

# CS388: Natural Language Processing

## Lecture 6: Neural Networks



Greg Durrett



## Administrivia

- ▶ Mini 1 graded, posted on Canvas
  - ▶ Xi Ye (88.0 F1), Quang Duong (87.3 F1), Uday Kusupati (87.2 F1)
  - 6 students in the 86 range, rest are 85 or below
  - ▶ Test F1s << dev F1
  - ▶ Changing thresholds / imbalanced classification
  - ▶ POS/chunk features
  - ▶ Someone got 86.3 with only 7 features total, classifier is a dictionary
- ▶ Project 1 due in 9 days
  - ▶ Small bug fixed in BadNerModel (no impact on the code you write)



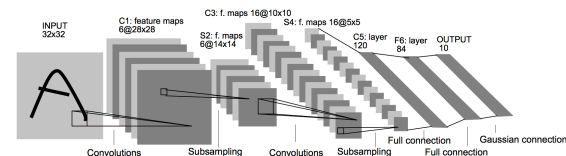
## This Lecture

- ▶ Beam search: in a few lectures
- ▶ Neural network history
- ▶ Neural network basics
- ▶ Feedforward neural networks + backpropagation
- ▶ Applications
- ▶ Implementing neural networks (if time)

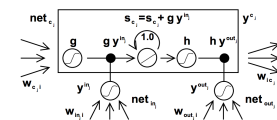


## History: NN “dark ages”

- ▶ Convnets: applied to MNIST by LeCun in 1998



- ▶ LSTMs: Hochreiter and Schmidhuber (1997)

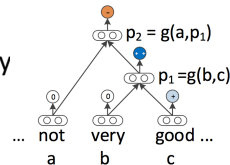
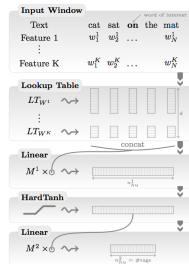


- ▶ Henderson (2003): neural shift-reduce parser, not SOTA



## 2008-2013: A glimmer of light...

- Collobert and Weston 2011: “NLP (almost) from scratch”
  - Feedforward neural nets induce features for sequential CRFs (“neural CRF”)
  - 2008 version was marred by bad experiments, claimed SOTA but wasn’t, 2011 version tied SOTA
- Krizhevsky et al. (2012): AlexNet for vision
- Socher 2011-2014: tree-structured RNNs working okay



## 2014: Stuff starts working

- Kim (2014) + Kalchbrenner et al. (2014): sentence classification / sentiment (convnets work for NLP?)
- Sutskever et al. + Bahdanau et al.: seq2seq for neural MT (LSTMs work for NLP?)
- Chen and Manning transition-based dependency parser (even feedforward networks work well for NLP?)
- 2015: explosion of neural nets for everything under the sun



## Why didn't they work before?

- Datasets too small:** for MT, not really better until you have 1M+ parallel sentences (and really need a lot more)
- Optimization not well understood:** good initialization, per-feature scaling + momentum (Adagrad / Adadelta / Adam) work best out-of-the-box
  - Regularization:** dropout is pretty helpful
  - Computers not big enough:** can't run for enough iterations
- Inputs:** need word representations to have the right continuous semantics

## Neural Net Basics



## Neural Networks

- ▶ Linear classification:  $\operatorname{argmax}_y w^\top f(x, y)$
- ▶ How can we do nonlinear classification? Kernels are too slow...
- ▶ Want to learn intermediate conjunctive features of the input

*the movie was **not** all that **good***

I[contains *not* & contains *good*]

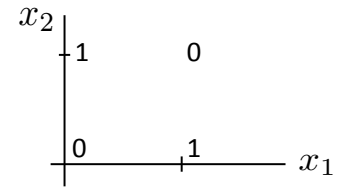


## Neural Networks: XOR

- ▶ Let's see how we can use neural nets to learn a simple nonlinear function

▶ Inputs  $x_1, x_2$   
(generally  $\mathbf{x} = (x_1, \dots, x_m)$ )

