# CS395T: Structured Models for NLP Lecture 14: Neural Network Implementation



## Greg Durrett



- Project 2 due today
- Project 3 out today
  - of RNNs or CNNs
  - today's lecture

## Sentiment analysis using feedforward neural networks plus your choice

Project zip contains sample Tensorflow code demonstrating what's in





# **Recall: Feedforward NNs**







## Recall: Backpropagation







## Implementation details

## Training

## Word representations

## This Lecture

Implementation Details



- Computing gradients is hard!

$$x = x * x \longrightarrow (x, dx) = codegen$$

computation graph library

## Automatic differentiation: instrument code to keep track of derivatives

### (x \* x, 2 \* x \* dx)

In practice: need other operations, want more control -> use an external

# Computation Graphs



http://tmmse.xyz/content/images/2016/02/theano-computation-graph.png



- Define computation abstractly, in terms of symbols
- Can compute gradients of c with respect to (x, y, z) easily
- Useful abstraction: supports both CPU and GPU implementations
- Disadvantage: higher-level specification, so hard to control memory allocation and low-level implementation details

# Tensorflow





http://tmmse.xyz/content/images/2016/02/theano-computation-graph.png

- x = tf.placeholder("x")
- y = tf.placeholder("y")
- z = tf.placeholder("z")
- a = tf.add(x, y)
- b = tf.multiply(a, z)
- c = tf.add(b, a)
- with tf.Session() as sess:

![](_page_9_Picture_0.jpeg)

# Computation Graph: FFNN

- $$\begin{split} P(\mathbf{y}|\mathbf{x}) &= \operatorname{softmax}(Wg(Vf(\mathbf{x}))) \\ & \text{fx} = \operatorname{tf.placeholder(tf.float32, feat_vec_size)} \\ & \text{V} = \operatorname{tf.get\_variable("V", [hidden_size, feat_vec_size])} \\ & z = \operatorname{tf.sigmoid(tf.tensordot(V, fx, 1))} \\ & \text{W} = \operatorname{tf.get\_variable("W", [num_classes, hidden_size])} \\ & \text{probs} = \operatorname{tf.nn.softmax(tf.tensordot(W, z, 1))} \end{split}$$
- Placeholder: input to the system; variable: parameter to learn

![](_page_10_Picture_0.jpeg)

![](_page_10_Picture_1.jpeg)

$P(\mathbf{y} \mathbf{x}) = \operatorname{softmax}(Wg(Vf))$
<pre>fx = tf.placeholder(tf.floa</pre>
<pre>V = tf.get_variable("V", [h</pre>
<pre>z = tf.sigmoid(tf.tensordot</pre>
<pre>W = tf.get_variable("W", [n</pre>
<pre>probs = tf.nn.softmax(tf.te</pre>
<pre>label = tf.placeholder(tf.i</pre>
<pre>loss = tf.negative(tf.log(t</pre>
Tensorflow can compute gradient
Chartaut halpar mathada aviat lik

Shortcut helper methods exist like tf.nn.softmax cross entropy with logits

# **Computation Graph: FFNN**

- $(\mathbf{x})))$
- at32, feat vec size)
- idden size, feat vec size]) :(V, fx, 1))
- num classes, hidden size])
- ensordot(W, z, 1))
- .nt32, num classes)
- f.tensordot(probs, label, 1))
- ts for W and V based on loss

![](_page_11_Picture_0.jpeg)

## Define a computation graph

- For each epoch:
  - For each example:
    - Evaluate the training operator on the example
- Decode test set

## Training a Model

### Define an operator that updates the parameters based on an example

![](_page_12_Picture_0.jpeg)

- leads to better learning outcomes too
- fx = tf.placeholder(tf.float32, [batch\_size, feat vec size])
- V = tf.get variable("V", [hidden size, feat vec size])
- • •
- loss = [sum over losses from batch]

# Batching

Batching data gives speedups due to more efficient matrix operations,

Need to make the computation graph process a batch at the same time

z = tf.sigmoid(tf.tensordot(V, fx, [1,1])) # batch size x hidden size

![](_page_12_Picture_16.jpeg)

![](_page_13_Picture_0.jpeg)

Define a computation graph to process a batch of data

Define an operator that updates the parameters based on a **batch** 

For each epoch:

For each **batch**:

Evaluate the training operator on the **batch** 

Decode test set in **batches** 

# Batch Training a Model

Training Tips

![](_page_15_Picture_0.jpeg)

# Training Basics

- Basic formula: compute gradients on batch, use first-order opt. method
- How to initialize? How to regularize? What optimizer to use?
- This lecture: some practical tricks. Take deep learning or optimization courses to understand this further

![](_page_16_Picture_0.jpeg)

![](_page_16_Figure_1.jpeg)

![](_page_16_Figure_2.jpeg)

How do we initialize V and W? What consequences does this have?

