

CS 6120/CS4120: Natural Language Processing

Instructor: Prof. Lu Wang
 College of Computer and Information Science
 Northeastern University
 Webpage: www.ccs.neu.edu/home/luwang

Outline

- Maximum Entropy
- Feedforward Neural Networks
- Recurrent Neural Networks

Introduction

- So far we've looked at "generative models"
 - Language models, Naive Bayes
- But there is now much use of conditional or discriminative probabilistic models in NLP, Speech, IR (and ML generally)
- Because:
 - They give high accuracy performance
 - They make it easy to incorporate lots of linguistically important features
 - They allow automatic building of language independent, retargetable NLP modules

Joint vs. Conditional Models

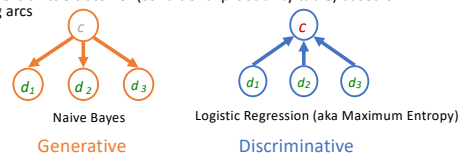
- We have some data $\{(d, c)\}$ of paired observations d and hidden classes c .
- **Joint (generative) models** place probabilities over both observed data and the hidden stuff (generate the observed data from hidden stuff): $P(c, d)$
 - All the classic statistic NLP models:
 - n -gram models, Naive Bayes classifiers, hidden Markov models, probabilistic context-free grammars, IBM machine translation alignment models

Joint vs. Conditional Models

- **Discriminative (conditional) models** take the data as given, and put a probability over hidden structure given the data: $P(c|d)$
 - Logistic regression, conditional loglinear or maximum entropy models, conditional random fields
 - Also, SVMs, (averaged) perceptron, etc. are discriminative classifiers (but not directly probabilistic)

Bayes Network/Graphical Models

- Bayes network diagrams draw circles for random variables, and lines for direct dependencies
- Some variables are observed; some are hidden
- Each node is a little classifier (conditional probability table) based on incoming arcs



Conditional vs. Joint Likelihood

- A *joint* model gives probabilities $P(d,c)$ and tries to maximize this joint likelihood.
 - It turns out to be trivial to choose weights: just relative frequencies.
- A *conditional* model gives probabilities $P(c|d)$. It takes the data as given and models only the conditional probability of the class.
 - We seek to maximize conditional likelihood.
 - Harder to do.
 - More closely related to classification error. (Easy to tune!)

Maximum Entropy (MaxEnt)

- Or logistic regression

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: $f: C \times D \rightarrow \mathbb{R}$

Example Task: Named Entity Type

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

Example features

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

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

Features

- In NLP uses, usually a feature specifies
 1. an indicator function – a yes/no boolean matching function – of properties of the input and a particular class

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

- Each feature picks out a data subset and suggests a label for it

Feature-Based Models

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

Data BUSINESS: Stocks hit a yearly low ...	Data ... to restructure bank:MONEY debt.	Data DT JJ NN ... The previous fall ...
Label: BUSINESS Features {..., stocks, hit, a, yearly, low, ...}	Label: MONEY Features {..., w ₋₁ =restructure, w ₋₁ =debt, L=12, ...}	Label: NN Features {w=fall, t ₋₁ =JJ w ₋₁ =previous}
Text Categorization	Word-Sense Disambiguation	POS Tagging

Other Maxent Classifier Examples

- You can use a maxent classifier whenever you want to assign data points to one of a number of classes:
 - Sentence boundary detection (Mikheev 2000)
 - Is a period end of sentence or abbreviation?
 - Sentiment analysis (Pang and Lee 2002)
 - Word unigrams, bigrams, POS counts, ...
 - Prepositional phrase attachment (Ratnaparkhi 1998)
 - Attach to verb or noun? Features of head noun, preposition, etc.
 - Parsing decisions in general (Ratnaparkhi 1997; Johnson et al. 1999, etc.)

Feature-Based Linear Classifiers

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

PERSON
in Québec
LOCATION
in Québec
DRUG
in Québec

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

Feature-Based Linear Classifiers

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

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

Maximum Entropy:

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

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

Makes votes positive
Normalizes votes

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

Maximum Entropy:

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

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

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

Feature-Based Linear Classifiers

- Exponential models:
 - Given this model form, we will choose parameters $\{\lambda_i\}$ that *maximize the conditional likelihood* of the data according to this model.
 - We construct not only classifications, but probability distributions over classifications.
 - There are other (good!) ways of discriminating classes – SVMs, boosting, even perceptrons – but these methods are not as trivial to interpret as distributions over classes.