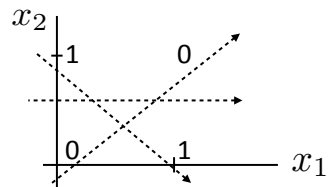
▶ Output  $y$   
(generally  $\mathbf{y} = (y_1, \dots, y_n)$ )



$x_1$	$x_2$	$y = x_1 \text{ XOR } x_2$
0	0	0
0	1	1
1	0	1
1	1	0



## Neural Networks: XOR



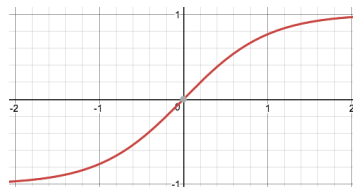
$$y = a_1 x_1 + a_2 x_2$$

✗

$$y = a_1 x_1 + a_2 x_2 + a_3 \tanh(x_1 + x_2)$$

“or” ✓

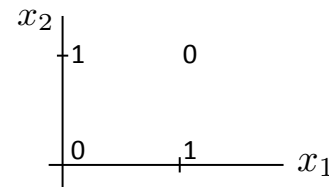
(looks like action potential in neuron)



$x_1$	$x_2$	$x_1 \text{ XOR } x_2$
0	0	0
0	1	1
1	0	1
1	1	0



## Neural Networks: XOR



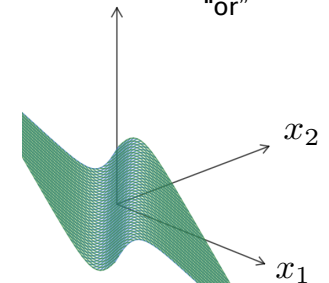
$$y = a_1 x_1 + a_2 x_2$$

✗

$$y = a_1 x_1 + a_2 x_2 + a_3 \tanh(x_1 + x_2)$$

“or” ✓

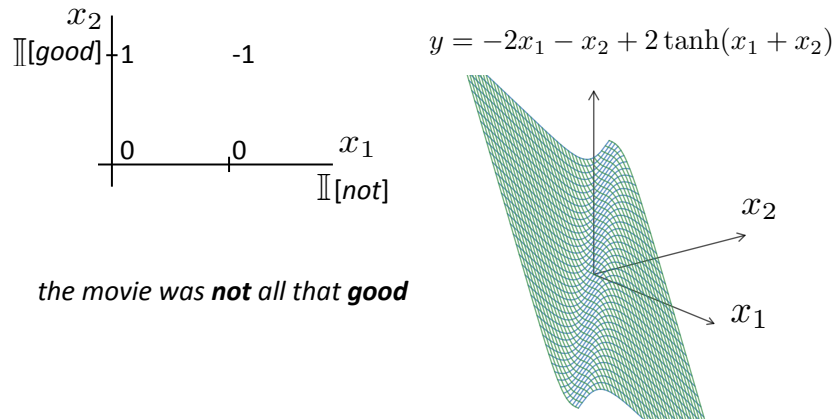
$$y = -x_1 - x_2 + 2 \tanh(x_1 + x_2)$$



$x_1$	$x_2$	$x_1 \text{ XOR } x_2$
0	0	0
0	1	1
1	0	1
1	1	0



## Neural Networks: XOR



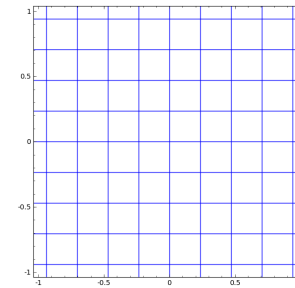
## Neural Networks

Linear model:  $y = \mathbf{w} \cdot \mathbf{x} + b$

$$y = g(\mathbf{w} \cdot \mathbf{x} + b)$$

$$\mathbf{y} = g(\mathbf{W}\mathbf{x} + \mathbf{b})$$

Nonlinear transformation    Warp space    Shift

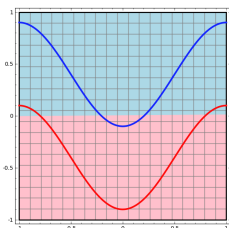


Taken from <http://colah.github.io/posts/2014-03-NN-Manifolds-Topology/>

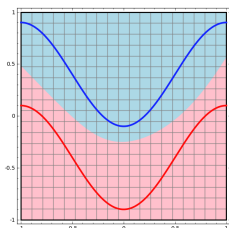


## Neural Networks

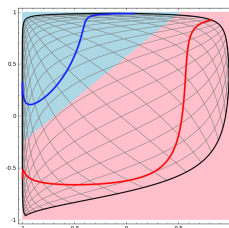
Linear classifier



Neural network



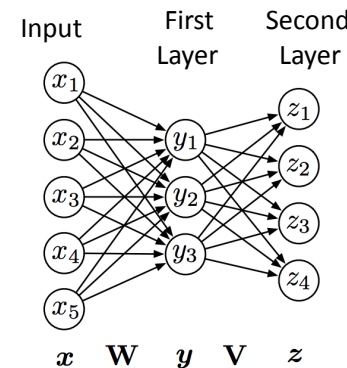
...possible because we transformed the space!



