

C224N: Stanford

Chris Manning:

## Lecture 1: Introduction and Word Vectors

- An understanding of the effective modern methods for deep learning
- Basics first, then key methods used in NLP:
  - Recurrent networks, attention, etc.
- A big picture understanding of human language and the difficulties in understanding and producing them
- An understanding of an ability to build systems for some of the major problems in NLP
  - word meaning, dependency parsing, machine translation, question answering
- Language isn't a formal system. It's a glorious chaos
- Humans are better because of the language and human communication network approach
- Manning vs LeCun (Humans are better)
  - Writing takes knowledge spatially throughout the word and temporally through time

1) How do you represent the meaning of a word?

Brigham and  
meaning

Definition: meaning

→ idea that is represented by a word, phrase etc.

→ idea that a person wants to express by using words, signs, etc

→ idea that is expressed in a work of writing, art etc.

Linguistic way of thinking of meaning:

Signifier (symbol)  $\leftrightarrow$  signified (idea of thing)

= denotational semantics

Common solutions to have usable meaning

Use Wordnet (nltk), a thesaurus containing lists of synonym sets and hypernyms

But that's not ~~terribly good at anything~~

problem: can't add pt

will specialize and refine

with a good library of books

## Representing words as discrete symbols

Traditional NLP (NLP till 2012) regards words as discrete symbols → localist rep.  
were represented as one-hot vectors  
⇒ but words are infinite.  
one possible solution to find the synonyms  
can be word similarity table

Distributional semantics: A word's meaning is given by the words that frequently appear close-by  
⇒ when a word  $w$  appears in a text, its context is the set of words that appear nearby (within a fixed-size window)  
Word Vectors: We will build a dense vector for each word, chosen so that it is similar to vectors of words that appear in similar contexts

Word vectors are sometimes called word embeddings or word representations. They are a distributed representation.

Word meaning as a neural word vector:

(Yoshua Bengio's work)

Word2Vec: (Mikolov et. al. 2013)

### Idea:

- i) We have a large corpus of text
- ii) Every word in a fixed vocabulary is represented by a vector
- iii) Go through each position  $t$  in the text, which has a center word  $c$  and context ("outside") words  $o$
- iv) Use the similarity of the word vectors for  $c$  and  $o$  to calculate the probability of  $o$  given  $c$  (or vice versa)  
$$P(w_{+j}/w_t)$$
- v) Keep adjusting the word vectors to maximise this probability

objection function (cost or loss)

for each position  $t = 1 \dots T$  predict context words within a window of fixed size  $m$ , give center word  $w_t$

$$\text{Likelihood} = L(\theta) = \prod_{t=1}^T \prod_{\substack{-m \leq j \leq m \\ j \neq 0}} P(w_{t+j} | w_t; \theta)$$

slide  $\rightarrow$  all variables to be optimized

The objective function  $J(\theta)$  is the (average) negative

log likelihood:

$$J(\theta) = -\frac{1}{T} \log L(\theta) = \theta \cdot \underbrace{\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}}}_{\substack{\text{minimizing} \\ \text{average}}} \log P(w_{t+j} | w_t; \theta)$$

Minimizing objective  $\Leftrightarrow$  Maximizing predictive accuracy

How to calculate  $P(w_{t+1} | w_t, \theta)$

Answers: Two vector representations

•  $v_w$  when  $w$  is center word

•  $u_w$  when  $w$  is a context word

center  $c$   
context  $w$

$$P(\theta | c) = \frac{\exp(u_w^T v_c)}{\sum_{w \in V} \exp(u_w^T v_c)} \rightarrow \begin{array}{l} \text{compare similarity of } \\ \theta \text{ and } c \\ \text{normalize over entire vocab} \end{array}$$

$$\text{softmax}(x_i) = \frac{\exp(x_i)}{\sum_{j=1}^n \exp(x_j)} = p_i$$

The softmax function maps arbitrary values  $x_i$  to a probability distribution  $p_i$

→ amplifier probability of largest  $x_i$  while still assigning some to smaller  $x$

### Training a model by optimizing parameters

$$\frac{\delta}{\delta V_c} \log \frac{\exp(u_0^T v_c)}{\sum_{w=1}^m \exp(u_w^T v_c)}$$

$$= \frac{\delta}{\delta V_c} \log \exp(u_0^T v_c) - \frac{\delta}{\delta V_c} \log \sum_{w=1}^m \exp(u_w^T v_c)$$

$$= \frac{\delta}{\delta V_c} u_0^T v_c +$$

$$= u_0^T v_c$$

negative word

gension

visualization. There's some

ambiguity with names.