Outline

- Maximum Entropy
- Feedforward Neural Networks
- Recurrent Neural Networks

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.

ARTIFICIAL NEURON

Topics: connection weights, bias, activation function

• Neuron pre-activation (or input activation):

$$a(\mathbf{x}) = b + \sum_i w_i x_i = b + \mathbf{w}^T \mathbf{x}$$

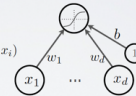
• Neuron (output) activation

$$h(\mathbf{x}) = g(a(\mathbf{x})) = g(b + \sum_i w_i x_i)$$

• \mathbf{w} are the connection weights

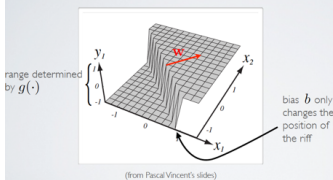
• b is the neuron bias

• $g(\cdot)$ is called the activation function



ARTIFICIAL NEURON

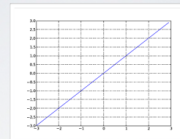
Topics: connection weights, bias, activation function



ACTIVATION FUNCTION

Topics: linear activation function

- Performs no input squashing
- Not very interesting...

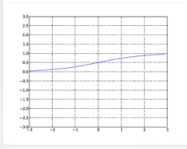


$$g(a) = a$$

ACTIVATION FUNCTION

Topics: sigmoid activation function

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

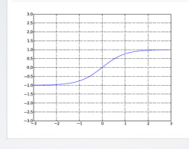


$$g(a) = \text{sigm}(a) = \frac{1}{1 + \exp(-a)}$$

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

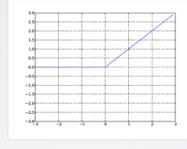


$$g(a) = \text{tanh}(a) = \frac{\exp(a) - \exp(-a)}{\exp(a) + \exp(-a)} = \frac{\exp(2a) - 1}{\exp(2a) + 1}$$

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

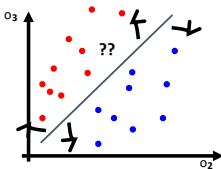


$$g(a) = \text{reclin}(a) = \max(0, a)$$

```
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
```

Linear Separator

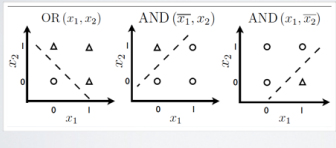
- Since one-layer neuron (aka perceptron) uses linear threshold function, it is searching for a linear separator that discriminates the classes.



ARTIFICIAL NEURON

Topics: capacity of single neuron

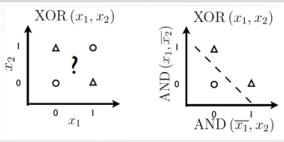
- Can solve linearly separable problems



ARTIFICIAL NEURON

Topics: capacity of single neuron

- Can't solve non linearly separable problems...



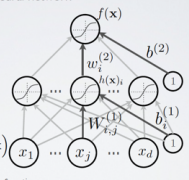
• ... unless the input is transformed in a better representation

NEURAL NETWORK

Topics: single hidden layer neural network

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

output activation function



NEURAL NETWORK

Topics: softmax activation function

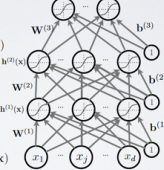
- For multi-class classification:
 - we need multiple outputs (1 output per class)
 - we would like to estimate the conditional probability $p(y = c|\mathbf{x})$
- We use the softmax activation function at the output:

$$\mathbf{o}(\mathbf{a}) = \text{softmax}(\mathbf{a}) = \left[\frac{\exp(a_1)}{\sum_c \exp(a_c)}, \dots, \frac{\exp(a_C)}{\sum_c \exp(a_c)} \right]^\top$$
 - strictly positive
 - sums to one
- Predicted class is the one with highest estimated probability

NEURAL NETWORK

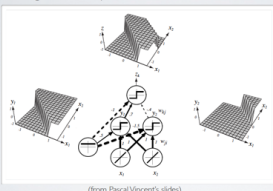
Topics: multilayer neural network

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



CAPACITY OF NEURAL NETWORK

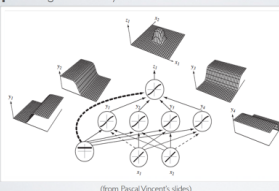
Topics: single hidden layer neural network



