BIG DATA

DATA SCIENCE DEEP LEARNING FOR STATISTICIANS

Ming Li @ Amazon

mli@alumni.iastate.edu

Hui Lin @ Netlify

hui@linhui.org

Links to Notebooks, Books and URLs

- 1. Course homepage: <u>https://course2019.netlify.com/</u>
- 2. Databrick free <u>community edition</u> account open: <u>link</u>
- 3. Perceptron notebook: link
- 4. Adaline notebook: link
- 5. Feedforward neural network notebook: link
- 6. Convolutional neural network notebook: link
- 7. Recurrent neural network notebook: link
- 8. Big Data Platform notebook: link
- 9. Data preprocessing notebook: link
- 10. Data wrangling notebook: link
- 11. Industry recommendations for academic data science programs: link
- 12. Deep Learning Using R, François Chollet with J. J. Allaire, ISBN 9781617295546 (2018)
- 13. Python Machine Learning by Sebastian Raschka, ISBN-13: 978-1787125933 (2018)
- 14. https://keras.rstudio.com/
- 15. http://spark.rstudio.com/
- 16. https://databricks.com/spark/about
- 17. https://github.com/onnx/onnx



THE NEW TOOL IN DATA SCIENTIST'S TOOL BOX

A Little Bit of History – Perceptron

- Fun video: <u>https://www.youtube.com/watch?v=cNxadbrN_al</u>
- Classification of N points into 2 classes: -1 and +1 (i.e. two different colors in the picture below)
- In this example below, only two features to use (X1 and X2)



A Little Bit of History – Perceptron

For i in 1:M: We set a maximum of epochs of M to run.

$$\phi_j = w_0 + w_1 x_{1,j} + w_2 x_{2,j}$$

$$Pred_{j} = \begin{cases} 1, if \ \phi_{j} > 0 \\ -1, if \ \phi_{j} \le 0 \end{cases}$$

$$w_{0} = w_{0} + \eta (Actual_{j} - Pred_{j})$$

$$w_{1} = w_{1} + \eta (Actual_{j} - Pred_{j})x_{1,j}$$

$$w_{2} = w_{2} + \eta (Actual_{j} - Pred_{j})x_{2,j}$$

For every data point, we update the weight based on the prediction correctness, learning rate and feature values.

Calculate accuracy for the entire dataset to see whether the criteria has meet after each epoch.

For not linearly separable dataset, we need to use some accuracy threshold to stop the algorithm.

Perceptron R notebook: link

A Little Bit of History – Adaline



Adaptive Linear Neuron (Adaline)

* From Sebastian Python Machine Learning

Adaline R notebook: link

• We can use the error for the entire dataset as the loss function (i.e. SSE): $\sum_{i=1}^{N} (Actual_i -$ • Very similar to Perceptron and the only difference is that the error is calculated based on $\frac{1}{2}(Actual_i -$

 $w_{0} = w_{0} + \eta (Actual_{j} - \phi_{j})$ $w_{1} = w_{1} + \eta (Actual_{j} - \phi_{j}) x_{1,j}$ $w_{2} = w_{2} + \eta (Actual_{j} - \phi_{j}) x_{2,j}$

When calculating prediction accuracy, we still based on whether ϕ_j is larger than zero or not with the final weight learned from the data.

FEED FORWARD NEURAL NETWORK

Simple Feed Forward Neural Network (FFNN)



 $[x_1, x_2, x_3]$: Input features vector for one observation $[y, y_2]$: Actual output outcome for one observation $f_1(\cdot), f_2(\cdot)$: Activation functions, usually non-linear $L(y, \hat{y}(\boldsymbol{b}, \boldsymbol{x}))$: Loss function where \hat{y} the model forecast responses and y actual observed responses Total number of parameters (i.e. size of a NN): $(3+1)^*4 + (4+1)^*2 = 26$

Typical Loss Functions

Two-class binary responses

✓ Binary cross-entropy:

$$\sum_{i=1}^{n} (-y_i \log(p_i) - (1 - y_i) \log(1 - p_i))$$

where y_i is actual value of 1 or 0, p_i is the predicted probability of being 1, and N is the total number of observations in the training data.

