

## CS 6120/CS 4120: Natural Language Processing

Instructor: Prof. Lu Wang  
 Northeastern University  
 Webpage: [www.ccs.neu.edu/home/luwang](http://www.ccs.neu.edu/home/luwang)

1

## Outline

- ➔ • Maximum Entropy
- Feedforward Neural Networks
- Recurrent Neural Networks

2

## Maximum Entropy (MaxEnt)

- Or logistic regression

3

## Features

- In these slides and most MaxEnt work: *features (or feature functions)  $f$*  are elementary pieces of evidence that link aspects of what we observe  $d$  with a category  $c$  that we want to predict
- A feature is a function with a **bounded** real value:  $f: C \times D \rightarrow \mathbb{R}$

4

## Example Task: Named Entity Type

LOCATION      LOCATION      DRUG      PERSON  
*in Arcadia*    *in Québec*    *taking Zantac*    *saw Sue*

5

## Example features

- $f_1(c, d) \equiv [c = \text{LOCATION} \wedge w_{-1} = \text{"in"} \wedge \text{isCapitalized}(w)]$
- $f_2(c, d) \equiv [c = \text{LOCATION} \wedge \text{hasAccentedLatinChar}(w)]$
- $f_3(c, d) \equiv [c = \text{DRUG} \wedge \text{ends}(w, \text{"c"})]$

LOCATION      LOCATION      DRUG      PERSON  
*in Arcadia*    *in Québec*    *taking Zantac*    *saw Sue*

- Models will assign to each feature a *weight*:
  - A positive weight votes that this configuration is likely correct
  - A negative weight votes that this configuration is likely incorrect

6

### Example features

- $f_1(c, d) \equiv [c = \text{LOCATION} \wedge w_1 = \text{"in"} \wedge \text{isCapitalized}(w)] \rightarrow \text{weight } 1.8$
- $f_2(c, d) \equiv [c = \text{LOCATION} \wedge \text{hasAccentedLatinChar}(w)] \rightarrow \text{weight } -0.6$
- $f_3(c, d) \equiv [c = \text{DRUG} \wedge \text{ends}(w, \text{"c"})] \rightarrow \text{weight } 0.3$

• Weights will be learned by training on a labeled dataset

7

### More about feature functions:

an indicator function – a yes/no boolean matching function – of properties of the input and a particular class

$$f_i(c, d) \equiv [\Phi(d) \wedge c = c_j] \quad [\text{Value is 0 or 1}]$$

8

### Feature-Based Models

• The decision about a data point is based only on the **features** active at that point.

<p><b>Data</b> BUSINESS: Stocks hit a yearly low ...</p> <p><b>Label: BUSINESS</b> Features {..., stocks, hit, a, yearly, low, ...}</p> <p>Text Classification</p>	<p><b>Data</b> ... to restructure bank: MONEY debt.</p> <p><b>Label: MONEY</b> Features {..., w<sub>1</sub>=restructure, w<sub>2</sub>=debt, L=12, ...}</p> <p>Word Sense Disambiguation</p>	<p><b>Data</b> DT JJ NN ... The previous fall ...</p> <p><b>Label: NN</b> Features {w<sub>1</sub>=fall, L<sub>1</sub>=JJ, w<sub>2</sub>=previous}</p> <p>POS Tagging</p>
--	--	--

9

### Feature-Based Linear Classifiers

- Linear classifiers at classification time:
  - Linear function from feature sets  $\{f_i\}$  to classes  $\{c\}$ .
  - Assign a weight  $\lambda_i$  to each feature  $f_i$ .
  - We consider each class for sample  $d$
  - For a pair  $(c, d)$ , features vote with their weights:
    - $\text{vote}(c) = \sum \lambda_i f_i(c, d)$

PERSON in Québec      LOCATION in Québec      DRUG in Québec

• Choose the class  $c$  which maximizes  $\sum \lambda_i f_i(c, d)$

10

### Maximum Entropy:

• Make a probabilistic model from the linear combination  $\sum \lambda_i f_i(c, d)$

$$P(c | d, \lambda) = \frac{\exp \sum \lambda_i f_i(c, d)}{\sum_{c'} \exp \sum \lambda_i f_i(c', d)}$$