(from Pascal Vincent's slides)

CAPACITY OF NEURAL NETWORK

Topics: single hidden layer neural network



(from Pascal Vincent's slides)

CAPACITY OF NEURAL NETWORK

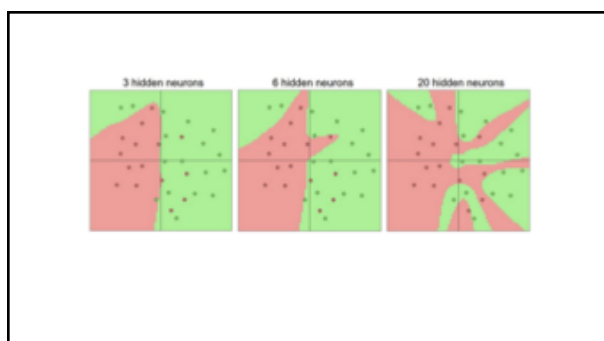
Topics: single hidden layer neural network

(from Pascal Vincent's slides)

CAPACITY OF NEURAL NETWORK

Topics: universal approximation

- Universal approximation theorem (Hornik, 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!



How to train a neural network?

Topics: multilayer neural network

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

Empirical Risk Minimization

Topics: empirical risk minimization, regularization

- Empirical risk minimization
 - framework to design learning algorithms

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

- $l(f(\mathbf{x}^{(l)}; \theta), y^{(l)})$ is a loss function
- $\Omega(\theta)$ is a regularizer (penalizes certain values of θ)
- 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)

The rest of the slides are for reference only (not covered by lecture or exam)

The Learning Algorithm

- Topics:** stochastic gradient descent (SGD)
- Algorithm that performs updates after each example
 - initialize θ ($\theta \equiv \{\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \dots, \mathbf{W}^{(L+1)}, \mathbf{b}^{(L+1)}\}$)
 - for N iterations
 - for each training example $(\mathbf{x}^{(i)}, y^{(i)})$
 - $\Delta = -\nabla_{\theta} l(f(\mathbf{x}^{(i)}; \theta), y^{(i)}) - \lambda \nabla_{\theta} \Omega(\theta)$
 - $\theta \leftarrow \theta + \alpha \Delta$
- } training epoch
= iteration over all examples
- To apply this algorithm to neural network training, we need
 - the loss function $l(f(\mathbf{x}^{(i)}; \theta), y^{(i)})$
 - a procedure to compute the parameter gradients $\nabla_{\theta} l(f(\mathbf{x}^{(i)}; \theta), y^{(i)})$
 - the regularizer $\Omega(\theta)$ (and the gradient $\nabla_{\theta} \Omega(\theta)$)
 - initialization method

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)
- $$l(f(\mathbf{x}), y) = -\sum_c 1_{(y=c)} \log f(\mathbf{x})_c = -\log f(\mathbf{x})_y$$
- we take the log to simplify for numerical stability and math simplicity
 - sometimes referred to as cross-entropy

The Learning Algorithm

- Topics:** stochastic gradient descent (SGD)
- Algorithm that performs updates after each example
 - initialize θ ($\theta \equiv \{\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \dots, \mathbf{W}^{(L+1)}, \mathbf{b}^{(L+1)}\}$)
 - for N iterations
 - for each training example $(\mathbf{x}^{(i)}, y^{(i)})$
 - $\Delta = -\nabla_{\theta} l(f(\mathbf{x}^{(i)}; \theta), y^{(i)}) - \lambda \nabla_{\theta} \Omega(\theta)$
 - $\theta \leftarrow \theta + \alpha \Delta$
- } training epoch
= iteration over all examples
- To apply this algorithm to neural network training, we need
 - the loss function $l(f(\mathbf{x}^{(i)}; \theta), y^{(i)})$
 - a procedure to compute the parameter gradients $\nabla_{\theta} l(f(\mathbf{x}^{(i)}; \theta), y^{(i)})$
 - the regularizer $\Omega(\theta)$ (and the gradient $\nabla_{\theta} \Omega(\theta)$)
 - initialization method

Gradient Computation

- Output layer gradient (o)
- Hidden layer gradient (h)
- Activation function gradient (a)
- Parameter gradient (W, b)

