

Natural Language Processing with Deep Learning

CS224N/Ling284



Lecture 4: Word Window Classification
and Neural Networks

Christopher Manning and Richard Socher

Overview Today:

- Classification background
- Updating word vectors for classification
- Window classification & cross entropy error derivation tips
- A single layer neural network!
- Max-Margin loss and **backprop**

This lecture will help a lot with PSet1 :)

Classification setup and notation

- Generally we have a training dataset consisting of samples

$$\{x_i, y_i\}_{i=1}^N$$

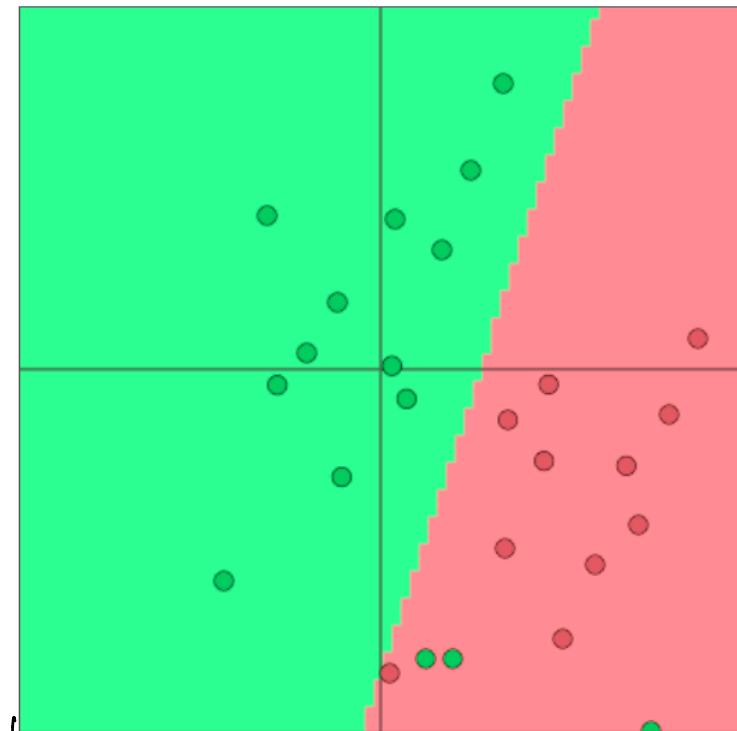
- x_i - inputs, e.g. words (indices or vectors!), context windows, sentences, documents, etc.

 Dimension ↕

- y_i - labels we try to predict, for example
 - class: sentiment, named entities, buy/sell decision,
 - other words
 - later: multi-word sequences

Classification intuition

- Training data: $\{x_i, y_i\}_{i=1}^N$
- Simple illustration case:
 - Fixed 2d word vectors to classify
 - Using logistic regression
 - → linear decision boundary →
- General ML: assume x is fixed,
train logistic regression weights $W \in \mathbb{R}^{C \times d}$
→ only modify the decision boundary
- Goal: predict for each x : $p(y|x)$ where $W \in \mathbb{R}^{C \times d}$



Visualizations with ConvNetJS by Karpathy!
<http://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>

$$p(y|x) = \frac{\exp(W_y \cdot x)}{\sum_{c=1}^C \exp(W_c \cdot x)}$$

Details of the softmax

- We can tease apart into two steps:

$$p(y|x) = \frac{\exp(W_y \cdot x)}{\sum_{c=1}^C \exp(W_c \cdot x)}$$

- Take the y 'th row of W and multiply that row with x :

$\mathbb{C} \begin{bmatrix} W \\ \vdots \\ \text{y'th row} \\ \vdots \\ W \end{bmatrix} \begin{bmatrix} x \\ \vdots \\ x \end{bmatrix} \Rightarrow \begin{bmatrix} \text{y'th row} \\ \vdots \\ \text{y'th row} \\ \vdots \\ \text{y'th row} \end{bmatrix} \begin{bmatrix} x \\ \vdots \\ x \end{bmatrix} \Rightarrow f_y$

Unnormalized Score.

$$W_y \cdot x = \sum_{i=1}^d W_{yi} x_i = f_y$$

Compute all f_c for $c=1, \dots, C$

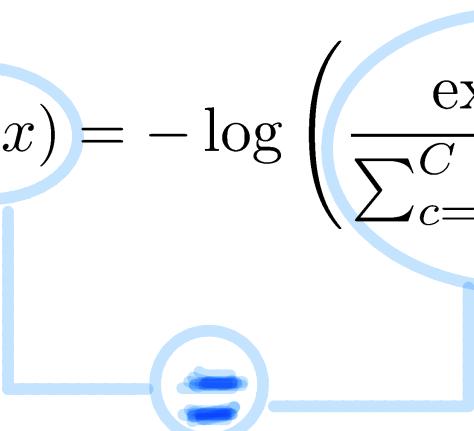
- Normalize to obtain probability with softmax function:

$$p(y|x) = \frac{\exp(f_y)}{\sum_{c=1}^C \exp(f_c)} = \text{softmax}(f)_y$$

Sums to One.

The softmax and cross-entropy error

- For each training example $\{x, y\}$, our objective is to maximize the probability of the correct class y
- Hence, we minimize the negative log probability of that class:

$$-\log p(y|x) = -\log \left(\frac{\exp(f_y)}{\sum_{c=1}^C \exp(f_c)} \right)$$


Background: Why “Cross entropy” error

- Assuming a ground truth (or gold or target) probability distribution that is 1 at the right class and 0 everywhere else: $p = [0, \dots, 0, 1, 0, \dots, 0]$ and our computed probability is q , then the cross entropy is:

$$H(p, q) = - \sum_{c=1}^C p(c) \log q(c)$$

↓ Ground label.
↓ Product through softmax.
↓ Predict \hat{y}_k .

- Because of one-hot p , the only term left is the negative log probability of the true class

$$\log p(y|x) = \log \left(\frac{\exp(f_y)}{\sum_{c=1}^C \exp(f_c)} \right)$$

Sidenote: The KL divergence

7

We don't have that much time today !! YES !!

- Cross-entropy can be re-written in terms of the entropy and *Kullback-Leibler* divergence between the two distributions:

$$H(p, q) = H(p) + D_{KL}(p||q)$$

- Because $H(p)$ is zero in our case (and even if it wasn't it would be fixed and have no contribution to gradient), to minimize this is equal to minimizing the KL divergence between p and q
- The KL divergence is **not a distance** but a non-symmetric measure of the difference between two probability distributions p and q

$$D_{KL}(p||q) = \sum_{c=1}^C p(c) \log \frac{p(c)}{q(c)}$$

026 Entropy

www.cs.rochester.edu/u/james/CSC248/Lec6.pdf

"Classification의 loss function으로는 cross entropy를 주로 사용합니다."
"두 분포 차이를 줄이기 위해 KL-divergence를 최소화 시킵니다."

- **Entropy** : 정보를 최적으로 인코딩하기 위해 필요한 bit의 수

ex) 오늘이 무슨 요일인지 bit로 전송한다. → 3bit 필요
월: 000, 화: 001, 수: 010, 목: 011, 금: 100, 토: 101, 일: 110 $\log_2 N$

- 만약 각각의 발생 확률이 다르다면?

ex) 40개의 문자 (A,B,C,D,...,Z,1,2,3,...14)를 bit로 전송한다. $\log_2 40 = 5.3$

그런데 A,B,C,D가 전체의 22.5%씩 전체 90% 확률로 발생한다

1st bit : A,B,C,D 인지 아닌지 YES → 추가로 2bit 더 필요 (ABCD 구분)

NO → 추가로 6bit 더 필요 ($\log_2 36$)

필요한 bit 수 = $0.9 * 3 + 0.1 * 6 = 3.3$ bit!

- Entropy는 각 label들의 확률분포의 함수!

$$H(y) = \sum_i y_i \log \frac{1}{y_i} = - \sum_i y_i \log y_i$$

$$\begin{aligned} H &= -4 * (0.225 * \log_2 0.225) \\ &\quad - 36 * (0.0028 * \log_2 0.0028) \\ &= 2.72 \text{ bit} \end{aligned}$$

027 Cross entropy & K-L Divergence

- Entropy

$$H(y) = \sum_i y_i \log \frac{1}{y_i} = - \sum_i y_i \log y_i$$

↑
획득 information gain

- N개의 구분을 위해 $\log_2 N$ 의 bit가 필요하다
- 빈번한 (Prob. ↑) 정보의 bit수를 줄이자

- Cross Entropy

$$H(p, q) = - \sum_x p(x) \log q(x)$$

- K-L divergence

$$\begin{aligned} D_{\text{KL}}(P \| Q) &= - \sum_x p(x) \log q(x) + \sum_x p(x) \log p(x) \\ &= H(P, Q) - H(P) \end{aligned}$$

- Why Cross Entropy minimization?