Taken from <http://colah.github.io/posts/2014-03-NN-Manifolds-Topology/>



## Deep Neural Networks



$$\mathbf{y} = g(\mathbf{W}\mathbf{x} + \mathbf{b})$$

$$\mathbf{z} = g(\mathbf{V}\mathbf{y} + \mathbf{c})$$

$$\mathbf{z} = g(\mathbf{V}g(\mathbf{W}\mathbf{x} + \mathbf{b}) + \mathbf{c})$$

output of first layer

“Feedforward” computation (not recurrent)

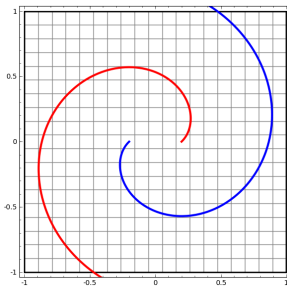
Check: what happens if no nonlinearity?  
More powerful than basic linear models?

$$\mathbf{z} = \mathbf{V}(\mathbf{W}\mathbf{x} + \mathbf{b}) + \mathbf{c}$$

Adopted from Chris Dyer



## Deep Neural Networks



Taken from <http://colah.github.io/posts/2014-03-NN-Manifolds-Topology/>

## Feedforward Networks, Backpropagation



## Logistic Regression with NNs

$$P(y|\mathbf{x}) = \frac{\exp(w^\top f(\mathbf{x}, y))}{\sum_{y'} \exp(w^\top f(\mathbf{x}, y'))}$$

▶ Single scalar probability

$$P(\mathbf{y}|\mathbf{x}) = \text{softmax}([w^\top f(\mathbf{x}, y)]_{y \in \mathcal{Y}})$$

▶ Compute scores for all possible labels at once (returns vector)

$$\text{softmax}(p)_i = \frac{\exp(p_i)}{\sum_{i'} \exp(p_{i'})}$$

▶ softmax: exps and normalizes a given vector

$$P(\mathbf{y}|\mathbf{x}) = \text{softmax}(W f(\mathbf{x}))$$

▶ Weight vector per class;  
W is [num classes x num feats]

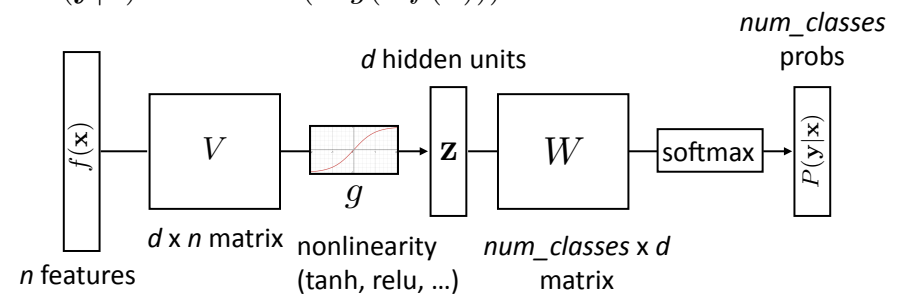
$$P(\mathbf{y}|\mathbf{x}) = \text{softmax}(W g(V f(\mathbf{x})))$$