Makes votes positive      Normalizes votes

11

### Feature-Based Linear Classifiers

- $f_1(c, d) \equiv [c = \text{LOCATION} \wedge w_1 = \text{"in"} \wedge \text{isCapitalized}(w)] \rightarrow \text{weight } 1.8$
- $f_2(c, d) \equiv [c = \text{LOCATION} \wedge \text{hasAccentedLatinChar}(w)] \rightarrow \text{weight } -0.6$
- $f_3(c, d) \equiv [c = \text{DRUG} \wedge \text{ends}(w, \text{"c"})] \rightarrow \text{weight } 0.3$

12

$f_1(c, d) \equiv [c = \text{LOCATION} \wedge w_{-1} = \text{"in"} \wedge \text{isCapitalized}(w)] \rightarrow \text{weight } 1.8$   
 $f_2(c, d) \equiv [c = \text{LOCATION} \wedge \text{hasAccentedLatinChar}(w)] \rightarrow \text{weight } -0.6$   
 $f_3(c, d) \equiv [c = \text{DRUG} \wedge \text{ends}(w, "c")] \rightarrow \text{weight } 0.3$

- Maximum Entropy:
  - Make a probabilistic model from the linear combination  $\sum \lambda_i f_i(c, d)$

$$P(c | d, \lambda) = \frac{\exp \sum_i \lambda_i f_i(c, d)}{\sum_{c'} \exp \sum_i \lambda_i f_i(c', d)}$$

Makes votes positive  
Normalizes votes

13

$f_1(c, d) \equiv [c = \text{LOCATION} \wedge w_{-1} = \text{"in"} \wedge \text{isCapitalized}(w)] \rightarrow \text{weight } 1.8$   
 $f_2(c, d) \equiv [c = \text{LOCATION} \wedge \text{hasAccentedLatinChar}(w)] \rightarrow \text{weight } -0.6$   
 $f_3(c, d) \equiv [c = \text{DRUG} \wedge \text{ends}(w, "c")] \rightarrow \text{weight } 0.3$

- Maximum Entropy:
  - Make a probabilistic model from the linear combination  $\sum \lambda_i f_i(c, d)$

$$P(c | d, \lambda) = \frac{\exp \sum_i \lambda_i f_i(c, d)}{\sum_{c'} \exp \sum_i \lambda_i f_i(c', d)}$$

Makes votes positive  
Normalizes votes

- $P(\text{LOCATION} | \text{in Québec}) = e^{1.8} e^{-0.6} / (e^{1.8} e^{-0.6} + e^{0.3} + e^0) = 0.586$
- $P(\text{DRUG} | \text{in Québec}) = e^{0.3} / (e^{1.8} e^{-0.6} + e^{0.3} + e^0) = 0.238$
- $P(\text{PERSON} | \text{in Québec}) = e^0 / (e^{1.8} e^{-0.6} + e^{0.3} + e^0) = 0.176$

The weights are the parameters of the probability model, combined via a "soft max" function

14

### Feature-Based Linear Classifiers

- Given this model form, we will choose parameters  $\{\lambda_i\}$  that *maximize the conditional likelihood* of the data according to this model.
- Parameter learning is omitted and not required for this course, but is often discussed in a machine learning class.
  - E.g. gradient descent for parameter learning

15

### Outline

- Maximum Entropy
- ➡ Feedforward Neural Networks
- Recurrent Neural Networks

16

### Neural Network Learning

- Learning approach based on modeling adaptation in biological neural systems.
- Perceptron:** Initial algorithm for learning simple neural networks (single layer) developed in the 1950's.
- Backpropagation:** More complex algorithm for learning multi-layer neural networks developed in the 1980's. (not required for this class)

17

### ARTIFICIAL NEURON

**Topics:** connection weights, bias, activation function

- Neuron pre-activation (or input activation):  
 $a(\mathbf{x}) = b + \sum_i w_i x_i = b + \mathbf{w}^T \mathbf{x}$
- Neuron (output) activation  
 $h(\mathbf{x}) = g(a(\mathbf{x})) = g(b + \sum_i w_i x_i)$
- $\mathbf{w}$  are the connection weights
- $b$  is the neuron bias
- $g(\cdot)$  is called the activation function

18

### ARTIFICIAL NEURON

**Topics:** connection weights, bias, activation function

range determined by  $g(\cdot)$

bias  $b$  only changes the position of the ruff

(from Pascal Vincent's slides)

19

### ACTIVATION FUNCTION

**Topics:** linear activation function

- Performs no input squashing
- Not very interesting..