$$C = -\frac{1}{n} \sum_x [y \ln a + (1 - y) \ln(1 - a)]$$

$$\begin{aligned} (\text{MSE}) \quad C &= \frac{(y - a)^2}{2} & a &= \sigma(z) \\ z &= wx + b & \frac{\partial C}{\partial w} &= (a - y)\sigma'(z)x \end{aligned}$$

$$\begin{aligned} \frac{\partial C}{\partial w_j} &= -\frac{1}{n} \sum_x \left(\frac{y}{\sigma(z)} - \frac{(1 - y)}{1 - \sigma(z)} \right) \frac{\partial \sigma}{\partial w_j} \\ &= -\frac{1}{n} \sum_x \left(\frac{y}{\sigma(z)} - \frac{(1 - y)}{1 - \sigma(z)} \right) \sigma'(z)x_j \end{aligned}$$

$$= \frac{1}{n} \sum_x \frac{\sigma'(z)x_j}{\sigma(z)(1 - \sigma(z))} (\sigma(z) - y)$$

sigmoid

$$\frac{\partial C}{\partial w_j} = \frac{1}{n} \sum_x x_j (\sigma(z) - y)$$

$\sigma(z)(1 - \sigma(z))$

Classification over a full dataset

- Cross entropy loss function over full dataset $\{x_i, y_i\}_{i=1}^N$

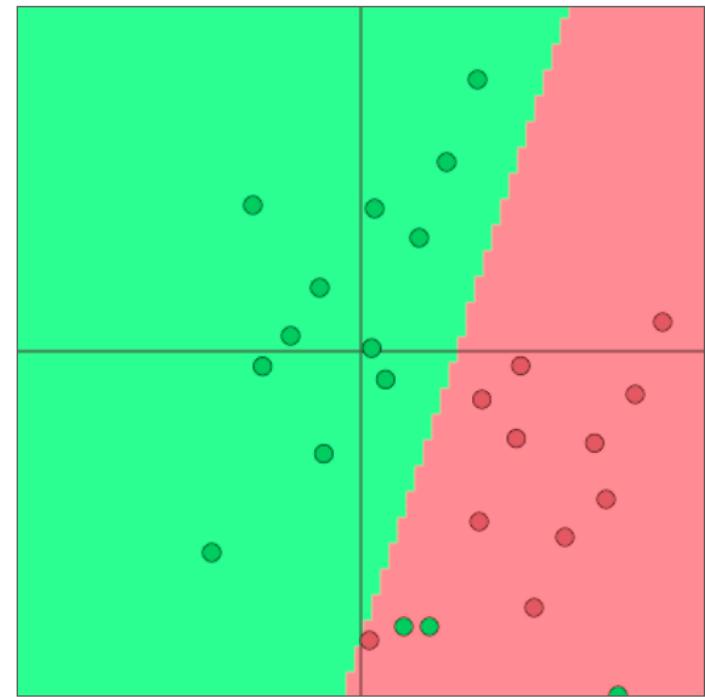
$$J(\theta) = \frac{1}{N} \sum_{i=1}^N -\log \left(\frac{e^{f_{y_i}}}{\sum_{c=1}^C e^{f_c}} \right)$$

what is P ? ... f ...

- Instead of

$$f_y = f_y(x) = W_y \cdot x = \sum_{j=1}^d W_{yj} x_j$$

- We will write f in matrix notation: $f = Wx$
- We can still index elements of it based on class

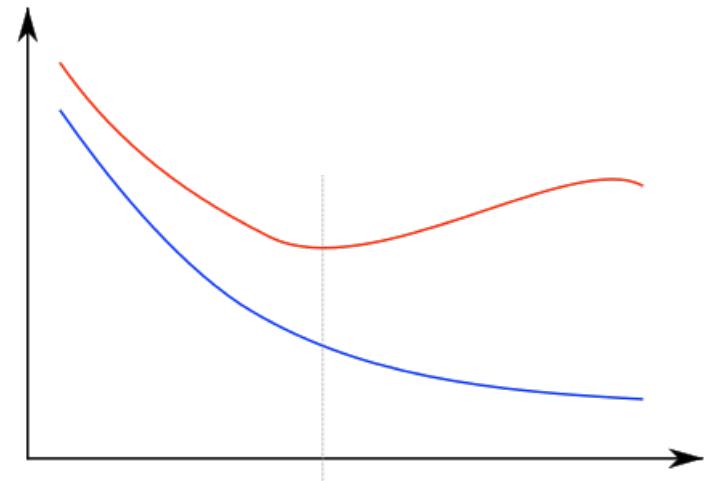


Classification: Regularization!

- Really full loss function over any dataset includes **regularization** over all parameters θ :

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N -\log \left(\frac{e^{f_{y_i}}}{\sum_{c=1}^C e^{f_c}} \right) + \lambda \sum_k \theta_k^2$$

- Regularization will prevent overfitting when we have a lot of features (or later a very powerful/deep model)
 - x-axis: more powerful model or more training iterations
 - Blue: training error, red: test error



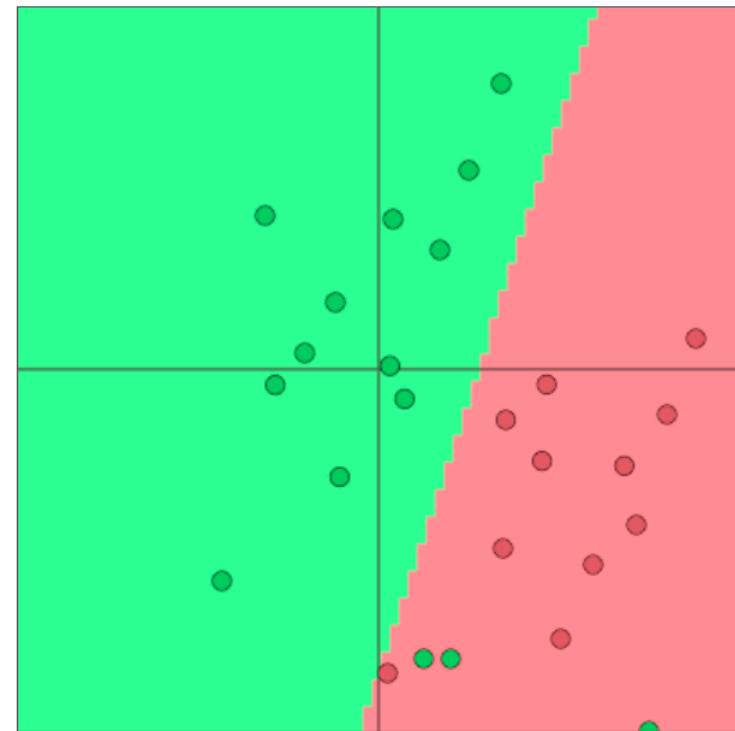
Details: General ML optimization

- For general machine learning θ usually only consists of columns of W :

$$\theta = \begin{bmatrix} W_{\cdot 1} \\ \vdots \\ W_{\cdot d} \end{bmatrix} = W(:) \in \mathbb{R}^{Cd}$$

- So we only update the decision boundary

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \nabla_{W_{\cdot 1}} \\ \vdots \\ \nabla_{W_{\cdot d}} \end{bmatrix} \in \mathbb{R}^{Cd}$$



Visualizations with ConvNetJS by Karpathy

Classification difference with word vectors



- Common in deep learning:
 - Learn both W and word vectors x

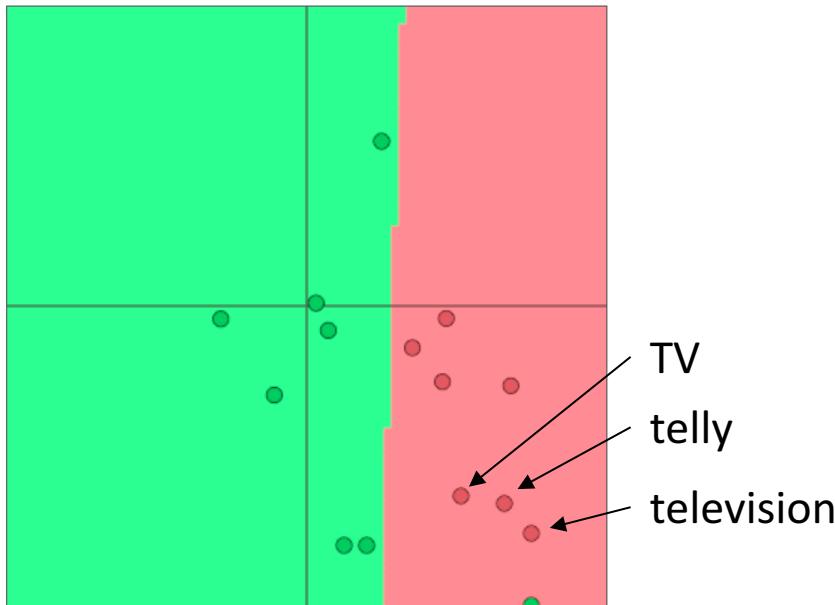
$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \nabla_{W_{\cdot 1}} \\ \vdots \\ \nabla_{W_{\cdot d}} \\ \nabla_{x_{aardvark}} \\ \vdots \\ \vdots \\ \nabla_{x_{zebra}} \end{bmatrix} \in \mathbb{R}^{Cd + Vd}$$

Very large!

Overfitting Danger!

A pitfall when retraining word vectors

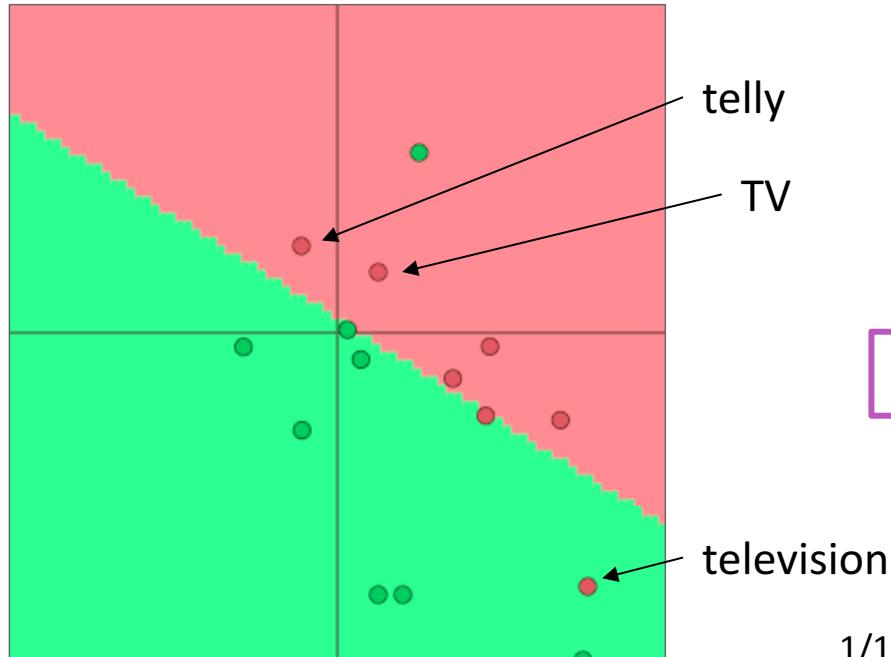
- Setting: We are training a logistic regression classification model for movie review sentiment using single words.
- In the **training data** we have “TV” and “telly”
- In the **testing data** we have “television”
- The **pre-trained word vectors** have all three similar:



- Question: What happens when we retrain the word vectors?

