

# Introduction to Neural Networks

Klint Kanopka

`kkanopka@stanford.edu`

Wednesday, May 15, 2019

## 1 Introduction

- Logistic Regression
- Maximum Likelihood Estimation
- Motivations for Neural Networks

## 2 Structure and Design

## 3 Estimation and Application

## 4 Other Neural Architectures

- Recurrent Neural Networks
- Convolutional Neural Network
- Deep Learning

# Logistic Regression

- Computer scientists call the logistic function the “sigmoid”

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

- Computer scientists call the logistic function the “sigmoid”

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

- Maps inputs to real numbers on the interval  $(0, 1)$

- Computer scientists call the logistic function the “sigmoid”

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

- Maps inputs to real numbers on the interval  $(0, 1)$
- Good for turning numbers into probabilities

- Computer scientists call the logistic function the “sigmoid”

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

- Maps inputs to real numbers on the interval  $(0, 1)$
- Good for turning numbers into probabilities
- Has a really simple derivative:

$$\frac{d\sigma}{dx} = \sigma(x)(1 - \sigma(x))$$

- Computer scientists call the logistic function the “sigmoid”

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

- Maps inputs to real numbers on the interval  $(0, 1)$
- Good for turning numbers into probabilities
- Has a really simple derivative:

$$\frac{d\sigma}{dx} = \sigma(x)(1 - \sigma(x))$$

- This is super useful for neural nets, as we'll discuss later.



# Maximum Likelihood Estimation

# Maximum Likelihood Estimation

- Models use parameters to relate data within observations

# Maximum Likelihood Estimation

- Models use parameters to relate data within observations
- MLE specifies a likelihood function that describes how likely (probable) it is to observe the data we have, given a set of model parameters

# Maximum Likelihood Estimation

- Models use parameters to relate data within observations
- MLE specifies a likelihood function that describes how likely (probable) it is to observe the data we have, given a set of model parameters
- MLE searches for model parameters that maximize the likelihood function

# Maximum Likelihood Estimation

- Models use parameters to relate data within observations
- MLE specifies a likelihood function that describes how likely (probable) it is to observe the data we have, given a set of model parameters
- MLE searches for model parameters that maximize the likelihood function
- Often work in terms of log likelihood



- Subfield of Artificial Intelligence

# Machine Learning

- Subfield of Artificial Intelligence
- Computers “learn” models from data



- Subfield of Artificial Intelligence
- Computers “learn” models from data
  - Instead of specifying a strict functional form (linear, logistic, etc.) and doing a regression, the computer learns the shape of the function and its parameters from the data

- Subfield of Artificial Intelligence
- Computers “learn” models from data
  - Instead of specifying a strict functional form (linear, logistic, etc.) and doing a regression, the computer learns the shape of the function and its parameters from the data
  - The data the computer learns from is called the “training set”

- Subfield of Artificial Intelligence
- Computers “learn” models from data
  - Instead of specifying a strict functional form (linear, logistic, etc.) and doing a regression, the computer learns the shape of the function and its parameters from the data
  - The data the computer learns from is called the “training set”
- Works great for prediction and classification tasks

- Subfield of Artificial Intelligence
- Computers “learn” models from data
  - Instead of specifying a strict functional form (linear, logistic, etc.) and doing a regression, the computer learns the shape of the function and its parameters from the data
  - The data the computer learns from is called the “training set”
- Works great for prediction and classification tasks
- Less useful for inference (and causal inference), but this is changing

# Why Neural Nets?

# Why Neural Nets?

- We have tons of data
- Compute time is getting increasingly cheap

# Why Neural Nets?

- We have tons of data
- Compute time is getting increasingly cheap
- The quality of an estimate is highly dependent on model specification

# Why Neural Nets?

- We have tons of data
- Compute time is getting increasingly cheap
- The quality of an estimate is highly dependent on model specification
- More flexible estimation techniques allow computers to learn functional forms



# Why Neural Nets?

- We have tons of data
- Compute time is getting increasingly cheap
- The quality of an estimate is highly dependent on model specification
- More flexible estimation techniques allow computers to learn functional forms
- Neural networks are “universal approximators” and can learn arbitrarily complex functional forms

# Guiding Idea

# Guiding Idea

- Neural networks are built out of *nodes*.

# Guiding Idea

- Neural networks are built out of *nodes*.
- These nodes are organized into *hidden layers*.