Gradient Computation

- Output layer gradient (o)
- Hidden layer gradient (h)
- Activation function gradient (a)
- Parameter gradient (W, b)

GRADIENT COMPUTATION

Topics: loss gradient at output

- Partial derivative:

$$\frac{\partial}{\partial f(\mathbf{x})_c} -\log f(\mathbf{x})_y = \frac{-1_{(y=c)}}{f(\mathbf{x})_y}$$
- Gradient:

$$\begin{aligned} \nabla_{f(\mathbf{x})} -\log f(\mathbf{x})_y &= \frac{-1_{(y=c)}}{f(\mathbf{x})_y} \begin{bmatrix} 1_{(y=0)} \\ \vdots \\ 1_{(y=C-1)} \end{bmatrix} \\ &= \frac{-\mathbf{e}(y)}{f(\mathbf{x})_y} \end{aligned}$$

GRADIENT COMPUTATION

Topics: loss gradient at output pre-activation

- Partial derivative:

$$\frac{\partial}{\partial a^{(L+1)}(x)_c} - \log f(x)_y$$

$$= -(1_{(y=c)} - f(x)_c)$$
- Gradient:

$$\nabla_{a^{(L+1)}(x)} - \log f(x)_y$$

$$= -(\mathbf{e}(y) - \mathbf{f}(x))$$

$$\frac{\partial}{\partial a^{(L+1)}(x)_c} - \log f(x)_y$$

$$\frac{\partial}{\partial a^{(L+1)}(x)_c} - \log f(x)_y$$

$$= \frac{-1}{f(x)_y} \frac{\partial}{\partial a^{(L+1)}(x)_c} f(x)_y$$

$$\frac{\partial}{\partial a^{(L+1)}(x)_c} - \log f(x)_y$$

$$= \frac{-1}{f(x)_y} \frac{\partial}{\partial a^{(L+1)}(x)_c} f(x)_y$$

$$= \frac{-1}{f(x)_y} \frac{\partial}{\partial a^{(L+1)}(x)_c} \text{softmax}(a^{(L+1)}(x))_y$$

$$\frac{\partial}{\partial a^{(L+1)}(x)_c} - \log f(x)_y$$

$$= \frac{-1}{f(x)_y} \frac{\partial}{\partial a^{(L+1)}(x)_c} f(x)_y$$

$$= \frac{-1}{f(x)_y} \frac{\partial}{\partial a^{(L+1)}(x)_c} \text{softmax}(a^{(L+1)}(x))_y$$

$$= \frac{-1}{f(x)_y} \frac{\partial}{\partial a^{(L+1)}(x)_c} \frac{\exp(a^{(L+1)}(x)_y)}{\sum_c \exp(a^{(L+1)}(x)_c)}$$

$$\frac{\partial}{\partial a^{(L+1)}(x)_c} - \log f(x)_y$$

$$= \frac{-1}{f(x)_y} \frac{\partial}{\partial a^{(L+1)}(x)_c} f(x)_y$$

$$= \frac{-1}{f(x)_y} \frac{\partial}{\partial a^{(L+1)}(x)_c} \text{softmax}(a^{(L+1)}(x))_y$$

$$= \frac{-1}{f(x)_y} \frac{\partial}{\partial a^{(L+1)}(x)_c} \frac{\exp(a^{(L+1)}(x)_y)}{\sum_c \exp(a^{(L+1)}(x)_c)}$$

$$= \frac{-1}{f(x)_y} \left(\frac{\frac{\partial}{\partial a^{(L+1)}(x)_c} \exp(a^{(L+1)}(x)_y)}{\sum_c \exp(a^{(L+1)}(x)_c)} - \frac{\exp(a^{(L+1)}(x)_y)}{(\sum_c \exp(a^{(L+1)}(x)_c))^2} \right)$$

$$\frac{\partial g(x)}{\partial x} = \frac{\partial g(x)}{\partial x} \frac{1}{h(x)} - \frac{g(x)}{h(x)^2} \frac{\partial h(x)}{\partial x}$$

GRADIENT COMPUTATION

Topics: chain rule

- If a function $p(a)$ can be written as a function of intermediate results $q_i(a)$ then we have:

$$\frac{\partial p(a)}{\partial a} = \sum_i \frac{\partial p(a)}{\partial q_i(a)} \frac{\partial q_i(a)}{\partial a}$$