A pitfall when retraining word vectors

- Question: What happens when we train the word vectors?
- Answer:
 - Those that are **in** the training data **move around**
 - “TV” and “telly”
 - Words **not** in the training data **stay**
 - “television”



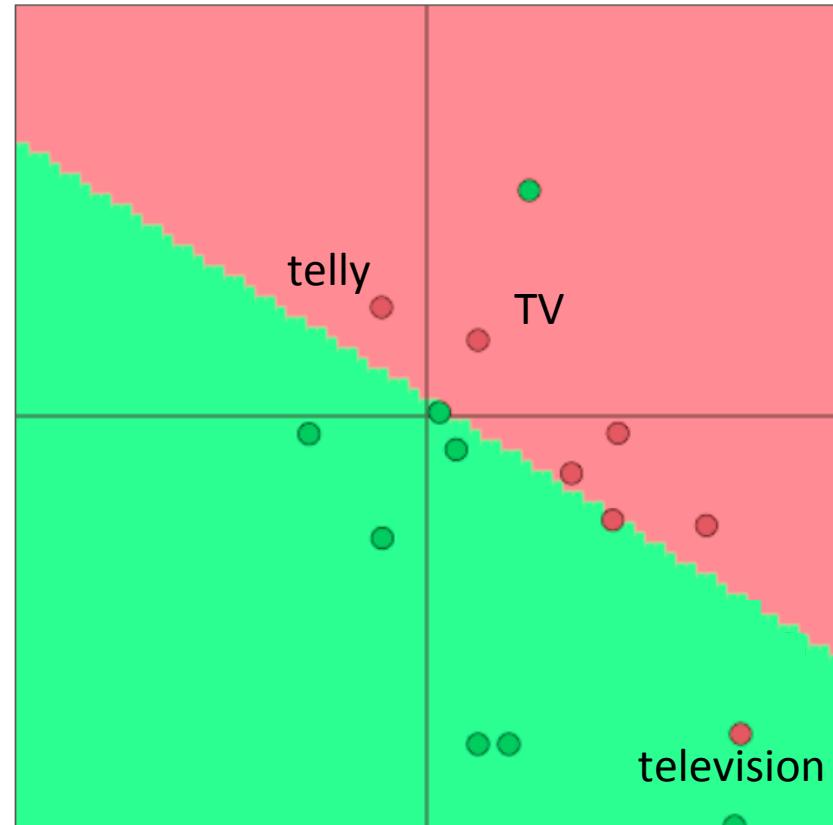
Losing generalization by re-training word vectors

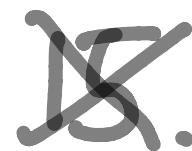
14

- Take home message:

If you only have a small training data set, don't train the word vectors.

If you have have a very large dataset, it may work better to train word vectors to the task.





Side note on word vectors notation

- The word vector matrix L is also called lookup table
- Word vectors = word embeddings = word representations (mostly)
- Mostly from methods like word2vec or Glove

$$L = d \begin{bmatrix} \vdots & \vdots & \cdots & \vdots & \vdots \\ \vdots & \vdots & \cdots & \vdots & \vdots \\ \vdots & \vdots & \cdots & \vdots & \vdots \\ \text{aardvark} & \text{a} & \dots & \text{meta} & \dots & \text{zebra} \end{bmatrix}$$

- These are the word features x_{word} from now on
- New development (later in the class): character models :o

Window classification

- Classifying single words is rarely done.
- Interesting problems like ambiguity arise in context!
- Example: auto-antonyms:
 - "To sanction" can mean "to permit" or "to punish."
 - "To seed" can mean "to place seeds" or "to remove seeds."
- Example: ambiguous named entities:
 - Paris → Paris, France vs Paris Hilton
 - Hathaway → Berkshire Hathaway vs Anne Hathaway

Window classification

- Idea: classify a word in its context window of neighboring words.
- For example named entity recognition into 4 classes:
 - Person, location, organization, none
- Many possibilities exist for classifying one word in context, e.g. averaging all the words in a window but that loses position information

Window classification

/8

- Train softmax classifier to classify a center word by taking concatenation of all word vectors surrounding it
 - Example: Classify “Paris” in the context of this sentence with window length 2:

... museums in Paris are amazing



$$X_{\text{window}} = [x_{\text{museums}} \quad x_{\text{in}} \quad x_{\text{Paris}} \quad x_{\text{are}} \quad x_{\text{amazing}}]^T$$

- Resulting vector $x_{\text{window}} = x \in \mathbb{R}^{5d}$, a column vector!

Simplest window classifier: Softmax

- With $x = x_{\text{window}}$ we can use the same softmax classifier as before

predicted model output probability $\hat{y}_y = p(y|x) = \frac{\exp(W_y \cdot x)}{\sum_{c=1}^C \exp(W_c \cdot x)}$

 예전 x_{window}
 텐서 \sim

- With cross entropy error as before:

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N -\log \left(\frac{e^{f_{y_i}}}{\sum_{c=1}^C e^{f_c}} \right)$$

- But how do you update the word vectors?

Updating concatenated word vectors

20

- Short answer: Just take derivatives as before
- Long answer: Let's go over steps together (helpful for PSet 1)
- Define:
 - \hat{y} : softmax probability output vector (see previous slide)
 - t : target probability distribution (all 0's except at ground truth index of class y , where it's 1)
 - $f = f(x) = Wx \in \mathbb{R}^C$ and $f_c = c^{\text{'}}\text{th element of the } f \text{ vector}$
Label A 은 A에 속한 4 Unnormalized score
- Hard, the first time, hence some tips now :)

$$W \in \mathbb{R}^{C \times D} \quad x \in \mathbb{R}^D \quad f = Wx \in \mathbb{R}^C$$

Normalized score = probability

Softmax

Updating concatenated word vectors

- Tip 1: Carefully define your variables and keep track of their dimensionality!

$$\begin{aligned} f &= f(x) = Wx \in \mathbb{R}^C \\ \hat{y} &\quad t \\ W &\in \mathbb{R}^{C \times 5d} \end{aligned}$$

- Tip 2: Chain rule! If $y = f(u)$ and $u = g(x)$, i.e. $y = f(g(x))$, then:
- Simple example: $\frac{dy}{dx} = \frac{dy}{du} \frac{du}{dx} = \frac{df(u)}{du} \frac{dg(x)}{dx}$

$$\frac{dy}{dx} = \frac{d}{dx} 5(x^3 + 7)^4$$

$$\begin{aligned} y &= f(u) = 5u^4 & u &= g(x) = x^3 + 7 \\ \frac{dy}{du} &= 20u^3 & \frac{du}{dx} &= 3x^2 \end{aligned}$$

$$\frac{dy}{dx} = 20(x^3 + 7)^3 \cdot 3x^2$$

Updating concatenated word vectors

22

$$f = f(x) = Wx \in \mathbb{R}^C$$
$$\hat{y} \quad t \quad W \in \mathbb{R}^{C \times 5d}$$

- Tip 2 continued: **Know thy chain rule**
- Don't forget which variables depend on what and that x appears inside all elements of f 's

$$\frac{\partial}{\partial x} - \log \text{softmax}(f_y(x)) = \sum_{c=1}^C - \left[\frac{\partial \log \text{softmax}(f_y(x))}{\partial f_c} \right] \cdot \left[\frac{\partial f_c(x)}{\partial x} \right]$$

1. $\frac{\partial}{\partial x} - \log \text{softmax}(f_y(x))$ 2. $\frac{\partial \log \text{softmax}(f_y(x))}{\partial f_c}$ 3. $\frac{\partial f_c(x)}{\partial x}$

- Tip 3: For the softmax part of the derivative: First take the derivative wrt f_c when $c=y$ (the correct class), then take derivative wrt f_c when $c \neq y$ (all the incorrect classes)

Updating concatenated word vectors

- Tip 4: When you take derivative wrt one element of f , try to see if you can create a gradient in the end that includes all partial derivatives:

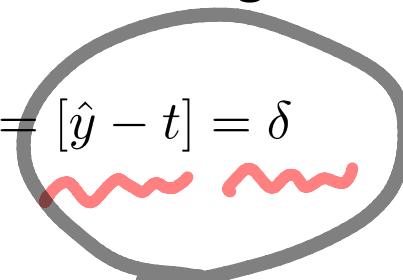
$$\frac{\partial}{\partial f} - \log \text{softmax}(f_y) =$$

$$\begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_y - 1 \\ \vdots \\ \hat{y}_C \end{bmatrix}$$

$$\hat{y} \quad t \\ f = f(x) = Wx \in \mathbb{R}^C$$

- Tip 5: To later not go insane & implementation! → results in terms of vector operations and define single indexable vectors:

$$\frac{\partial}{\partial f} - \log \text{softmax}(f_y) = [\hat{y} - t] = \delta$$



Updating concatenated word vectors

24

- Tip 6: When you start with the chain rule, first use explicit sums and look at partial derivatives of e.g. x_i or W_{ij}

$$\hat{y} \quad t \\ f = f(x) = Wx \in \mathbb{R}^C$$

$$\sum_{c=1}^C -\frac{\partial \log \text{softmax}(f_y(x))}{\partial f_c} \cdot \frac{\partial f_c(x)}{\partial x} = \sum_{c=1}^C \delta_c W_c^T$$


- Tip 7: To clean it up for even more complex functions later: Know dimensionality of variables & simplify into matrix notation

$$\frac{\partial}{\partial x} - \log p(y|x) = \sum_{c=1}^C \delta_c W_c^T = W^T \delta$$


- Tip 8: Write this out in full sums if it's not clear!