▶ Now one hidden layer



## Neural Networks for Classification

$$P(\mathbf{y}|\mathbf{x}) = \text{softmax}(W g(V f(\mathbf{x})))$$





## Training Neural Networks

$$P(\mathbf{y}|\mathbf{x}) = \text{softmax}(W\mathbf{z}) \quad \mathbf{z} = g(Vf(\mathbf{x}))$$

- ▶ Maximize log likelihood of training data

$$\mathcal{L}(\mathbf{x}, i^*) = \log P(y = i^*|\mathbf{x}) = \log (\text{softmax}(W\mathbf{z}) \cdot e_{i^*})$$

- ▶  $i^*$ : index of the gold label
- ▶  $e_i$ : 1 in the  $i$ th row, zero elsewhere. Dot by this = select  $i$ th index

$$\mathcal{L}(\mathbf{x}, i^*) = W\mathbf{z} \cdot e_{i^*} - \log \sum_j \exp(W\mathbf{z}) \cdot e_j$$



## Computing Gradients

$$\mathcal{L}(\mathbf{x}, i^*) = W\mathbf{z} \cdot e_{i^*} - \log \sum_j \exp(W\mathbf{z}) \cdot e_j$$

- ▶ Gradient with respect to  $W$

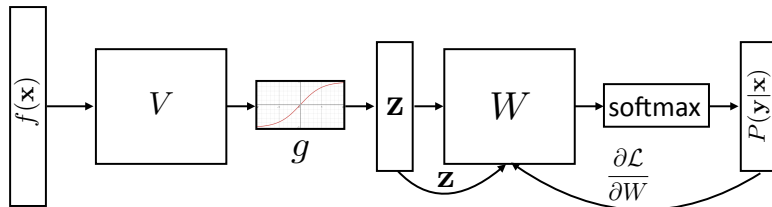
$$\frac{\partial}{\partial W_{ij}} \mathcal{L}(\mathbf{x}, i^*) = \begin{cases} \mathbf{z}_j - P(y = i|\mathbf{x})\mathbf{z}_j & \text{if } i = i^* \\ -P(y = i|\mathbf{x})\mathbf{z}_j & \text{otherwise} \end{cases}$$

- ▶ Looks like logistic regression with  $\mathbf{z}$  as the features!



## Neural Networks for Classification

$$P(\mathbf{y}|\mathbf{x}) = \text{softmax}(Wg(Vf(\mathbf{x})))$$



## Computing Gradients: Backpropagation

$$\mathcal{L}(\mathbf{x}, i^*) = W\mathbf{z} \cdot e_{i^*} - \log \sum_j \exp(W\mathbf{z}) \cdot e_j \quad \mathbf{z} = g(Vf(\mathbf{x}))$$

Activations at hidden layer

- ▶ Gradient with respect to  $V$ : apply the chain rule

$$\frac{\partial \mathcal{L}(\mathbf{x}, i^*)}{\partial V_{ij}} = \frac{\partial \mathcal{L}(\mathbf{x}, i^*)}{\partial \mathbf{z}} \frac{\partial \mathbf{z}}{\partial V_{ij}}$$

[some math...]