$g(a) = a$

20

### ACTIVATION FUNCTION

**Topics:** sigmoid activation function

- Squashes the neuron's pre-activation between 0 and 1
- Always positive
- Bounded
- Strictly increasing

$g(a) = \text{sigm}(a) = \frac{1}{1 + \exp(-a)}$

21

### ACTIVATION FUNCTION

**Topics:** hyperbolic tangent ("tanh") activation function

- Squashes the neuron's pre-activation between -1 and 1
- Can be positive or negative
- Bounded
- Strictly increasing

$g(a) = \text{tanh}(a) = \frac{\exp(a) - \exp(-a)}{\exp(a) + \exp(-a)} = \frac{\exp(2a) - 1}{\exp(2a) + 1}$

22

### ACTIVATION FUNCTION

**Topics:** rectified linear activation function

- Bounded below by 0 (always non-negative)
- Not upper bounded
- Strictly increasing
- Tends to give neurons with sparse activities

$g(a) = \text{reclin}(a) = \max(0, a)$

23

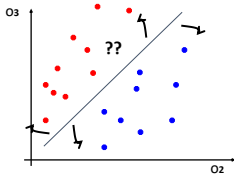
```

class Neuron(object):
    # ...
    def forward(inputs):
        """ assume inputs and weights are 1-D numpy arrays and bias is a number """
        cell_body_sum = np.sum(inputs * self.weights) + self.bias
        firing_rate = 1.0 / (1.0 + math.exp(-cell_body_sum)) # sigmoid activation function
        return firing_rate
    
```

24

### Linear Separator

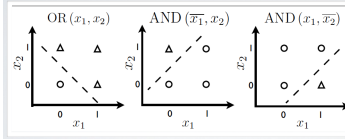
- Since one-layer neuron (aka perceptron) uses linear threshold function, it is searching for a linear separator that discriminates the classes.



25

### ARTIFICIAL NEURON

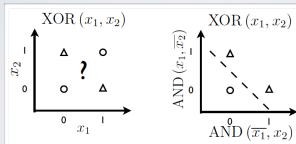
- **Topics:** capacity of single neuron
- Can solve linearly separable problems



26

### ARTIFICIAL NEURON

- **Topics:** capacity of single neuron
- Can't solve non linearly separable problems...

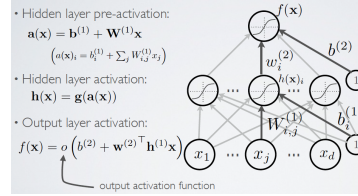


• ... unless the input is transformed in a better representation

27

### NEURAL NETWORK

- **Topics:** single hidden layer neural network



- Hidden layer pre-activation:  
 $\mathbf{a}(x) = \mathbf{b}^{(1)} + \mathbf{W}^{(1)}\mathbf{x}$   
 $(a(x)_i = b_i^{(1)} + \sum_j W_{ij}^{(1)}x_j)$
- Hidden layer activation:  
 $\mathbf{h}(x) = \mathbf{g}(\mathbf{a}(x))$
- Output layer activation:  
 $f(x) = \sigma(\mathbf{b}^{(2)} + \mathbf{w}^{(2)\top}\mathbf{h}^{(1)}\mathbf{x})$

28

### NEURAL NETWORK

- **Topics:** softmax activation function
- For multi-class classification:
  - we need multiple outputs (1 output per class)
  - we would like to estimate the conditional probability  $p(y = c|\mathbf{x})$

• We use the softmax activation function at the output:

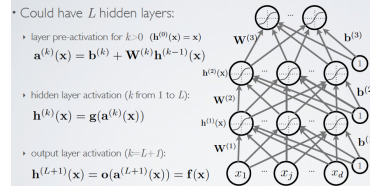
$$\mathbf{o}(\mathbf{a}) = \text{softmax}(\mathbf{a}) = \left[ \frac{\exp(a_1)}{\sum_c \exp(a_c)}, \dots, \frac{\exp(a_c)}{\sum_c \exp(a_c)} \right]^T$$

- strictly positive
- sums to one
- Predicted class is the one with highest estimated probability