- We can invoke it by setting
 - a to a unit in layer
 - $q_i(a)$ to a pre-activation in the layer above
 - $p(a)$ is the loss function

GRADIENT COMPUTATION

Topics: loss gradient at hidden layers

- Partial derivative:

$$\begin{aligned} & \frac{\partial}{\partial h^{(k)}(x)_j} - \log f(x)_y \\ &= \sum_i \frac{\partial -\log f(x)_y}{\partial a^{(k+1)}(x)_i} \frac{\partial a^{(k+1)}(x)_i}{\partial h^{(k)}(x)_j} \\ &= \sum_i \frac{\partial -\log f(x)_y}{\partial a^{(k+1)}(x)_i} W_{i,j}^{(k+1)} \\ &= (\mathbf{W}^{(k+1)})^\top (\nabla_{\mathbf{a}^{(k+1)}(\mathbf{x})} - \log f(\mathbf{x})_y) \end{aligned}$$

REMINDER: $a^{(k)}(x)_i = h_i^{(k)} + \sum_j W_{ij}^{(k)} h_j^{(k-1)}(x)$

GRADIENT COMPUTATION

Topics: loss gradient at hidden layers

- Gradient:

$$\nabla_{h^{(k)}(\mathbf{x})} - \log f(\mathbf{x})_y = (\mathbf{W}^{(k+1)})^\top (\nabla_{\mathbf{a}^{(k+1)}(\mathbf{x})} - \log f(\mathbf{x})_y)$$

REMINDER: $a^{(k)}(x)_i = h_i^{(k)} + \sum_j W_{ij}^{(k)} h_j^{(k-1)}(x)$

GRADIENT COMPUTATION

Topics: loss gradient at hidden layers pre-activation

- Partial derivative:

$$\begin{aligned} & \frac{\partial}{\partial a^{(k)}(x)_j} - \log f(x)_y \\ &= \frac{\partial -\log f(x)_y}{\partial h^{(k)}(x)_j} \frac{\partial h^{(k)}(x)_j}{\partial a^{(k)}(x)_j} \\ &= \frac{\partial -\log f(x)_y}{\partial h^{(k)}(x)_j} g'(a^{(k)}(x)_j) \end{aligned}$$

REMINDER: $h^{(k)}(x)_j = y(a^{(k)}(x)_j)$

GRADIENT COMPUTATION

Topics: loss gradient at hidden layers pre-activation

- Gradient:

$$\begin{aligned} & \nabla_{a^{(k)}(\mathbf{x})} - \log f(\mathbf{x})_y \\ &= (\nabla_{h^{(k)}(\mathbf{x})} - \log f(\mathbf{x})_y)^\top \nabla_{a^{(k)}(\mathbf{x})} h^{(k)}(\mathbf{x}) \\ &= (\nabla_{h^{(k)}(\mathbf{x})} - \log f(\mathbf{x})_y) \odot [\dots, g'(a^{(k)}(x)_j), \dots] \end{aligned}$$

↑
element-wise product

REMINDER: $h^{(k)}(\mathbf{x})_j = g(a^{(k)}(\mathbf{x})_j)$

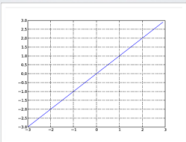
Gradient Computation

- Output layer gradient (o)
- Hidden layer gradient (h)
- Activation function gradient (a)
- Parameter gradient (W, b)