Updating concatenated word vectors

- What is the dimensionality of the window vector gradient?

$$\frac{\partial}{\partial \mathbf{x}} - \log p(y|x) = \sum_{c=1}^C \delta_c W_c. = W^T \delta$$

$\xrightarrow{\text{Window}} \text{Window} = [x_{\text{museum}} \quad x_{\text{in}} \quad x_{\text{parts}} \quad x_{\text{is}} \quad x_{\text{amazing}}]$

- x is the entire window, 5 d-dimensional word vectors, so the derivative wrt to x has to have the same dimensionality:

$$\nabla_x J = W^T \delta \in \mathbb{R}^{5d}$$

Updating concatenated word vectors

- The gradient that arrives at and updates the word vectors can simply be split up for each word vector:
- Let $\nabla_x J = W^T \delta = \delta_{x_{window}}$
- With $x_{window} = [x_{museums} \quad x_{in} \quad x_{Paris} \quad x_{are} \quad x_{amazing}]$
- We have

$$\delta_{window} = \begin{bmatrix} \nabla x_{museums} \\ \nabla x_{in} \\ \nabla x_{Paris} \\ \nabla x_{are} \\ \nabla x_{amazing} \end{bmatrix} \in \mathbb{R}^{5d}$$

각 word vector의 Gradient .

Updating concatenated word vectors

- This will push word vectors into areas such they will be helpful in determining named entities.
- For example, the model can learn that seeing x_{in} as the word just before the center word is indicative for the center word to be a location

What's missing for training the window model? 28



- The gradient of J wrt the softmax weights W !
- Similar steps, write down partial wrt W_{ij} first!
- Then we have full

이해못함 T^T

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \nabla_{W_{\cdot 1}} \\ \vdots \\ \nabla_{W_{\cdot d}} \\ \nabla_{x_{aardvark}} \\ \vdots \\ \nabla_{x_{zebra}} \end{bmatrix} \in \mathbb{R}^{Cd+Vd}$$

A note on matrix implementations

29

- There are two expensive operations in the softmax:
- The matrix multiplication $f = Wx$ and the exp
- A for loop is never as efficient when you implement it compared to a large matrix multiplication!
- Example code →

A note on matrix implementations



- Looping over word vectors instead of concatenating them all into one large matrix and then multiplying the softmax weights with that matrix

```
from numpy import random
N = 500 # number of windows to classify
d = 300 # dimensionality of each window
C = 5 # number of classes
W = random.rand(C,d)
wordvectors_list = [random.rand(d,1) for i in range(N)]
wordvectors_one_matrix = random.rand(d,N)

%timeit [W.dot(wordvectors_list[i]) for i in range(N)]
%timeit W.dot(wordvectors_one_matrix)
```

- 1000 loops, best of 3: 639 μ s per loop
10000 loops, best of 3: 53.8 μ s per loop

A note on matrix implementations



```
from numpy import random
N = 500 # number of windows to classify
d = 300 # dimensionality of each window
C = 5 # number of classes
W = random.rand(C,d)
wordvectors_list = [random.rand(d,1) for i in range(N)]
wordvectors_one_matrix = random.rand(d,N)

%timeit [W.dot(wordvectors_list[i]) for i in range(N)]
%timeit W.dot(wordvectors_one_matrix)
```

- Result of faster method is a C x N matrix:
 - Each column is an $f(x)$ in our notation (unnormalized class scores)
- Matrices are awesome!
- You should speed test your code a lot too

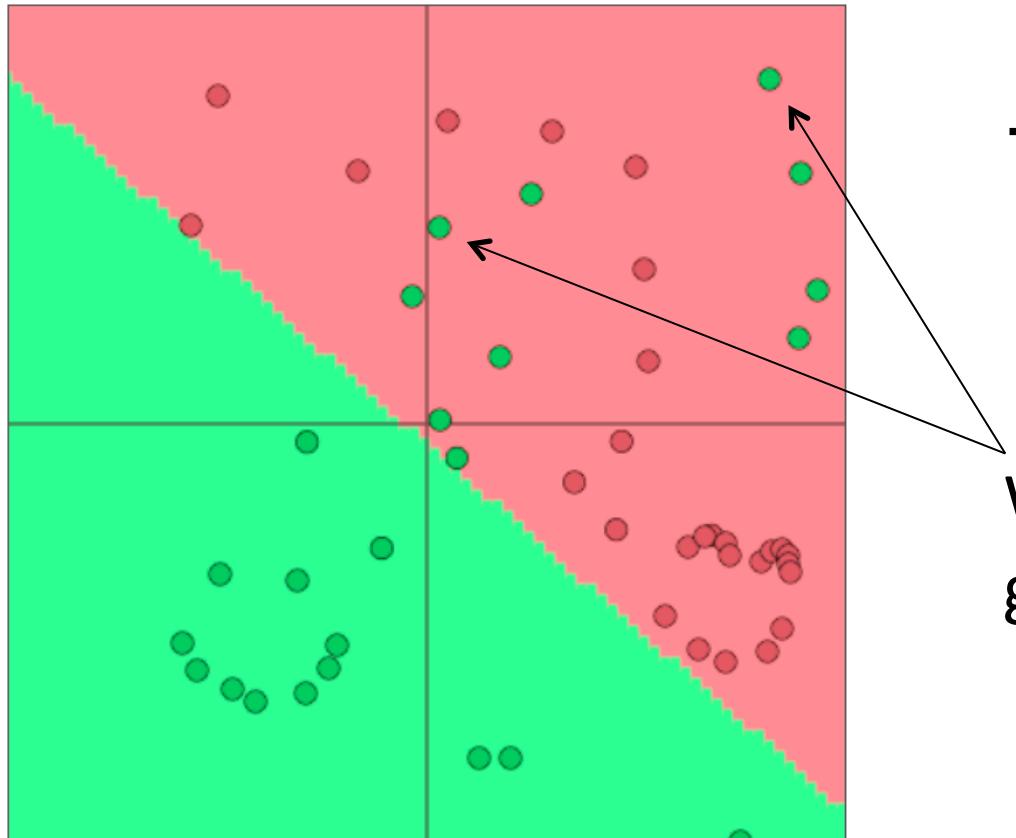
Softmax (= logistic regression) alone not very powerful

- Softmax only gives linear decision boundaries in the original space.
- With little data that can be a good regularizer
- With more data it is very limiting!

Softmax (= logistic regression) is not very powerful

33

- Softmax only linear decision boundaries



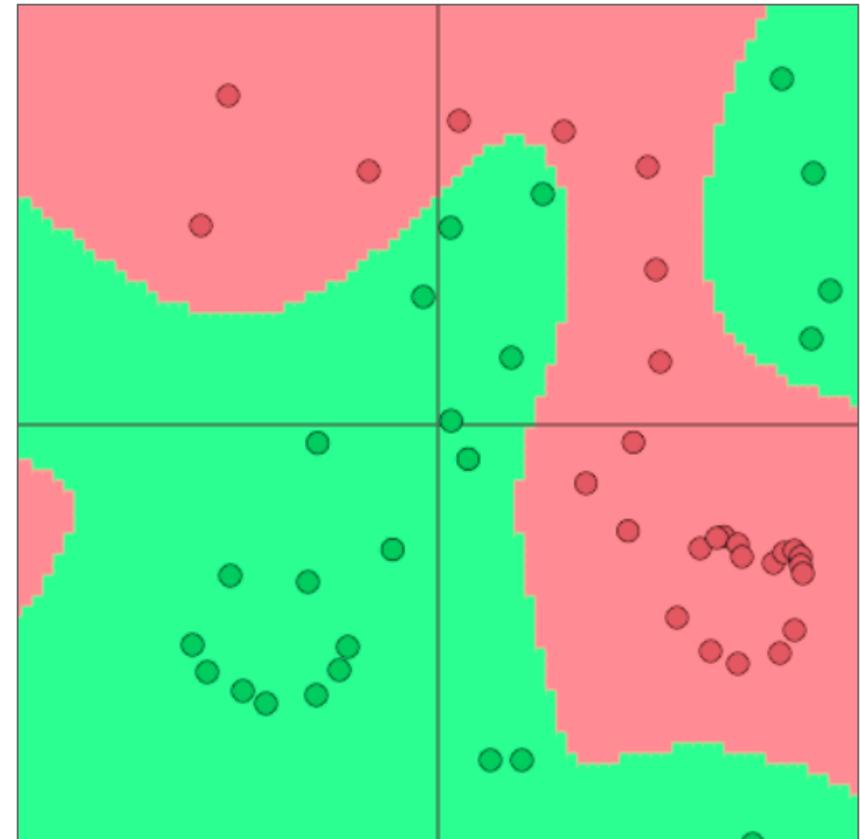
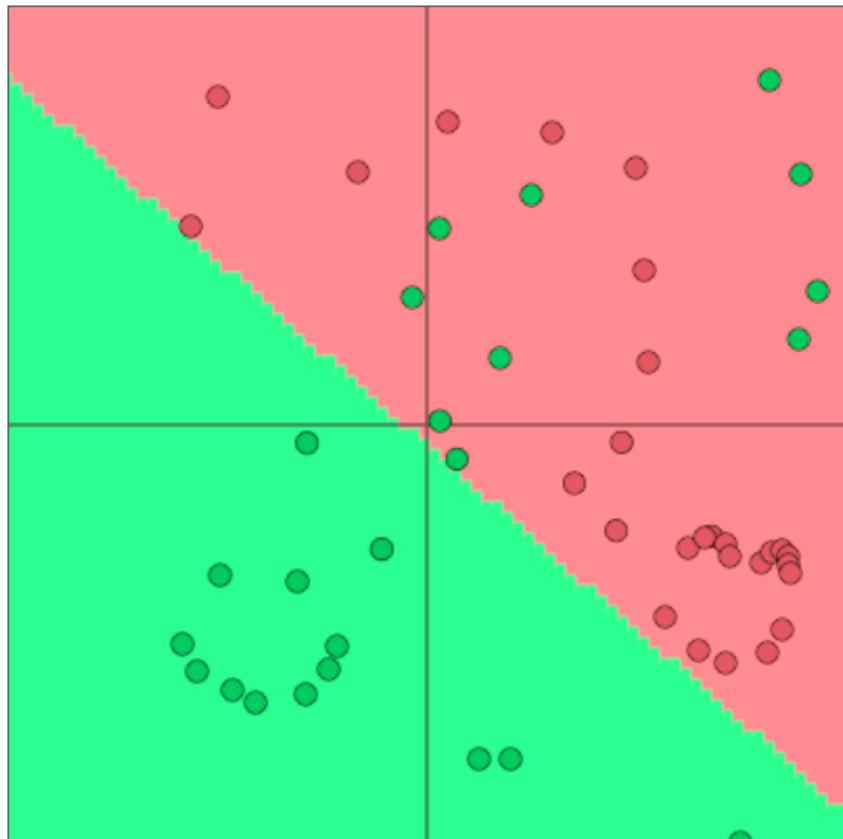
→ Lame when problem
is complex
~~이거 예상할 수 있다 ~~~

Wouldn't it be cool to
get these correct?