29

### NEURAL NETWORK

- **Topics:** multilayer neural network



- Could have  $L$  hidden layers:
- layer pre-activation for  $k > 0$  ( $\mathbf{h}^{(0)}(\mathbf{x}) = \mathbf{x}$ ):  
 $\mathbf{a}^{(k)}(\mathbf{x}) = \mathbf{b}^{(k)} + \mathbf{W}^{(k)}\mathbf{h}^{(k-1)}(\mathbf{x})$
- hidden layer activation ( $k$  from 1 to  $L$ ):  
 $\mathbf{h}^{(k)}(\mathbf{x}) = \mathbf{g}(\mathbf{a}^{(k)}(\mathbf{x}))$
- output layer activation ( $k=L+1$ ):  
 $\mathbf{h}^{(L+1)}(\mathbf{x}) = \mathbf{o}(\mathbf{a}^{(L+1)}(\mathbf{x})) = \mathbf{f}(\mathbf{x})$

30

```

# forward-pass of a 3-layer neural network:
f = lambda x: 1.0/(1.0 + np.exp(-x)) # activation function (use sigmoid)
x = np.random.randn(3, 1) # random input vector of three numbers (3x1)
h1 = f(np.dot(W1, x) + b1) # calculate first hidden layer activations (4x1)
h2 = f(np.dot(W2, h1) + b2) # calculate second hidden layer activations (4x1)
out = np.dot(W3, h2) + b3 # output neuron (1x1)
    
```

31

### CAPACITY OF NEURAL NETWORK

Topics: single hidden layer neural network

(from Pascal Vincent's slides)

32

### CAPACITY OF NEURAL NETWORK

Topics: single hidden layer neural network

(from Pascal Vincent's slides)

33

### CAPACITY OF NEURAL NETWORK

Topics: single hidden layer neural network

(from Pascal Vincent's slides)

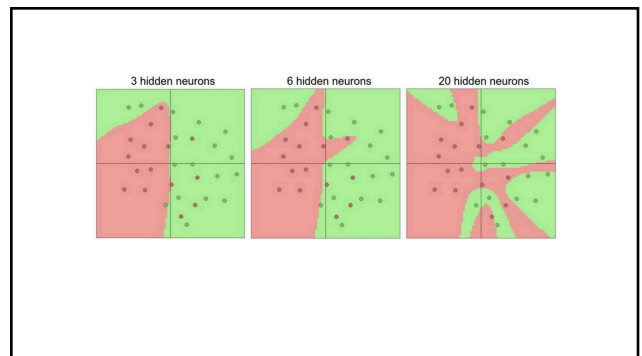
34

### CAPACITY OF NEURAL NETWORK

Topics: universal approximation

- Universal approximation theorem (Kornik, 1991):
  - "a single hidden layer neural network with a linear output unit can approximate any continuous function arbitrarily well, given enough hidden units"
- The result applies for sigmoid, tanh and many other hidden layer activation functions
- This is a good result, but it doesn't mean there is a learning algorithm that can find the necessary parameter values!

35



36

### How to train a neural network? (Not covered in this course, only for reference)

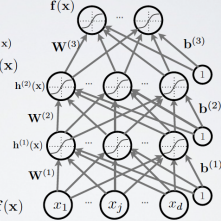
**Topics:** multilayer neural network

- Could have  $L$  hidden layers:

- layer input activation for  $k > 0$  ( $\mathbf{h}^{(0)}(\mathbf{x}) = \mathbf{x}$ )  
 $\mathbf{a}^{(k)}(\mathbf{x}) = \mathbf{b}^{(k)} + \mathbf{W}^{(k)}\mathbf{h}^{(k-1)}(\mathbf{x})$

- hidden layer activation ( $k$  from 1 to  $L$ ):  
 $\mathbf{h}^{(k)}(\mathbf{x}) = \mathbf{g}(\mathbf{a}^{(k)}(\mathbf{x}))$

- output layer activation ( $k=L+1$ ):  
 $\mathbf{h}^{(L+1)}(\mathbf{x}) = \mathbf{o}(\mathbf{a}^{(L+1)}(\mathbf{x})) = \mathbf{f}(\mathbf{x})$