ACTIVATION FUNCTION

Topics: linear activation function gradient

• Partial derivative:
 $g'(a) = 1$

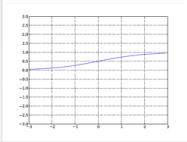


$g(a) = a$

ACTIVATION FUNCTION

Topics: sigmoid activation function gradient

• Partial derivative:
 $g'(a) = g(a)(1 - g(a))$

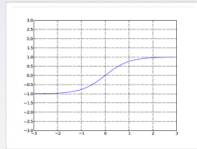


$g(a) = \text{sigm}(a) = \frac{1}{1 + \exp(-a)}$

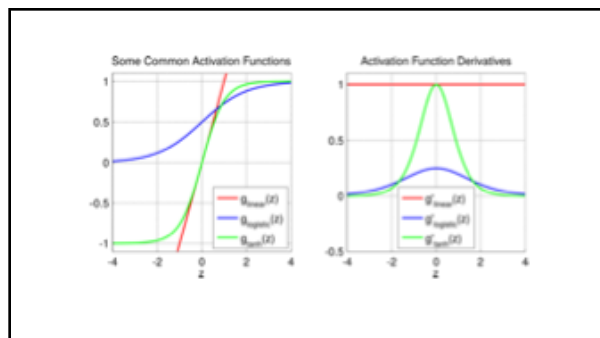
ACTIVATION FUNCTION

Topics: tanh activation function gradient

• Partial derivative:
 $g'(a) = 1 - g(a)^2$



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



- ### Gradient Computation
- Output layer gradient (o)
 - Hidden layer gradient (h)
 - Activation function gradient (a)
 - Parameter gradient (W, b)

GRADIENT COMPUTATION

Topics: loss gradient of parameters

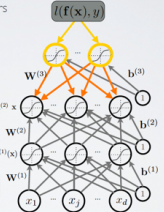
• Partial derivative (weights):

$$\frac{\partial}{\partial W_{i,j}^{(k)}} - \log f(\mathbf{x})_p$$

$$= \frac{\partial - \log f(\mathbf{x})_p}{\partial a^{(k)}(\mathbf{x})_i} \frac{\partial a^{(k)}(\mathbf{x})_i}{\partial W_{i,j}^{(k)}}$$

$$= \frac{\partial - \log f(\mathbf{x})_p}{\partial a^{(k)}(\mathbf{x})_i} j_i^{(k-1)}(\mathbf{x})$$

REMEMBER: $a^{(k)}(\mathbf{x})_i = b_i^{(k)} + \sum_j W_{ij}^{(k)} j_j^{(k-1)}(\mathbf{x})$



GRADIENT COMPUTATION

Topics: loss gradient of parameters

- Gradient (weights):

$$\nabla_{W^{(k)}} - \log f(\mathbf{x})_y = (\nabla_{a^{(k)}(\mathbf{x})} - \log f(\mathbf{x})_y) \mathbf{h}^{(k-1)\top}(\mathbf{x})$$

$$= \frac{\partial -\log f(\mathbf{x})_y}{\partial a^{(k)}(\mathbf{x})} \frac{\partial a^{(k)}(\mathbf{x})}{\partial h^{(k-1)}(\mathbf{x})}$$

$$= \frac{\partial -\log f(\mathbf{x})_y}{\partial a^{(k)}(\mathbf{x})} \sum_j W_{ij}^{(k)} h_j^{(k-1)}(\mathbf{x})$$

REMINDER
 $a^{(k)}(\mathbf{x}) = b_i^{(k)} + \sum_j W_{ij}^{(k)} h_j^{(k-1)}(\mathbf{x})$

GRADIENT COMPUTATION

Topics: loss gradient of parameters

- Partial derivative (biases):

$$\frac{\partial}{\partial b_i^{(k)}} - \log f(\mathbf{x})_y = \frac{\partial -\log f(\mathbf{x})_y}{\partial a^{(k)}(\mathbf{x})} \frac{\partial a^{(k)}(\mathbf{x})}{\partial b_i^{(k)}}$$

$$= \frac{\partial -\log f(\mathbf{x})_y}{\partial a^{(k)}(\mathbf{x})} \frac{\partial a^{(k)}(\mathbf{x})}{\partial h^{(k-1)}(\mathbf{x})}$$

$$= \frac{\partial -\log f(\mathbf{x})_y}{\partial a^{(k)}(\mathbf{x})} \sum_j W_{ij}^{(k)} h_j^{(k-1)}(\mathbf{x})$$

REMINDER
 $a^{(k)}(\mathbf{x}) = b_i^{(k)} + \sum_j W_{ij}^{(k)} h_j^{(k-1)}(\mathbf{x})$

GRADIENT COMPUTATION

Topics: loss gradient of parameters

- Gradient (biases):

$$\nabla_{b^{(k)}} - \log f(\mathbf{x})_y = \nabla_{a^{(k)}(\mathbf{x})} - \log f(\mathbf{x})_y$$

$$= \frac{\partial -\log f(\mathbf{x})_y}{\partial a^{(k)}(\mathbf{x})}$$

REMINDER
 $a^{(k)}(\mathbf{x}) = b_i^{(k)} + \sum_j W_{ij}^{(k)} h_j^{(k-1)}(\mathbf{x})$