# Guiding Idea

- Neural networks are built out of *nodes*.
- These nodes are organized into *hidden layers*.
- For now, each node consists of an individual sigmoid function (logistic regression)

# Guiding Idea

- Neural networks are built out of *nodes*.
- These nodes are organized into *hidden layers*.
- For now, each node consists of an individual sigmoid function (logistic regression)
- *Examples on the board...*

# Training/Test Data

# Training/Test Data

- Conventional wisdom says to split your data 80/20 into training/test data



# Training/Test Data

- Conventional wisdom says to split your data 80/20 into training/test data
- This is often terrible advice

# Training/Test Data

- Conventional wisdom says to split your data 80/20 into training/test data
- This is often terrible advice
- The purpose of this is to hold out some data to test how well your model generalizes

# Training/Test Data

- Conventional wisdom says to split your data 80/20 into training/test data
- This is often terrible advice
- The purpose of this is to hold out some data to test how well your model generalizes
- If you only have a little bit of data, this will work. If you have more, use a *dev set*.

# Training/Test Data

- Conventional wisdom says to split your data 80/20 into training/test data
- This is often terrible advice
- The purpose of this is to hold out some data to test how well your model generalizes
- If you only have a little bit of data, this will work. If you have more, use a *dev set*.
- The dev set helps you tune *hyperparameters* like number of nodes, number of hidden layers, and learning rate

# Training/Test Data

- Conventional wisdom says to split your data 80/20 into training/test data
- This is often terrible advice
- The purpose of this is to hold out some data to test how well your model generalizes
- If you only have a little bit of data, this will work. If you have more, use a *dev set*.
- The dev set helps you tune *hyperparameters* like number of nodes, number of hidden layers, and learning rate
- Good splits (train/dev/test) are 80/10/10, 90/5/5, or even 98/1/1 if you have a ton of data

# Training/Test Data

- Conventional wisdom says to split your data 80/20 into training/test data
- This is often terrible advice
- The purpose of this is to hold out some data to test how well your model generalizes
- If you only have a little bit of data, this will work. If you have more, use a *dev set*.
- The dev set helps you tune *hyperparameters* like number of nodes, number of hidden layers, and learning rate
- Good splits (train/dev/test) are 80/10/10, 90/5/5, or even 98/1/1 if you have a ton of data
- More training data translates into a more robust model

# Loss Functions

# Loss Functions

- A function that defines how wrong your current predictions are



# Loss Functions

- A function that defines how wrong your current predictions are
- For binary classification (1 or 0):

$$L = - \sum_{i=1}^N y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)$$

# Loss Functions

- A function that defines how wrong your current predictions are
- For binary classification (1 or 0):

$$L = - \sum_{i=1}^N y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)$$

- For continuous variables:

$$L = \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

# Loss Functions

- A function that defines how wrong your current predictions are
- For binary classification (1 or 0):

$$L = - \sum_{i=1}^N y_i \log \hat{y}_i + (1 - y_i) \log(1 - \hat{y}_i)$$

- For continuous variables:

$$L = \sum_{i=1}^N (y_i - \hat{y}_i)^2$$

- General goal: Minimize how wrong you are

- Estimation is done with an algorithm called *backpropagation*
- Similar to MLE, except now we search over parameters to *minimize* the loss function
- Updates to the individual parameters in the model are typically done via *gradient descent*
- Similar to Newton-Raphson, but doesn't require inverting a Hessian and requires more iterations to converge - in big data applications, often faster (smaller but quicker steps)

Clone or download the files from the github repo:

[https://github.com/klintkanopka/nn\\_workshop](https://github.com/klintkanopka/nn_workshop)

# Feed-Forward Neural Networks

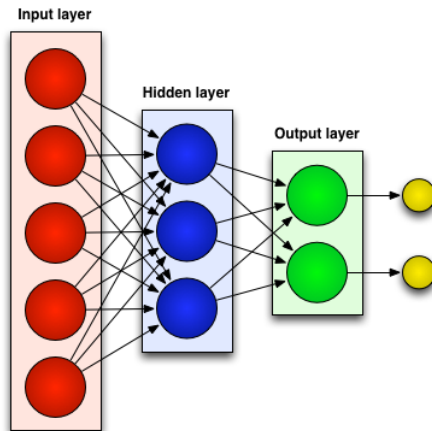


Figure 1: Feed-Forward Neural Network. Image Credit: Joseph Wilk

# Recurrent Neural Networks

- Operate on series or sequence of data



# Recurrent Neural Networks

- Operate on series or sequence of data
- The same network is applied at each time step

# Recurrent Neural Networks

- Operate on series or sequence of data
- The same network is applied at each time step
- Each time step also takes the output of the previous time step as input