Multiple-class responses

✓ Categorical cross-entropy for *M* classes:

$$\sum_{i=1}^{N} \left(-\sum_{j=1}^{M} y_{i,j} log(p_{i,j}) \right)$$

where $y_{i,j}$ is actual value of 1 or 0 for a class of j, $p_{i,j}$ is the predicted probability of being at class j and N is the total number of observations in the training data.

Continuous responses

- ✓ Mean square error
- Mean absolute error
- ✓ Mean absolute percentage error

From Slow Progress to Wide Adoption

1940s - 1980s, very slow progress due to:

- Computation hardware capacity limitation
- Number of observations with labeled results in the dataset
- Lack efficient algorithm to estimate the parameters in the model

1980s – 2010, a few applications in real word due to:

- Moore's Law + GPU
- Internet generated large labeled dataset
- **Efficient** algorithm for optimization (i.e. SGD + Backpropagation)
- Better activation functions (i.e. Relu)

2010-Now, very broad application in various areas:

- Near-human-level image classification
- Near-human-level handwriting transcription
- Much improved speech recognition
- Much improved machine translation

Now we know working neural network models usually contains many layers (i.e. the depth of the network is deep), and to achieve near-human-level accuracy the deep neural network need huge training dataset (for example millions of labeled pictures for image classification problem).

Optimization Methods

- Mini-batch Stochastic Gradient Descent (SGD)
 - □ Use a small segment of data (i.e. 128 or 256) to update the SGD parameters: $b = b - \alpha \nabla_b L(b, x^s, y^s)$ where α is the learning rate which is a hyper parameter that can be changed
 - Gradients are efficiently calculated using **backpropagation** method
 - When the entire dataset are used to updated the SGD, it is called one epoch and multiple epochs are needed to run to reach convergence
 - □ An updated version with 'momentum' for quick convergence:

$$v = \gamma v + \alpha \nabla_b L(b, x^s, y^s)$$
$$b = b - v$$

The optimal number for epoch is determined by when the model is not overfitted (i.e. validation accuracy reaches the best performance).



More Optimization Methods

- With Adaptive Learning Rates
 - Adagrad: Scale learning rate inversely proportional to the square root of the sum of all historical squared values of the gradient (stored for every weight)
 - RMSProp: Exponentially weighted moving average to accumulate gradients
 - AdaDelta: Uses sliding window to accumulate gradients
 - Adam (adaptive moments): momentum integrated, bias correction in decay
- A good summary: <u>http://ruder.io/optimizing-gradient-descent/</u>

Activation Functions

$$y = f_n(f_{n-1}(f_{n-2}(\dots f_1(x, b_1) \dots, b_{n-2}), b_{n-1}), b_n)$$

Gradient:
$$\frac{\partial y}{\partial x} = \frac{\partial y}{\partial f_n} \frac{\partial f_n}{\partial f_{n-1}} \frac{\partial f_{n-2}}{\partial f_{n-1}} \dots \frac{\partial f_1}{\partial x}$$

Intermediate layers

- Relu (i.e. rectified linear unit) is usually a good choice which has the following good properties: (1) fast computation; (2) non-linear; (3) reduced likelihood of the gradient to vanish; (4) Unconstrained response
- Sigmoid, studied in the past, not as good as Relu in deep learning, due to the gradient vanishing problem when there are many layers
- Last layer which connects to the output
 - □ Binary classification: sigmoid with binary cross entropy as loss function
 - Multiple class, single-label classification: softmax with categorical cross entropy for loss function

□ Continuous responses: identity function (i.e. y = x)