Backpropagation

FLOW GRAPH

Topics: flow graph

- Forward propagation can be represented as an acyclic flow graph
- It's a nice way of implementing forward propagation in a modular way
- each box could be an object with an fprop method, that computes the value of the box given its children
- calling the fprop method of each box in the right order yield forward propagation

```

# 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)
    
```

FLOW GRAPH

Topics: automatic differentiation

- Each object also has a bprop method
 - it computes the gradient of the loss with respect to each children
 - fprop depends on the fprop of a box's children, while bprop depends the bprop of a box's parents
- By calling bprop in the reverse order, we get back-propagation
 - only need to reach the parameters

FLOW GRAPH

Topics: automatic differentiation

- Each object also has a bprop method
 - it computes the gradient of the loss with respect to each children
 - fprop depends on the fprop of a box's children, while bprop depends the bprop of a box's parents
- By calling bprop in the reverse order, we get back-propagation
 - only need to reach the parameters

BACKPROPAGATION

Topics: backpropagation algorithm

- This assumes a forward propagation has been made before
 - compute output gradient (before activation)

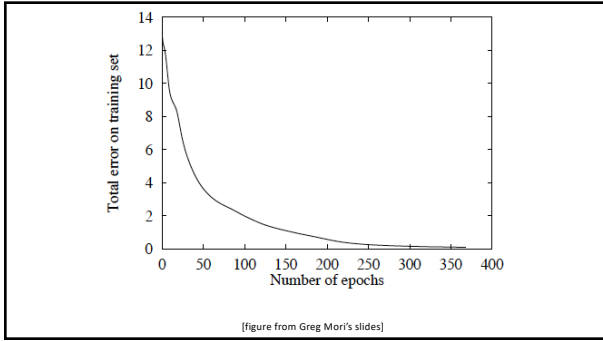
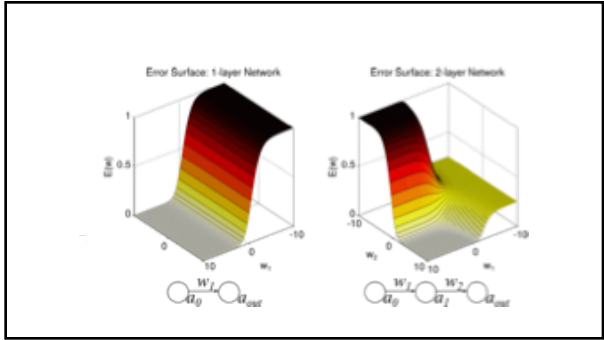
$$\nabla_{a^{(L+1)(x)}} - \log f(x)_y \leftarrow -(\epsilon(y) - f(x))$$
 - for k from $L-1$ to 1
 - compute gradients of hidden layer parameter

$$\nabla_{W^{(k)}} - \log f(x)_y \leftarrow (\nabla_{a^{(k+1)(x)}} - \log f(x)_y) \cdot h^{(k-1)(x)T}$$

$$\nabla_{b^{(k)}} - \log f(x)_y \leftarrow \nabla_{a^{(k+1)(x)}} - \log f(x)_y$$
 - compute gradient of hidden layer below

$$\nabla_{a^{(k-1)(x)}} - \log f(x)_y \leftarrow W^{(k)T} (\nabla_{a^{(k+1)(x)}} - \log f(x)_y)$$
 - compute gradient of hidden layer before (before activation)

$$\nabla_{a^{(k-1)(x)}} - \log f(x)_y \leftarrow (\nabla_{a^{(k-1)(x)}} - \log f(x)_y) \odot [\dots, g'(a^{(k-1)(x)}_j), \dots]$$



The Learning Algorithm

Topics: stochastic gradient descent (SGD)

- Algorithm that performs updates after each example
 - initialize θ ($\theta \equiv \{W^{(1)}, b^{(1)}, \dots, W^{(L+1)}, b^{(L+1)}\}$)
 - for N iterations
 - for each training example $(x^{(i)}, y^{(i)})$

$$\Delta = -\nabla_{\theta} l(f(x^{(i)}; \theta), y^{(i)}) - \lambda \nabla_{\theta} \Omega(\theta)$$
 - $\theta \leftarrow \theta + \alpha \Delta$
- To apply this algorithm to neural network training, we need
 - the loss function $l(f(x^{(i)}; \theta), y^{(i)})$
 - a procedure to compute the parameter gradients $\nabla_{\theta} l(f(x^{(i)}; \theta), y^{(i)})$
 - the regularizer $\Omega(\theta)$ (and the gradient $\nabla_{\theta} \Omega(\theta)$)
 - initialization method