# Recurrent Neural Networks

- Operate on series or sequence of data
- The same network is applied at each time step
- Each time step also takes the output of the previous time step as input
- Specialized versions of the RNN work to help “remember” older data:

# Recurrent Neural Networks

- Operate on series or sequence of data
- The same network is applied at each time step
- Each time step also takes the output of the previous time step as input
- Specialized versions of the RNN work to help “remember” older data:
  - Gated Recurrent Unit (GRU)

# Recurrent Neural Networks

- Operate on series or sequence of data
- The same network is applied at each time step
- Each time step also takes the output of the previous time step as input
- Specialized versions of the RNN work to help “remember” older data:
  - Gated Recurrent Unit (GRU)
  - Long Short Term Memory (LSTM)

# Recurrent Neural Networks

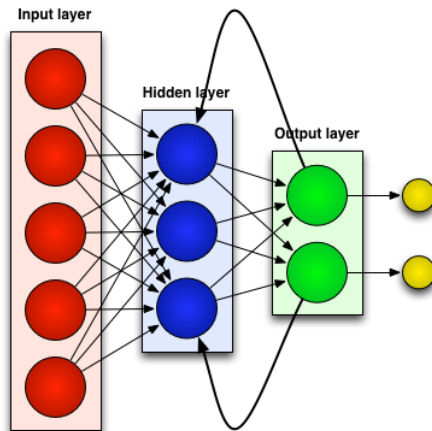


Figure 2: Recurrent Neural Network. Image Credit: Joseph Wilk

# Recurrent Neural Networks

Useful for:

# Recurrent Neural Networks

Useful for:

- Sequences of data



# Recurrent Neural Networks

Useful for:

- Sequences of data
- Item responses

Useful for:

- Sequences of data
- Item responses
- Text

Useful for:

- Sequences of data
- Item responses
- Text
- Time series data

# Recurrent Neural Networks

Useful for:

- Sequences of data
- Item responses
- Text
- Time series data
- Panel data

# Convolutional Neural Networks

# Convolutional Neural Networks

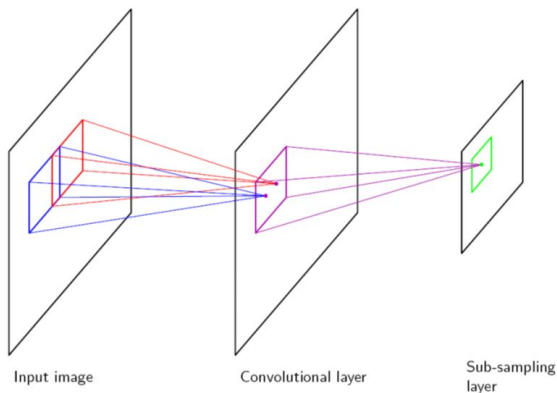
- Uses a window (or kernel) to perform arithmetic reductions on subsets of the input

# Convolutional Neural Networks

- Uses a window (or kernel) to perform arithmetic reductions on subsets of the input
- These windows are called filters and can learn to detect features in the input

# Convolutional Neural Networks

## Convolutional neural networks in 2-D (from Le Cun et al, 1989)





# Convolutional Neural Networks

Useful for:

# Convolutional Neural Networks

Useful for:

- Arrays of data

# Convolutional Neural Networks

Useful for:

- Arrays of data
- Images

# Convolutional Neural Networks

Useful for:

- Arrays of data
- Images
- Sequences of data

# Convolutional Neural Networks

Useful for:

- Arrays of data
- Images
- Sequences of data
- Text

# Convolutional Neural Networks

Useful for:

- Arrays of data
- Images
- Sequences of data
- Text
- Signal processing



- Tech bro buzzword for large neural networks



# Deep Learning

- Tech bro buzzword for large neural networks
- "Deep" refers lots of hidden layers

- Tech bro buzzword for large neural networks
- "Deep" refers lots of hidden layers
- Leverages huge amounts of data

- Tech bro buzzword for large neural networks
- "Deep" refers lots of hidden layers
- Leverages huge amounts of data
- Trained/run on cloud servers with GPUs/TPUs

- Tech bro buzzword for large neural networks
- "Deep" refers lots of hidden layers
- Leverages huge amounts of data
- Trained/run on cloud servers with GPUs/TPUs
- Good word to stick into things you want to get published or funded

Thank you for your attention.

Klint Kanopka  
kkanopka@stanford.edu

To download today's materials:  
[https://github.com/klintkanopka/nn\\_workshop](https://github.com/klintkanopka/nn_workshop)