# How does initialization affect learning?

![](_page_17_Picture_0.jpeg)

# How does initialization affect learning?

![](_page_17_Figure_3.jpeg)

- If cell activations are too large in absolute value, gradients are small too
- big values, can break down if everything is too negative

ReLU: larger dynamic range (all positive numbers), but can produce

![](_page_17_Picture_7.jpeg)

![](_page_17_Figure_8.jpeg)

# Initialization

![](_page_18_Picture_1.jpeg)

- 1) Can't use zeroes for parameters to produce hidden layers: all values in that hidden layer are always 0 and have gradients of 0, never change 2) Initialize too large and cells are saturated
- Can do random uniform / normal initialization with appropriate scale Glorot initializer:  $U = \sqrt{\frac{1}{\text{fan-in}}}$ 
  - Want variance of inputs and gradients for each layer to be the same
- Batch normalization (Ioffe and Szegedy, 2015): periodically shift+rescale each layer to have mean 0 and variance 1 over a batch (useful if net is deep)

$$\frac{6}{+ \text{ fan-out}}, +\sqrt{\frac{6}{\text{ fan-in + fan-out}}}$$

![](_page_18_Picture_7.jpeg)

![](_page_19_Picture_0.jpeg)

- Probabilistically zero out parts of the network during training to prevent overfitting, use whole network at test time
- Form of stochastic regularization
- Similar to benefits of ensembling: network needs to be robust to missing signals, so it has redundancy

![](_page_19_Picture_5.jpeg)

## Dropout

(a) Standard Neural Net

![](_page_19_Picture_8.jpeg)

(b) After applying dropout.

Srivastava et al. (2014)

![](_page_19_Figure_11.jpeg)

![](_page_20_Picture_0.jpeg)

## In tensorflow: implemented as an additional layer in a network

hidden dropped out = tf.nn.dropout(hidden, dropout keep prob)

Often use low dropout (keep a value with probability 0.8) at the input and moderate dropout (keep with probability 0.5) internally in feedforward networks (not in RNNs)

![](_page_21_Picture_0.jpeg)

## Adam (Kingma and Ba, ICLR 2015) is very widely used Adaptive step size like Adagrad, incorporates momentum

![](_page_21_Figure_3.jpeg)

# Optimizer

![](_page_21_Figure_5.jpeg)

![](_page_22_Picture_0.jpeg)

- Wilson et al. NIPS 2017: adaptive methods can actually perform badly at test time (Adam is in pink, SGD in black)
- Check dev set periodically, decrease learning rate if not making progress

![](_page_22_Figure_4.jpeg)

(e) Generative Parsing (Training Set)

# Optimizer

(f) Generative Parsing (Development Set)

![](_page_22_Figure_8.jpeg)

![](_page_22_Figure_9.jpeg)

![](_page_23_Picture_0.jpeg)

# Visualization with Tensorboard

## Visualize the computation graph and logs of the objective over time

![](_page_23_Figure_3.jpeg)

![](_page_23_Figure_4.jpeg)

![](_page_24_Picture_0.jpeg)

- Four elements of a structured machine learning method:
- Model: feedforward, RNNs, CNNs can be defined in a uniform framework
- Objective: many loss functions look similar, just changes the last layer of the neural network
- Inference: define the network, Tensorflow takes care of it (mostly...)
- Training: lots of choices for optimization/hyperparameters

# Structured Prediction

![](_page_24_Figure_8.jpeg)

![](_page_24_Picture_9.jpeg)

## Word Representations

![](_page_26_Picture_0.jpeg)

# Word Representations

- Continuous model <-> expects continuous semantics from input
- Neural networks work very well at continuous data, but words are discrete

![](_page_26_Picture_5.jpeg)

![](_page_27_Picture_0.jpeg)

Part-of-speech tagging with FFNNs

<u>?</u>?

Fed raises interest rates in order to ...