$$\frac{d}{dv_c} \log \sum_{w=1}^V \exp(u_w^T v_c)$$

$$= \frac{1}{\sum_{w=1}^V \exp(u_w^T v_c)} \cdot \sum_{x=1}^V \frac{d}{dv_c} \exp(u_x^T v_c)$$

$$\sum_{x=1}^V \exp(u_x^T v_c) \frac{d}{dv_c} u_x^T v_c$$

$$\sum_{x=1}^V \exp(u_x^T v_c) u_x$$

$$\frac{\delta}{\delta v_c} \log P(O|C) = u_o - \frac{\sum_{x=1}^V \exp(u_x^T v_c) u_x}{\sum_{w=1}^V \exp(u_w^T v_c)}$$

$$= u_o - \sum_{x=1}^V \frac{\exp(u_x^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)} u_x$$

actual context word

$u \rightarrow$  outside  
 $v \rightarrow$  center

$v, v^T$   
softmax( $v, v^T$ )  
probs

$$= u_o - \sum_{x=1}^V P(x|C) u_x$$

1 observed representation of each word of context word

expected context word

probability of it in current model

gem sim (word similarity packet)

LDA

negative similarity

## Lecture 2 (Word Vectors & Word Senses) (4/10/2020) (2)

The word vectors are represented as rows

We want a model that gives a reasonably high probability estimate to all words that occur in the context (that's why we are using softmax)

Optimization: Gradient Descent

minimize  $J(\theta)$

$$\theta^{\text{new}} = \theta^{\text{old}} - \alpha \nabla J(\theta)$$

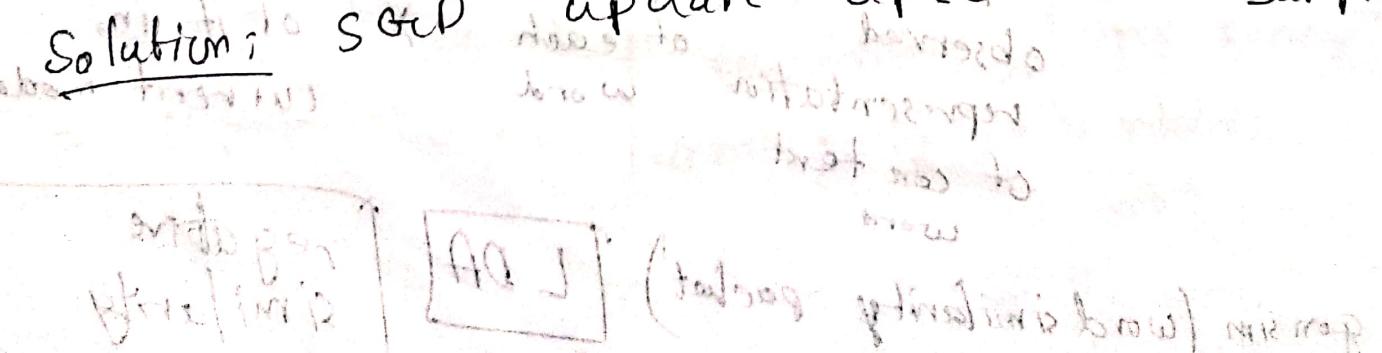
Stochastic Gradient Descent

Since  $J(\theta)$  is very expensive to compute here,

→ Very bad idea for pretty much all uses

SGD update after each sample

Solution:



## Stochastic gradient with wordvecs

We might only update the word vectors that actually appear.

Solution: Either we need sparse matrix update operations to only update certain rows of full embeddings matrices  $U$  and  $V$ , or we need to keep a hash of word vectors if you have millions of words and do distributed computing. It is important to not have to send gigantic updates around.

Why two vectors? Average both at the end

Easier optimization. Average both at the end

Two model variants:

1) Skip-Gram ( $SR$ )  $\Rightarrow$  predict context words given center word

2) Continuous Bag of Words ( $(CBOW)$ )  $\Rightarrow$  predict context words from (bag of) context words

Naive Softmax

Skip-gram with Neg-sampling: Normalization factor is too computationally expensive

Idea of Negative sampling:

Train binary logistic regression for a P

true pair (center word and word in its context)

versus several noise pairs (center word paired with a random word) → then try to give them as lower prob as possible

Paper: "Distributed Representations of Words and Phrases and their Compositionalities" (Mikolov 2013)

$$J(\theta) = \frac{1}{T} \sum_{t=1}^T J_t(\theta)$$

$$J_t(\theta) = \log \delta(u_b^T v_c) + \sum_{j=1}^k p(w_j) \left[ \log \delta(-u_j^T v_c) \right]$$

→ We maximize the probability of two words co-occurring in the first log

$$J_{\text{neg-sample}}(\theta, v_c, u) = -\log (\delta(u_b^T v_c)) - \sum_{k=1}^K \log (\delta(-u_k^T v_c))$$

We take k negative samples (using word probs)

- Maximize probability that word appears outside random words appear around center word

$$P(w) = \frac{W(w)}{\text{Unigram dist}}$$

count of word / total words

makes less frequent word to be sampled more often

But why not capture co-occurrence counts directly?