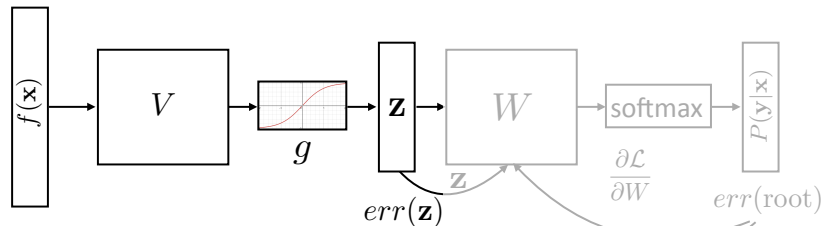
$$\text{err}(\text{root}) = e_{i^*} - P(\mathbf{y}|\mathbf{x}) \quad \frac{\partial \mathcal{L}(\mathbf{x}, i^*)}{\partial \mathbf{z}} = \text{err}(\mathbf{z}) = W^\top \text{err}(\text{root})$$

dim = m                      dim = d



## Backpropagation: Picture

$$P(\mathbf{y}|\mathbf{x}) = \text{softmax}(Wg(Vf(\mathbf{x})))$$



- ▶ Can forget everything after  $\mathbf{z}$ , treat it as the output and keep backpropping



## Backpropagation: Takeaways

- ▶ Gradients of output weights  $W$  are easy to compute — looks like logistic regression with hidden layer  $\mathbf{z}$  as feature vector
- ▶ Can compute derivative of loss with respect to  $\mathbf{z}$  to form an “error signal” for backpropagation
- ▶ Easy to update parameters based on “error signal” from next layer, keep pushing error signal back as backpropagation
- ▶ Need to remember the values from the forward computation

## Applications



## NLP with Feedforward Networks

- ▶ Part-of-speech tagging with FFNNs

??

Fed raises interest rates in order to ...

previous word

- ▶ Word embeddings for each word form input

- ▶ ~1000 features here — smaller feature vector than in sparse models, but every feature fires on every example

curr word

next word

- ▶ Weight matrix learns position-dependent processing of the words

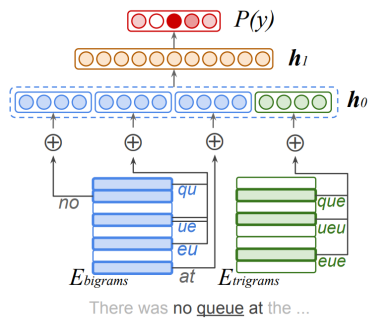
other words, feats, etc.



Botha et al. (2017)



## NLP with Feedforward Networks



- Hidden layer mixes these different signals and learns feature conjunctions

Botha et al. (2017)



## NLP with Feedforward Networks

- Multilingual tagging results:

Model	Acc.	Wts.	MB	Ops.
Gillick et al. (2016)	95.06	900k	-	6.63m
Small FF	94.76	241k	0.6	0.27m
+Clusters	95.56	261k	1.0	0.31m
$\frac{1}{2}$ Dim.	95.39	143k	0.7	0.18m

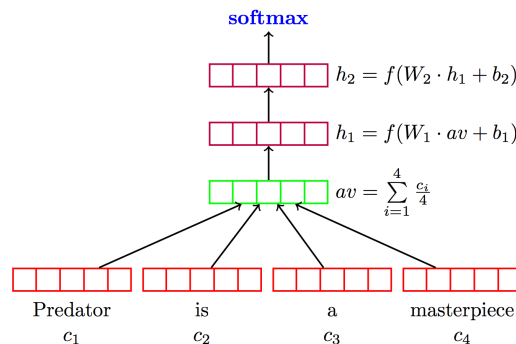
- Gillick used LSTMs; this is smaller, faster, and better

Botha et al. (2017)



## Sentiment Analysis

- Deep Averaging Networks: feedforward neural network on average of word embeddings from input



Iyyer et al. (2015)