Word embeddings for each word form input

# Word Embeddings

![](_page_27_Figure_7.jpeg)

![](_page_27_Picture_8.jpeg)

![](_page_28_Picture_0.jpeg)

### the movie was great

 $\approx$ 

the movie was good

![](_page_28_Picture_6.jpeg)

# Word Embeddings

Want a vector space where similar words have similar embeddings

![](_page_29_Picture_0.jpeg)

# Word Representations

- Neural networks work very well at continuous data, but words are discrete
- Continuous model <-> expects continuous semantics from input
- "Can tell a word by the company it keeps" Firth 1957

	theore	sidentso		
(	president	the of		
1	president	the said		
	governor	the of		
	governor	the app		
	said	sources		
	said	president		
	reported	sources		

![](_page_29_Figure_7.jpeg)

[Finch and Chater 92, Shuetze 93, many others]

![](_page_29_Picture_9.jpeg)

![](_page_30_Picture_0.jpeg)

## Continuous Bag-of-Words

## Predict word from context

![](_page_30_Figure_3.jpeg)

![](_page_31_Picture_0.jpeg)

## Predict one word of context from word

![](_page_31_Figure_3.jpeg)

Another training example: bit -> the Parameters: d x |V| vectors, |V| x d output parameters (W)

## Skip-Gram

![](_page_31_Figure_6.jpeg)

## Mikolov et al. (2013)

![](_page_31_Picture_10.jpeg)

![](_page_32_Picture_0.jpeg)

- Problem: want to train on 1B+ words, the dog bit the man multiplying by |V| x d matrix for each is too expensive
- Solution: take (word, context) pairs and classify them as "real" or not. Create random negative examples by sampling
  - (bit, the) => +1(bit, a) => -1words in similar (bit, fish) => -1contexts select for similar c vectors
  - (*bit, dog*) => +1
- $P(\mathrm{pos}|w,c) = \frac{e^{w \cdot c}}{e^{w \cdot c} + 1}$

 $\mathbf{V} = d \times |\mathbf{V}|$  vectors,  $d \times |\mathbf{V}|$  context vectors (same # of params as before) Mikolov et al. (2013)

## Skip-Gram with Negative Sampling

![](_page_32_Picture_9.jpeg)

![](_page_33_Picture_0.jpeg)

## (king - man) + woman = queen

king + (woman - man) = queen

- Why would this be?
- woman man captures the difference in the contexts that these occur in
- Dominant change: more "he" with man and "she" with woman — similar to difference between king and queen

## **Regularities in Vector Space**

![](_page_33_Figure_8.jpeg)

![](_page_34_Picture_0.jpeg)

## Word co-occurrences are what matter directly

Probability and Ratio	k = solid	k = gas	k = water	k = fashion
P(k ice)	1.9 × 10 <sup>-4</sup>	6.6 × 10 <sup>-5</sup>	$3.0 \times 10^{-3}$	$1.7 \times 10^{-5}$
P(k steam)	$2.2 \times 10^{-5}$	$7.8 \times 10^{-4}$	$2.2 \times 10^{-3}$	$1.8 \times 10^{-5}$
P(k ice)/P(k steam)	8.9	$8.5 \times 10^{-2}$	1.36	0.96

Weighted least-squares problem to directly predict word co-occurrence matrix (like matrix factorization)

## Pennington et al. (2014)

![](_page_34_Picture_7.jpeg)

# Using Word Embeddings

![](_page_35_Picture_1.jpeg)

# Indexed sentence of length sent len, e.g.: [12, 36, 47, 8] input words = tf.placeholder(tf.int32, [sent len]) encoder = tf.get variable("embed", [voc size, embedding\_size]) # embedded input words: sent len x embedding size tensor

- Approach 1: learn embeddings as parameters from your data
- Approach 2: initialize using GloVe/CBOW/SGNS, keep fixed
  - Faster because no need to update these parameters
- Approach 3: initialize using GloVe/CBOW/SGNS, fine-tune
  - Typically works best

embedded input words = tf.nn.embedding lookup(encoder, input words)

![](_page_35_Picture_9.jpeg)

![](_page_36_Picture_0.jpeg)

- Lots to tune with neural networks
  - Training: optimizer, initializer, regularization (dropout), ...
  - Hyperparameters: dimensionality of word embeddings, layers, ...
- Word vectors: various choices of pre-trained vectors work well as initializers
- Next time: RNNs / LSTMs / GRUs

## Takeaways