REGULARIZATION

Topics: L2 regularization

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

- Gradient: $\nabla_{\mathbf{W}^{(k)}} \Omega(\theta) = 2\mathbf{W}^{(k)}$
- Only applied on weights, not on biases (weight decay)
- Can be interpreted as having a Gaussian prior over the weights

REGULARIZATION

Topics: L1 regularization

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

- Gradient: $\nabla_{\mathbf{W}^{(k)}} \Omega(\theta) = \text{sign}(\mathbf{W}^{(k)})$
 - where $\text{sign}(\mathbf{W}^{(k)})_{i,j} = 1_{W_{i,j}^{(k)} > 0} - 1_{W_{i,j}^{(k)} < 0}$
- Also only applied on weights
- Unlike L2, L1 will push certain weights to be exactly 0
- Can be interpreted as having a Laplacian prior over the weights

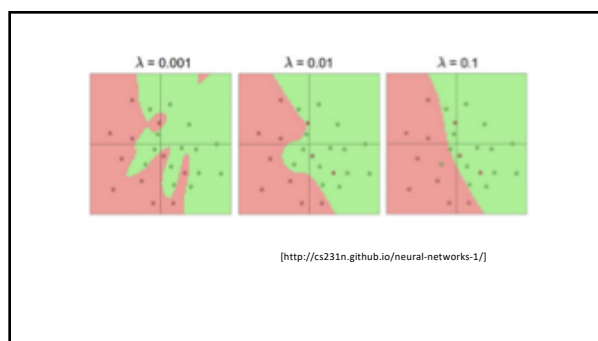
Empirical Risk Minimization

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 θ)
- 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)



The Learning Algorithm

Topics: stochastic gradient descent (SGD)

- Algorithm that performs updates after each example
 - initialize θ ($\theta \equiv \{\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \dots, \mathbf{W}^{(L+1)}, \mathbf{b}^{(L+1)}\}$)
 - for N iterations
 - for each training example $(\mathbf{x}^{(t)}, y^{(t)})$
 - $\Delta = -\nabla_{\theta} l(f(\mathbf{x}^{(t)}; \theta), y^{(t)}) - \lambda \nabla_{\theta} \Omega(\theta)$
 - $\theta \leftarrow \theta + \alpha \Delta$

} training epoch
= iteration over all examples

- To apply this algorithm to neural network training, we need
 - the loss function $l(f(\mathbf{x}^{(t)}; \theta), y^{(t)})$
 - a procedure to compute the parameter gradients $\nabla_{\theta} l(f(\mathbf{x}^{(t)}; \theta), y^{(t)})$
 - the regularizer $\Omega(\theta)$ (and the gradient $\nabla_{\theta} \Omega(\theta)$)
 - initialization method

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
 - Recipes sample $\mathbf{W}_{i,j}^{(k)}$ 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)

The Learning Algorithm

Topics: stochastic gradient descent (SGD)

- Algorithm that performs updates after each example
 - initialize θ ($\theta \equiv \{\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \dots, \mathbf{W}^{(L+1)}, \mathbf{b}^{(L+1)}\}$)
 - for N iterations
 - for each training example $(\mathbf{x}^{(i)}, y^{(i)})$
 - $\Delta = -\nabla_{\theta} l(f(\mathbf{x}^{(i)}; \theta), y^{(i)}) - \lambda \nabla_{\theta} \Omega(\theta)$
 - $\theta \leftarrow \theta + \alpha \Delta$
- } training epoch
= iteration over **all** examples
- To apply this algorithm to neural network training, we need
 - the loss function $l(f(\mathbf{x}^{(i)}; \theta), y^{(i)})$
 - a procedure to compute the parameter gradients $\nabla_{\theta} l(f(\mathbf{x}^{(i)}; \theta), y^{(i)})$
 - the regularizer $\Omega(\theta)$ (and the gradient $\nabla_{\theta} \Omega(\theta)$)
 - initialization method

Toolkits

- TensorFlow
 - <https://www.tensorflow.org/>
- Theano (not maintained any more)
 - <http://deeplearning.net/software/theano/>
- PyTorch
 - <http://pytorch.org/>