Neural Nets for the Win!

34

- Neural networks can learn much more complex functions and nonlinear decision boundaries!



From logistic regression to neural nets

Demystifying neural networks

35

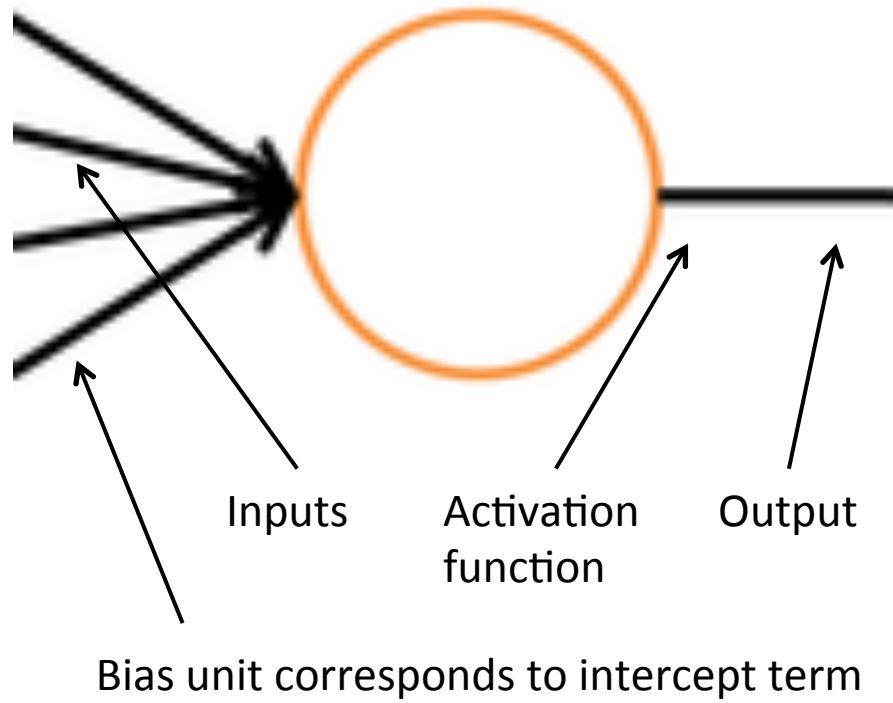
정리를 알아볼까요...!

Neural networks come with their own terminological baggage

But if you understand how softmax models work

Then **you already understand** the operation of a basic neuron!

A single neuron
A computational unit with n (3) inputs and 1 output and parameters W, b



A neuron is essentially a binary logistic regression unit

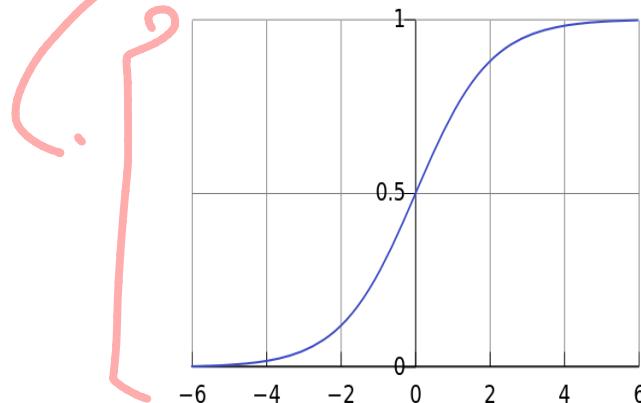
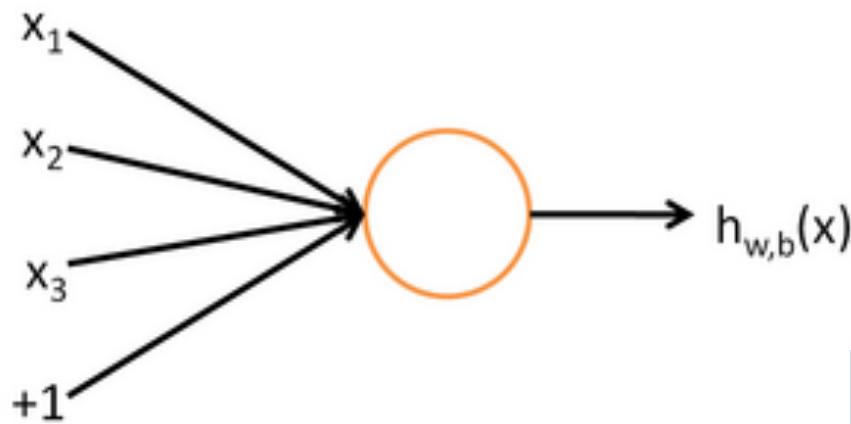
$$\boxed{w} \boxed{x} + \boxed{b} \Rightarrow \text{activation function } f \Rightarrow h_{w,b}(x)$$

sigmoid

$$h_{w,b}(x) = f(w^T x + b)$$

b : We can have an “always on” feature, which gives a class prior, or separate it out, as a bias term

$$f(z) = \frac{1}{1 + e^{-z}}$$

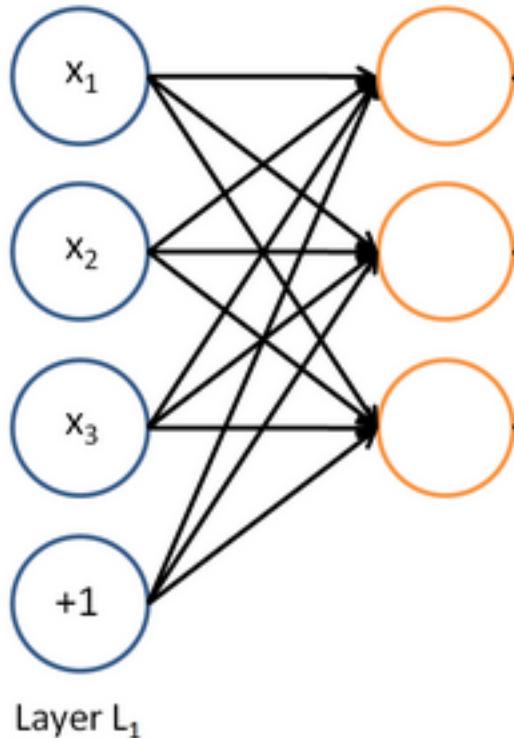


w, b are the parameters of this neuron i.e., this logistic regression model

A neural network

= running several logistic regressions at the same time

If we feed a vector of inputs through a bunch of logistic regression functions, then we get a vector of outputs ...

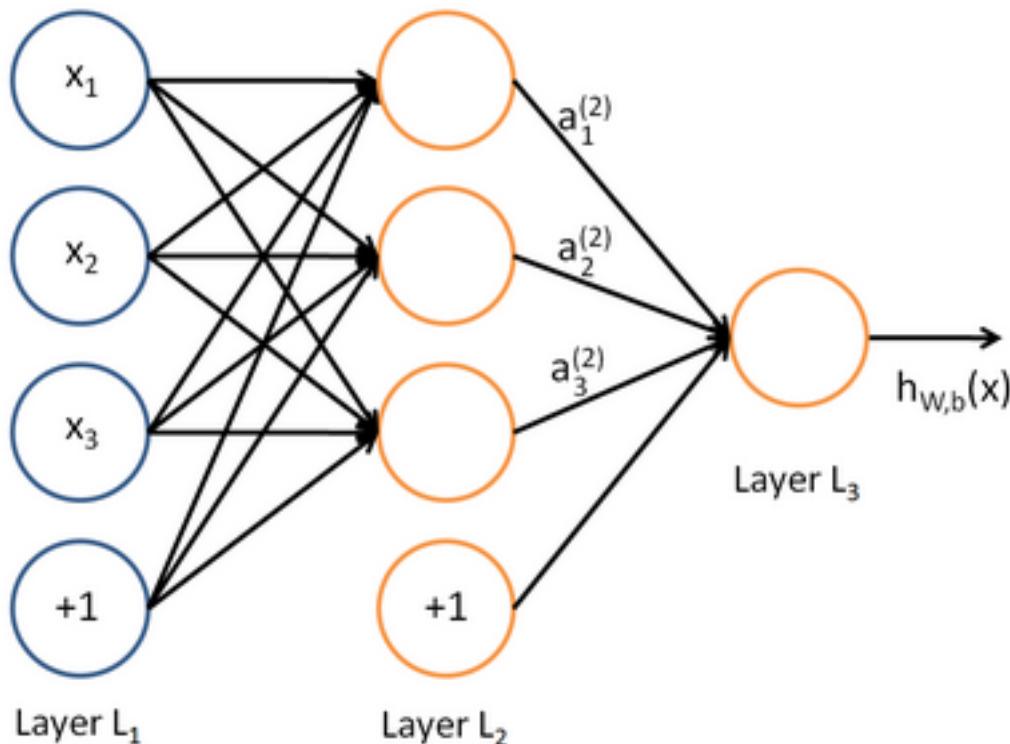


But we don't have to decide ahead of time what variables these logistic regressions are trying to predict!