37

### Empirical Risk Minimization

**Topics:** empirical risk minimization, regularization

- Empirical risk minimization
  - framework to design learning algorithms

$$\arg \min_{\theta} \frac{1}{T} \sum_t l(f(\mathbf{x}^{(t)}; \theta), y^{(t)}) + \lambda \Omega(\theta)$$

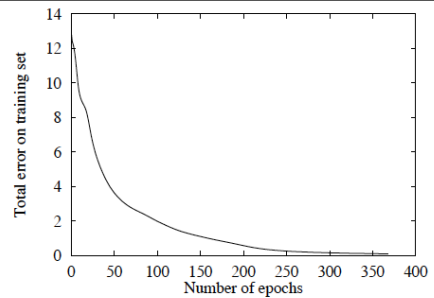
- $l(f(\mathbf{x}^{(t)}; \theta), y^{(t)})$  is a loss function
- $\Omega(\theta)$  is a regularizer (penalizes certain values of  $\theta$ )
- Learning is cast as optimization
  - ideally, we'd optimize classification error, but it's not smooth
  - loss function is a surrogate for what we truly should optimize (e.g. upper bound)

38

### LOSS FUNCTION

**Topics:** loss function for classification

- Neural network estimates  $f(\mathbf{x})_c = p(y=c|\mathbf{x})$ 
  - we could maximize the probabilities of  $y^{(i)}$  given  $\mathbf{x}^{(i)}$  in the training set
- To frame as minimization, we minimize the negative log-likelihood
  - natural log (ln)
$$l(f(\mathbf{x}), y) = -\sum_c \mathbf{1}_{(y=c)} \log f(\mathbf{x})_c = -\log f(\mathbf{x})_y$$
  - we take the log to simplify for numerical stability and math simplicity
  - sometimes referred to as cross-entropy



[figure from Greg Mori's slides]

39

40

### REGULARIZATION

**Topics:** L2 regularization

$$\Omega(\theta) = \sum_k \sum_i \sum_j (W_{i,j}^{(k)})^2 = \sum_k \|\mathbf{W}^{(k)}\|_F^2$$

41

### Empirical Risk Minimization

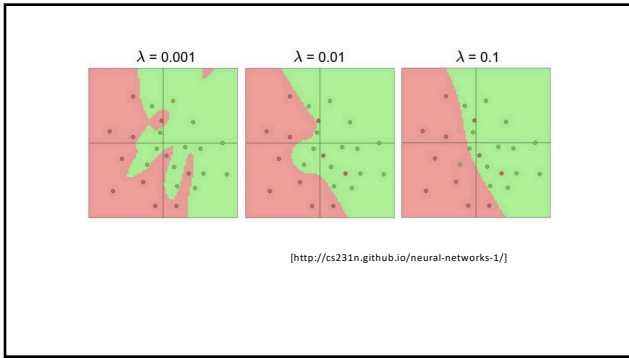
**Topics:** empirical risk minimization, regularization

- Empirical risk minimization
  - framework to design learning algorithms

$$\arg \min_{\theta} \frac{1}{T} \sum_t l(f(\mathbf{x}^{(t)}; \theta), y^{(t)}) + \lambda \Omega(\theta)$$

- $l(f(\mathbf{x}^{(t)}; \theta), y^{(t)})$  is a loss function
- $\Omega(\theta)$  is a regularizer (penalizes certain values of  $\theta$ )
- Learning is cast as optimization
  - ideally, we'd optimize classification error, but it's not smooth
  - loss function is a surrogate for what we truly should optimize (e.g. upper bound)

42



43

### INITIALIZATION

**Topics:** initialization

- For biases
  - initialize all to 0
- For weights
  - Can't initialize weights to 0 with tanh activation
    - we can show that all gradients would then be 0 (saddle point)
  - Can't initialize all weights to the same value
    - we can show that all hidden units in a layer will always behave the same
    - need to break symmetry
  - Recipe: sample  $\mathbf{W}_{ij}^{(k)}$  from  $U[-b, b]$  where  $b = \frac{\sqrt{6}}{\sqrt{n_k + n_{k-1}}}$ 
    - the idea is to sample around 0 but break symmetry
    - other values of  $b$  could work well (not an exact science) (see Glorot & Bengio, 2010)

