993) A feedforward network with one hidden layer with a "squashing" activation or rectified linear activation and a linear output layer can approximate any function to within any given error bound, given enough hidden units.
- function we're trying to learn!
- **Question:** Why bother with multiple layers? (i.e., depth > 1)

• So a wide but shallow feedforward network can represent any

Neural Network Parameters



A neural network is just a **supervised model**

- parameters θ
- Question: What is θ in a feedforward neural network?

• It is a function that takes inputs \mathbf{x} , and computes an output y based on

Training Neural Networks

- Specify a loss L and a set of training examples:
- Training by gradient descent: \bullet

 - Compute gradient of loss:
 - 3. Update parameters to make loss smaller:

 $S = (\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(n)}, y^{(n)})$ Loss function (e.g., squared error) Compute loss on training data: $L(\mathbf{W}, \mathbf{b}) = \sum \ell(h(\mathbf{x}^{(i)}; \mathbf{W}, \mathbf{b}), y^{(i)})$ Prediction Target i=1 $\nabla L(\mathbf{W}, \mathbf{b})$ (Subsequent lecture)

$$\begin{bmatrix} \mathbf{W}^{new} \\ \mathbf{b}^{new} \end{bmatrix} = \begin{bmatrix} \mathbf{W}^{old} \\ \mathbf{b}^{old} \end{bmatrix} - \eta \nabla L(\mathbf{W}^{old}, \mathbf{b}^{old})$$

Hidden Unit Activations

- Default choice: Rectified linear units (ReLU) $g(z) = \max\{0, z\}$
- Other common types: \bullet
 - tanh(z)

•
$$\frac{1}{1+e^{-z}}$$
 (sigmoid)

• Sigmoid suffers from vanishing gradients; ReLU does not



Torch: Representating Layers

Input will be a tensor with dimension `(n_datapoints, n_channels, height, width)` # So `(64, 1, 30, 30)` in the case of MNIST dataset self.layers = Sequential(

Flatten(), Linear(in_features=30*30, out_features=512), $\# \rightarrow (64, 512)$ ReLU(), Linear(in_features=512, out_features=10), $\# \rightarrow (64, 10)$ LogSoftmax(dim=1)

Torch thinks about **operations** rather than **units**

- We've been thinking about a unit as "weighted sum and then activation" \bullet
- Torch specifies the weighting (Linear) and then the activation (ReLU) separately
- This is especially handy for output layers, where you often want to normalize by the sum of all the outputs (e.g., LogSoftmax)

-> (64, 900) # -> (64, 512) # -> (64, 10)



Summary

- Generalized linear models are insufficiently expressive for many applications
- Composing GLMs into a network is arbitrarily expressive
 - A neural network with a single hidden layer can approximate any function
 - But the network might need to be impractically large, prone to overfitting, or inefficient to train
- Neural networks are trained using variants of gradient descent
- Architectural choices can make a network easier to train, less prone to overfitting