A neural network

= running several logistic regressions at the same time

... which we can feed into another logistic regression function



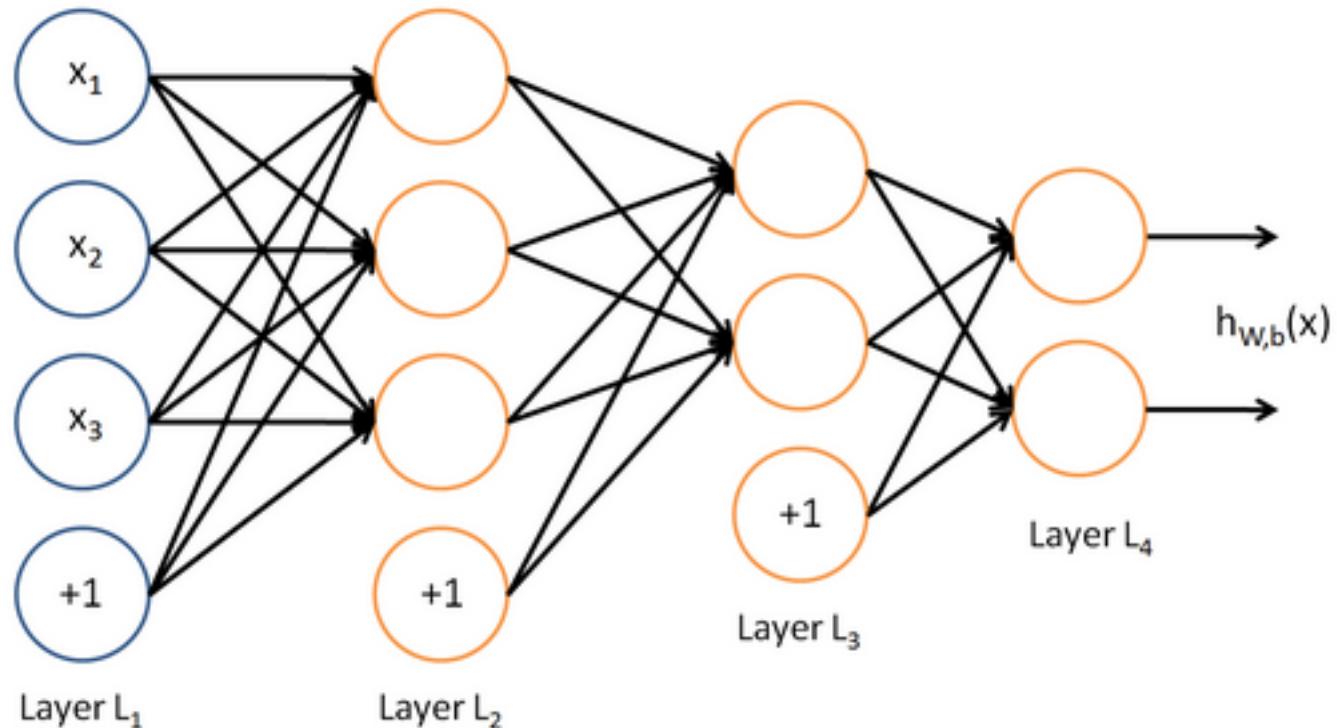
It is the loss function that will direct what the intermediate hidden variables should be, so as to do a good job at predicting the targets for the next layer, etc.

A neural network

39

= running several logistic regressions at the same time

Before we know it, we have a multilayer neural network....



Matrix notation for a layer

We have

$$a_1 = f(W_{11}x_1 + W_{12}x_2 + W_{13}x_3 + b_1)$$

$$a_2 = f(W_{21}x_1 + W_{22}x_2 + W_{23}x_3 + b_2)$$

etc.

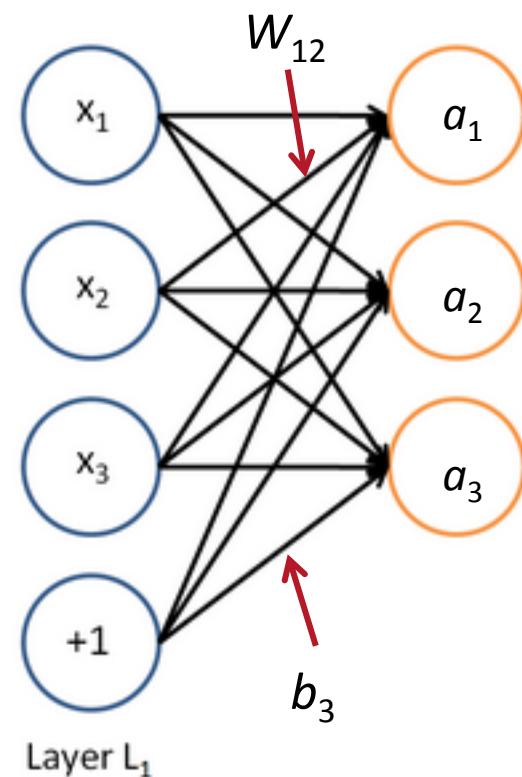
In matrix notation

$$z = Wx + b$$

$$a = f(z)$$

where f is applied element-wise:

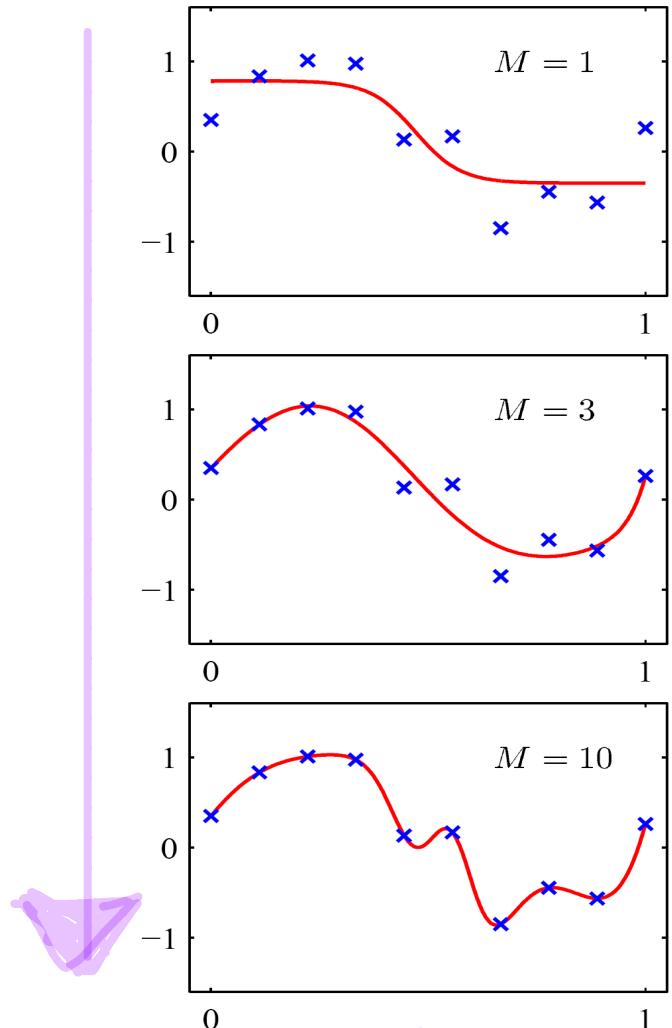
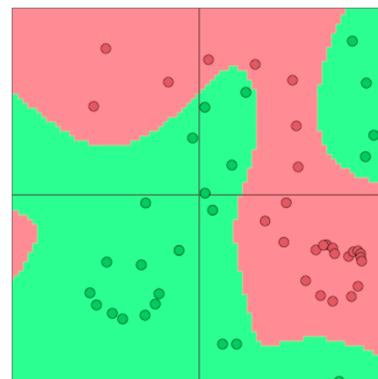
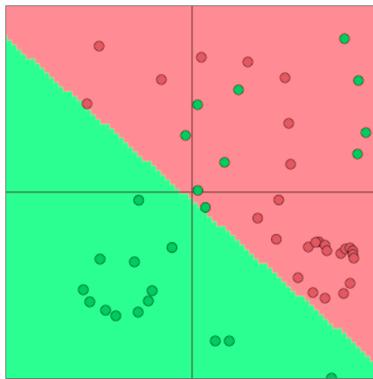
$$f([z_1, z_2, z_3]) = [f(z_1), f(z_2), f(z_3)]$$



41

Non-linearities (f): Why they're needed

- Example: function approximation, e.g., regression or classification
 - Without non-linearities, deep neural networks can't do anything more than a linear transform
 - Extra layers could just be compiled down into a single linear transform:
 $W_1 W_2 x = Wx \rightarrow \text{linear only} = \text{linear}$
 - With more layers, they can approximate more complex functions!



more non-linearities.
more fit to complex!

A more powerful, neural net window classifier 42

- Revisiting
- $X_{\text{window}} = [x_{\text{museums}} \quad x_{\text{in}} \quad x_{\text{Paris}} \quad x_{\text{are}} \quad x_{\text{amazing}}]$
- Assume we want to classify whether the center word is a
location or not

⇒ Logistic.

Binary classification with unnormalized scores

- Revisiting our previous example:
 $X_{\text{window}} = [x_{\text{museums}} \quad x_{\text{in}} \quad x_{\text{Paris}} \quad x_{\text{are}} \quad x_{\text{amazing}}]$
 - Assume we want to classify whether the center word is a Location (Named Entity Recognition)
 - Similar to word2vec, we will go over all positions in a corpus. But this time, it will be supervised and only some positions should get a high score.
 - The positions that have an actual NER location in their center are called “true” positions.