44

### Model Learning

- Backpropagation (BP) algorithm (not required for this course)
- Further reading on BP:
  - <https://towardsdatascience.com/understanding-backpropagation-algorithm-7bb3aa2f95fd>
  - <https://matmazur.com/2015/03/17/a-step-by-step-backpropagation-example/>

45

### Outline

- Maximum Entropy
- Feedforward Neural Networks
- ➔ • Recurrent Neural Networks

46

### Long Distance Dependencies

- It is very difficult to train NNs to retain information over many time steps
- This makes it very difficult to handle long-distance dependencies, such as subject-verb agreement.
- E.g. Jane walked into the room. John walked in too. It was late in the day. Jane said hi to \_?\_

47

### Recurrent Neural Networks

Feed-forward NN

$$\mathbf{h} = g(\mathbf{V}\mathbf{x} + \mathbf{c})$$

$$\hat{\mathbf{y}} = \mathbf{W}\mathbf{h} + \mathbf{b}$$

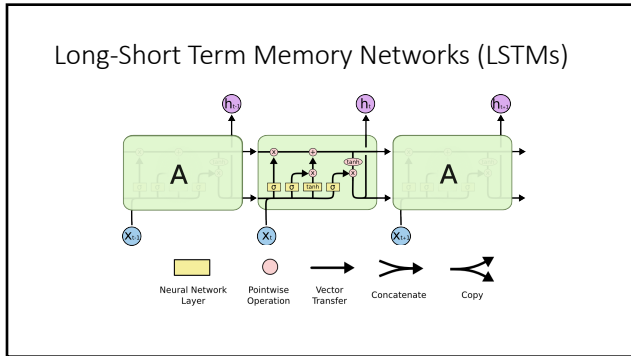
Recurrent NN

$$\mathbf{h}_t = g(\mathbf{V}\mathbf{x}_t + \mathbf{U}\mathbf{h}_{t-1} + \mathbf{c})$$

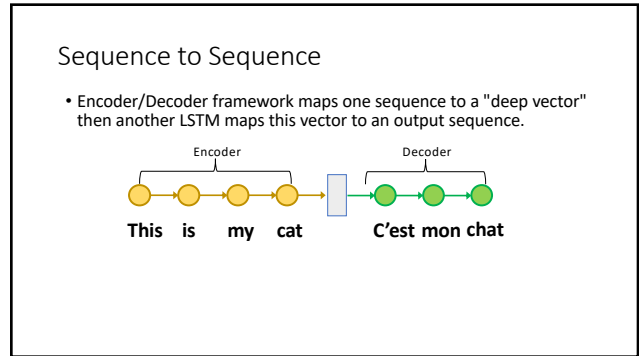
$$\hat{\mathbf{y}}_t = \mathbf{W}\mathbf{h}_t + \mathbf{b}$$

48

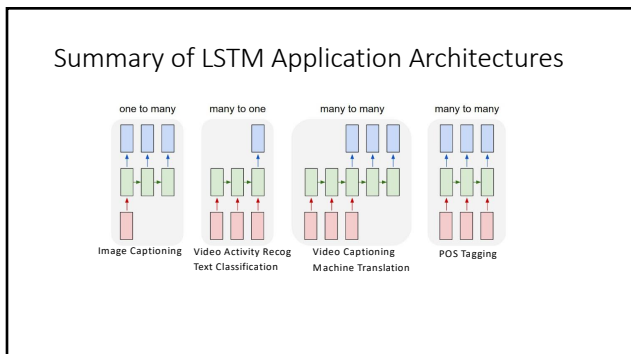




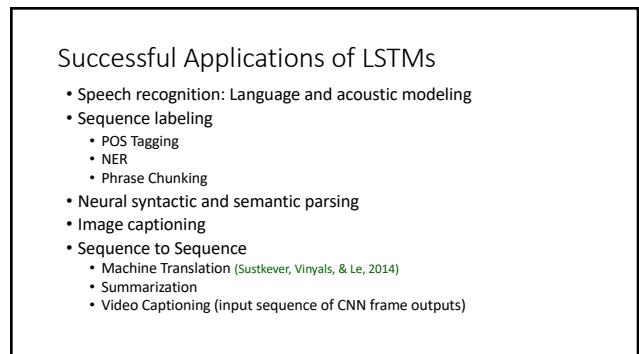
49



50



51



52