## Sentiment Analysis

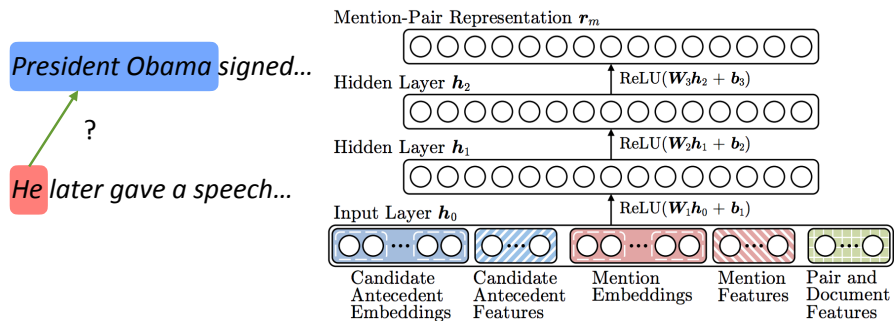
	Model	RT	SST fine	SST bin	IMDB	Time (s)	
Bag-of-words	DAN-ROOT	—	46.9	85.7	—	31	Iyyer et al. (2015)
	DAN-RAND	77.3	45.4	83.2	88.8	136	
	DAN	80.3	47.7	86.3	89.4	136	
	NBOW-RAND	76.2	42.3	81.4	88.9	91	Wang and Manning (2012)
	NBOW	79.0	43.6	83.6	89.0	91	
Tree RNNs / CNNs / LSTMs	BiNB	—	41.9	83.1	—	—	
	NBSVM-bi	79.4	—	—	91.2	—	
	RecNN*	77.7	43.2	82.4	—	—	Kim (2014)
	RecNTN*	—	45.7	85.4	—	—	
	DRecNN	—	49.8	86.6	—	431	
	TreeLSTM	—	50.6	86.9	—	—	
	DCNN*	—	48.5	86.9	89.4	—	
	PVEC*	—	48.7	87.8	92.6	—	
	CNN-MC	81.1	47.4	88.1	—	2,452	
	WRRBM*	—	—	—	89.2	—	





## Coreference Resolution

- Feedforward networks identify coreference arcs



Clark and Manning (2015), Wiseman et al. (2015)

## Implementation Details



## Computation Graphs

- Computing gradients is hard!
  - Automatic differentiation: instrument code to keep track of derivatives
- Example:  $y = x * x$   $\xrightarrow{\text{codegen}}$   $(y, dy) = (x * x, 2 * x * dx)$
- Computation is now something we need to reason about symbolically
  - Use a library like Pytorch or Tensorflow. This class: Pytorch



## Computation Graphs in Pytorch

- Define forward pass for  $P(y|x) = \text{softmax}(Wg(Vf(x)))$

```
class FFNN(nn.Module):
    def __init__(self, inp, hid, out):
        super(FFNN, self).__init__()
        self.V = nn.Linear(inp, hid)
        self.g = nn.Tanh()
        self.W = nn.Linear(hid, out)
        self.softmax = nn.Softmax(dim=0)

    def forward(self, x):
        return self.softmax(self.W(self.g(self.V(x))))
```



## Computation Graphs in Pytorch

$P(\mathbf{y}|\mathbf{x}) = \text{softmax}(Wg(Vf(\mathbf{x})))$

$\mathbf{e}_i$ : one-hot vector of the label  
 (e.g., [0, 1, 0])

$\swarrow$

```

ffnn = FFNN()
def make_update(input, gold_label):
    ffnn.zero_grad() # clear gradient variables
    probs = ffnn.forward(input)
    loss = torch.neg(torch.log(probs)).dot(gold_label)
    loss.backward()
    optimizer.step()
  
```



## Training a Model

Define a computation graph

For each epoch:

For each batch of data:

Compute loss on batch

Autograd to compute gradients and take step

Decode test set



## Batching

- ▶ Batching data gives speedups due to more efficient matrix operations
- ▶ Need to make the computation graph process a batch at the same time

```

# input is [batch_size, num_feats]
# gold_label is [batch_size, num_classes]
def make_update(input, gold_label)
    ...
    probs = ffnn.forward(input) # [batch_size, num_classes]
    loss = torch.sum(torch.neg(torch.log(probs)).dot(gold_label))
    ...
  
```

- ▶ Batch sizes from 1-100 often work well



## Next Time

- ▶ More implementation details: practical training techniques
- ▶ Word representations / word vectors
- ▶ word2vec, GloVe