A Single Layer Neural Network

- A single layer is a combination of a linear layer and a nonlinearity:

$$\begin{aligned} z &= Wx + b \\ a &= f(z) \end{aligned}$$

- The neural activations a can then be used to compute some output

- For instance, a probability via softmax

$$P(y|x) = \text{softmax}(Wa)$$

- Or an unnormalized score (even simpler)

$$\text{score}(x) = U^T a \in \mathbb{R}$$

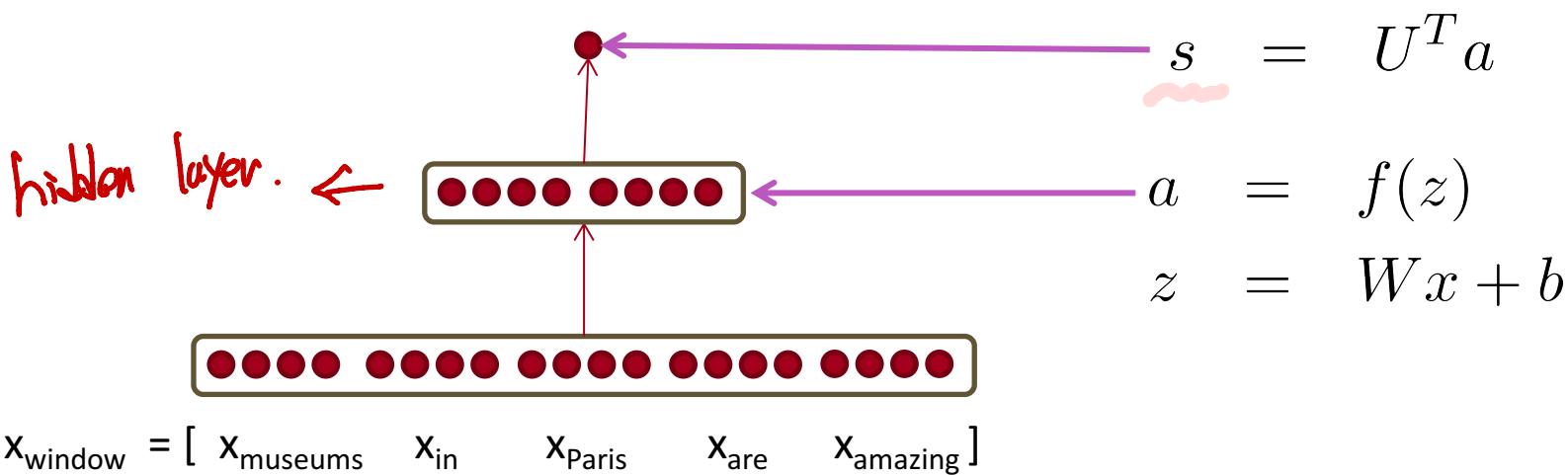
Summary: Feed-forward Computation

We compute a window's **score** with a 3-layer neural net:

- $s = \text{score}(\text{"museums in Paris are amazing"})$

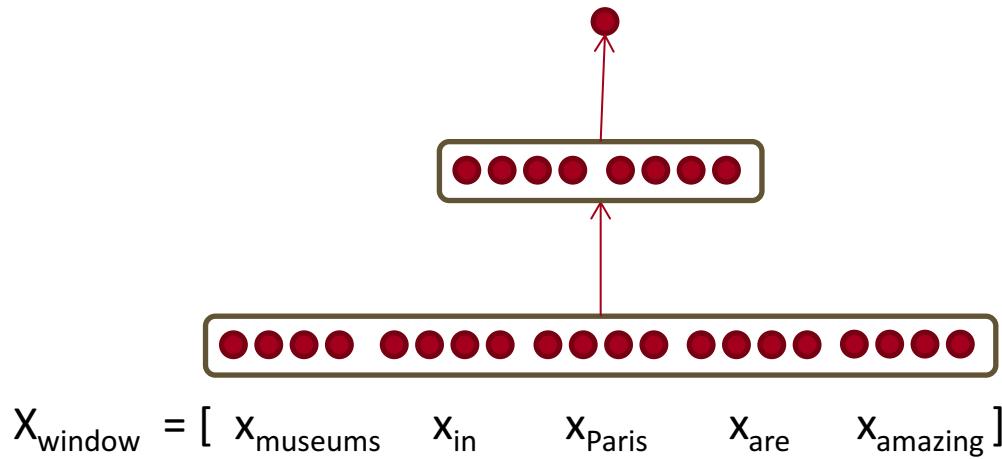
$$s = U^T f(Wx + b)$$

$$x \in \mathbb{R}^{20 \times 1}, W \in \mathbb{R}^{8 \times 20}, U \in \mathbb{R}^{8 \times 1}$$



Main intuition for extra layer

The layer learns **non-linear interactions** between the input word vectors.

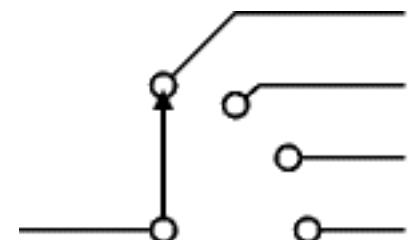


Example: only if “*museums*” is first vector should it matter that “*in*” is in the second position

The max-margin loss

- s = score(museums in Paris are amazing)
- s_c = score(Not all museums in Paris)
- Idea for training objective: make score of true window larger and corrupt window's score lower (until they're good enough): minimize

$$J = \max(0, 1 - s + s_c)$$



- This is continuous --> we can use SGD

Max-margin Objective function

- Objective for a single window:

$$J = \max(0, 1 - s + s_c)$$

- Each window with a location at its center should have a score +1 higher than any window without a location at its center
- 
- For full objective function: Sample several corrupt windows per true one. Sum over all training windows

Training with Backpropagation

$$J = \max(0, 1 - s + s_c)$$

$$\begin{aligned}s &= U^T f(Wx + b) \\ s_c &= U^T f(Wx_c + b)\end{aligned}$$

Assuming cost J is > 0 ,
compute the derivatives of s and s_c wrt all the
involved variables: U, W, b, x

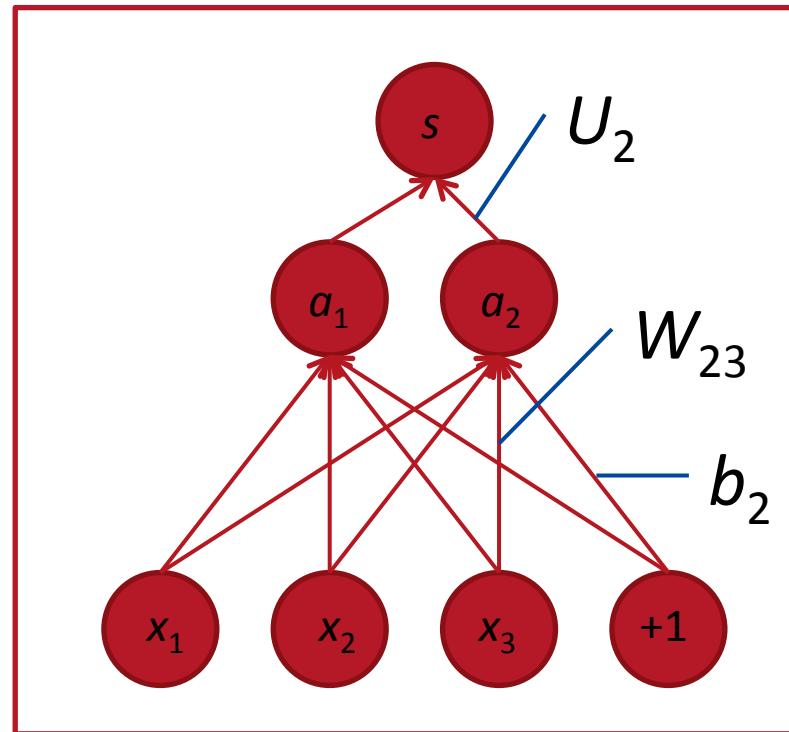
$$\frac{\partial s}{\partial U} = \frac{\partial}{\partial U} U^T a \qquad \frac{\partial s}{\partial U} = a$$

Training with Backpropagation

- Let's consider the derivative of a single weight W_{ij}

$$\frac{\partial s}{\partial W} = \frac{\partial}{\partial W} U^T a = \frac{\partial}{\partial W} U^T f(z) = \frac{\partial}{\partial W} U^T f(Wx + b)$$

- This only appears inside a_i
- For example: W_{23} is only used to compute a_2



Training with Backpropagation

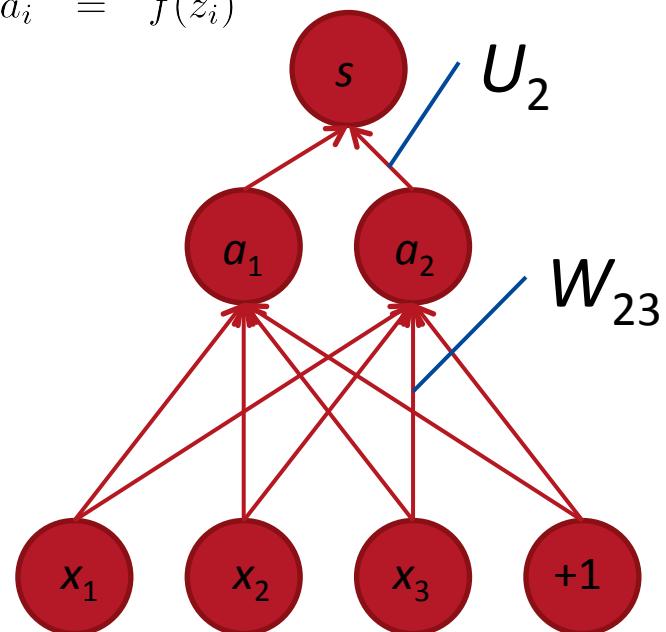
$$\frac{\partial s}{\partial W} = \frac{\partial}{\partial W} U^T a = \frac{\partial}{\partial W} U^T f(z) = \frac{\partial}{\partial W} U^T f(Wx + b)$$