Rectified linear
unit (ReLU)^[10]
$$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \ge 0 \end{cases}$$
Softmax $f_i(\vec{x}) = \frac{e^{x_i}}{\sum_{j=1}^J e^{x_j}}$ for $i = 1, ..., J$ Logistic (a.k.a.
Sigmoid or Soft
step) $f(x) = \sigma(x) = \frac{1}{1 + e^{-x}}$ Softmax $f_i(\vec{x}) = \frac{e^{x_i}}{\sum_{j=1}^J e^{x_j}}$ for $i = 1, ..., J$

Deal With Overfitting

- Huge number of parameters, even with large amount of training data, there is a potential of overfitting
 - Overfitting due to size of the NN (i.e. total number of parameters)
 - Overfitting due to using the training data for too many epochs

□ Solution for overfitting due to NN size

- Dropout: randomly dropout some proportion (such as 0.3 or 0.5) of nodes at each layer, which is similar to random forest concept
- □ Using L1 or L2 regularization in the activation function at each layer

□ Solution for overfitting due to using too many epochs

- Run NN with large number of epochs to reach overfitting region through cross validation from training/validation vs. epoch curve
- Choose the model with number of epochs that have the minimum validation accuracy as the final NN model

Recap of A Few Key Concepts

- **Data**: Require large well-labeled dataset
- **Computation**: intensive matrix-matrix operation
- □ Structure of fully connected feedforward NN
 - □ Size of the NN: total number of parameters
 - Depth: total number of layers (this is where deep learning comes from)
 - Width of a particular layer: number of nodes (i.e. neurons) in that layer

Activation function Intermediate layers

Last layer connecting to outputs

Loss function

- □ Classification (i.e. categorical response)
- Regression (i.e. continuous response)

Optimization methods (SGD)

- Batch size
- Learning rate
- Epoch

Deal with overfitting

- Dropout
- Regularization (L1 or L2)

THE MNIST DATASET

.

-

MNIST Dataset

Originally created by NIST, then modified for machine leaning training purpose

- Contains 70000 handwritten digit images and the label of the digit in the image where 60000 images are the training dataset and the rest 10000 images are the testing dataset.
- **Census Bureau employees and American high school students wrote these digits**
- **Each image is 28x28 pixel in greyscale**
- Yann LeCun used convolutional network LeNet to achieve < 1% error rate at 1990s



.

-

IMDB Dataset

□ Raw data: 50,000 movie review text (X) and it's corresponding sentiment of positive (1, 50%) or negative (0, 50%) (Y).

Included in Keras package, can be easily loaded and preprocessed

□ Preprocessing includes:

- Set size of the vocabulary (i.e. N most frequently occurring words)
- Set length of the review by padding using '0' by default or truncating as we have to have same length for all reviews for modeling
- Any words not in the chosen vocabulary replaced by '2' by default
- Words are indexed by overall frequency in the chosen vocabulary

Once the dataset is preprocessed, we can apply encoding or embedding and then feed the data to FFNN or RNN

IMDB Dataset - Tokenization

Algorithm cannot deal with raw text and we have to convert text into numbers for machine learning methods.

	This movie is great ! Great movie ? Are you kidding me ! Not worth the money. Love it		
Raw Text			
Tokenize	[23, 55, 5, 78, 9] [78, 55, 8, 17, 12, 234, 33, 9, 14, 78, 32, 77, 4]	Suppose we cap unique	
Ļ	[65, 36] 	each one of these unique word is replaced	
Pad & Truncate	[23, 55, 5, 78, 9, 0, 0, 0, 0, 0] [78, 55, 8, 17, 12, 234, 33, 9, 14, 78] [65, 36, 0, 0, 0, 0, 0, 0, 0, 0]	with an integer. "2" will be used for any other words, and "0" will be used for padding.	

Now we have a typical data frame, each row is an observation, and each column is a feature. Here we have 10 columns by designing after the padding and truncating stage. We have converted raw text into categorical integers.

DEEP LEARNING MODELS ACROSS PLATFORMS

Open NN Exchange Format (ONNX)