Derivative of weight W_{ij} :

$$\frac{\partial}{\partial W_{ij}} U^T a \rightarrow \frac{\partial}{\partial W_{ij}} U_i a_i$$

$$\begin{aligned} U_i \frac{\partial}{\partial W_{ij}} a_i &= U_i \frac{\partial a_i}{\partial z_i} \frac{\partial z_i}{\partial W_{ij}} \\ &= U_i \frac{\partial f(z_i)}{\partial z_i} \frac{\partial z_i}{\partial W_{ij}} \\ &= U_i f'(z_i) \frac{\partial z_i}{\partial W_{ij}} \\ &= U_i f'(z_i) \frac{\partial W_i \cdot x + b_i}{\partial W_{ij}} \end{aligned}$$

$$\begin{aligned} z_i &= W_i \cdot x + b_i = \sum_{j=1}^3 W_{ij} x_j + b_i \\ a_i &= f(z_i) \end{aligned}$$



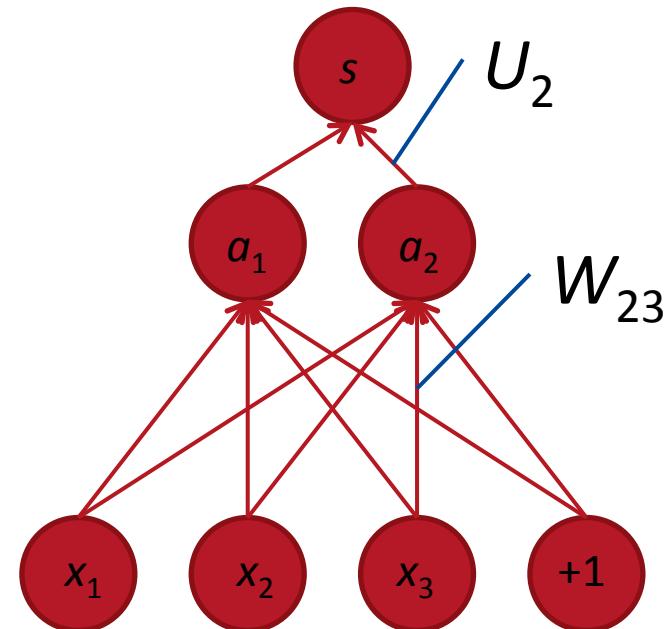
Training with Backpropagation

Derivative of single weight W_{ij} :

$$\begin{aligned}
 U_i \frac{\partial}{\partial W_{ij}} a_i &= U_i f'(z_i) \frac{\partial W_i \cdot x + b_i}{\partial W_{ij}} \\
 &= U_i f'(z_i) \frac{\partial}{\partial W_{ij}} \sum_k W_{ik} x_k \\
 &= \underbrace{U_i f'(z_i)}_{\delta_i} x_j \\
 &= \underbrace{\delta_i}_{\text{Local error signal}} \underbrace{x_j}_{\text{Local input signal}}
 \end{aligned}$$

where $f'(z) = f(z)(1 - f(z))$ for logistic f

$$\begin{aligned}
 z_i &= W_i \cdot x + b_i = \sum_{j=1}^3 W_{ij} x_j + b_i \\
 a_i &= f(z_i)
 \end{aligned}$$



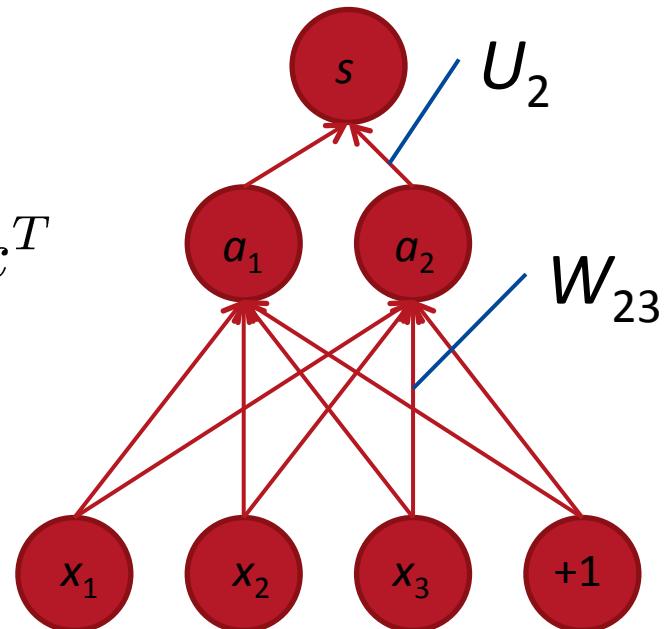
Training with Backpropagation

- From single weight W_{ij} to full W :

$$\begin{aligned}\frac{\partial s}{\partial W_{ij}} &= \underbrace{U_i f'(z_i)}_{\delta_i} x_j \\ &= \delta_i \quad x_j\end{aligned}$$

$$\begin{aligned}z_i &= W_i \cdot x + b_i = \sum_{j=1}^3 W_{ij} x_j + b_i \\ a_i &= f(z_i)\end{aligned}$$

- We want all combinations of $i = 1, 2$ and $j = 1, 2, 3 \rightarrow ?$
- Solution: Outer product: $\frac{\partial s}{\partial W} = \delta x^T$ where $\delta \in \mathbb{R}^{2 \times 1}$ is the “responsibility” or error signal coming from each activation a

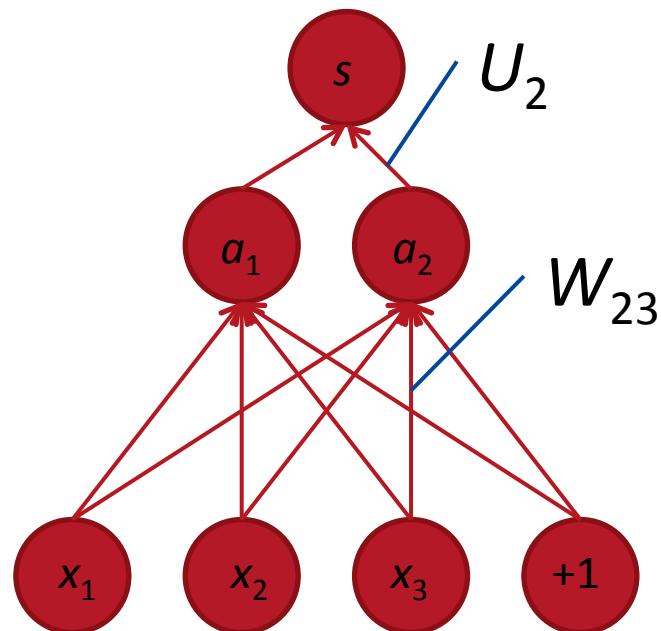


Training with Backpropagation

- For biases b , we get:

$$\begin{aligned} z_i &= W_i \cdot x + b_i = \sum_{j=1}^3 W_{ij} x_j + b_i \\ a_i &= f(z_i) \end{aligned}$$

$$\begin{aligned} & U_i \frac{\partial}{\partial b_i} a_i \\ = & U_i f'(z_i) \frac{\partial W_i \cdot x + b_i}{\partial b_i} \\ = & \delta_i \end{aligned}$$



Training with Backpropagation

That's almost backpropagation

It's taking derivatives and using the chain rule

Remaining trick: we can **re-use** derivatives computed for higher layers in computing derivatives for lower layers!

Example: last derivatives of model, the word vectors in x

Training with Backpropagation

- Take derivative of score with respect to single element of word vector
- Now, we cannot just take into consideration one a_i , because each x_j is connected to all the neurons above and hence x_j influences the overall score through all of these, hence:

$$\begin{aligned}\frac{\partial s}{\partial x_j} &= \sum_{i=1}^2 \frac{\partial s}{\partial a_i} \frac{\partial a_i}{\partial x_j} \\ &= \sum_{i=1}^2 \frac{\partial U^T a}{\partial a_i} \frac{\partial a_i}{\partial x_j} \\ &= \sum_{i=1}^2 U_i \frac{\partial f(W_i \cdot x + b)}{\partial x_j} \\ &= \sum_{i=1}^2 \underbrace{U_i f'(W_i \cdot x + b)}_{\delta_i} \frac{\partial W_i \cdot x}{\partial x_j} \\ &= \sum_{i=1}^2 \delta_i W_{ij} \\ &= W_{\cdot j}^T \delta\end{aligned}$$

Re-used part of previous derivative

Training with Backpropagation

- With $\frac{\partial s}{\partial x_j} = W_{\cdot j}^T \delta$, what is the full gradient? \rightarrow

$$\frac{\partial s}{\partial x} = W^T \delta$$

- Observations: The error message δ that arrives at a hidden layer has the same dimensionality as that hidden layer

Putting all gradients together:

- Remember: Full objective function for each window was:

$$J = \max(0, 1 - s + s_c)$$

$$s = U^T f(Wx + b)$$
$$s_c = U^T f(Wx_c + b)$$

- For example: gradient for U :

$$\frac{\partial J}{\partial U} = 1\{1 - s + s_c > 0\} (-f(Wx + b) + f(Wx_c + b))$$

$$\frac{\partial J}{\partial U} = 1\{1 - s + s_c > 0\} (-a + a_c)$$

Summary

Congrats! Super useful basic components and real model

- Word vector training
- Windows
- Softmax and cross entropy error → PSet1
- Scores and max-margin loss
- Neural network → PSet1

One more half of a math-heavy lecture

Then the rest will be easier and more applied :)

Next lecture:

Project advice

Taking more and **deeper derivatives** → Full **Backprop**

Then we have all the basic tools in place to learn about more complex models and have some fun :)