With co-occurrence matrix  $X$

- 2 options: windows vs. full doc
- Windows: similar to word2vec, use window around each word → captures both syntactic (POS) and semantic information
- Word-document co-occurrence matrix will give general topics (all sports terms will have similar entries) leading to "Latent Semantic Analysis"

- Example: window base co-occurrence matrix
- Window length 1 (more common - 10)
  - Symmetric (irrelevant whether left or right context)

### Problems with simple co-occurrence vecs

- Increase the size with vocab
- very high dimensional: requires a lot of storage
- subsequent classification model have sparsity issues
- models are less robust at scaling up.

Solution: Low dimensional vectors

store "most" of the important info in a fixed small number of dimensions; a dense vector usually 25-100 dimensions, similar to word2vec

Reducing dimensionality

## Method 1: Dimensionality Reduction of $X$ (HW1)

Singular Value Decomposition of co-occurrence matrix  $X$ .

Factorizes  $X$  to  $U\Sigma V^T$  where  $U$  and  $V$  are orthogonal

If we take  $k$  values then, we ignore the least frequent context vectors

$\hat{X}$  is the best rank approx to  $X$ , in terms of least squares

Problems: Expensive to compute for large matrices.

[Play around with it in homework]

Hacks to  $X$  (Rohde 2005)

Scaling the counts in the cells help a lot.

- Problems function word (the, he, has) are too frequent  $\rightarrow$  syntax has too much impact.  
→ Some fixes  $\min(X, t) + \approx 100$   
→ ignore all

- Ramped windows that count closer words more
- Use pearson correlation instead of counts, then set negatives to 0

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2} \sqrt{\sum (y_i - \bar{y})^2}}$$

### Count-based vs direct prediction

- LSA, HAL (Lund & Burgess)
- CoAE, Hellinger-PCA (Rohde) (Lebret & Collobert)

- Skip-gram/CBOW (Mikolov)
- NNLM (~~RNN~~) HLBL (Collobert)  
RNN (Bengio, Huang, Mnih)

Pros	<ul style="list-style-type: none"> <li>Fast training</li> <li>Efficient usages of statistics</li> </ul>
Cons	<ul style="list-style-type: none"> <li>Primarily used to capture word similarity</li> <li>Disproportionate importance given to large counts</li> </ul>

Cons	<ul style="list-style-type: none"> <li>Scales with corpus size</li> <li>Inefficient usage of stats</li> </ul>
Pros	<ul style="list-style-type: none"> <li>Generate improved performance</li> <li>Can capture complex patterns beyond word similarity</li> </ul>

Encoding meaning in vector differences [EMNLP 2014]

(Pennington)

[Combining the best of both worlds]

Crucial insight: ratios of co-occurrence probability  
can encode meaning components

Capturing ratios of co-occurrence probs as  
linear meaning

$$w_i w_j = \log P(i/j)$$

$$w_x (w_a - w_b) = \log \frac{P(x/a)}{P(x/b)}$$

$$J = \sum_{ij=1}^V f(x_{ij}) (w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log x_{ij})^2 \quad [\text{GloVe}]$$

- fast training

- Scalable huge corpora

- Good performance even with small corpus and  
small vectors

## Evaluate Word Vectors

### General Evaluation in NLP: Intrinsic vs Extrinsic

#### Intrinsic:

- Evaluation of a specific/intermediate subtask
- Fast to compute
- Helps to understand that system
- Not clear if really helpful unless correlation to real task is established.

#### Extrinsic:

- Evaluation on a real task
- Can take long time to compute accuracy
- Unclear if the subsystem is the problem or its interaction ~~is established or~~ <sup>other</sup> subsystems
- If replacing exactly one subsystem with another improves accuracy → winning

## Intrinsic word vector evaluation

- Word vector analogies

$a:b :: c:?$

man: woman :: king: ?

$$d = \operatorname{argmax}_i \frac{(x_b - x_a + x_c)^T x_i}{\|x_b - x_a + x_c\|}$$

- Evaluate word vectors by how well their cosine distance after addition captures intuitive semantic and syntactical analogy questions.

- Discarding the input words from the search!
- Problem: what if the info is there but not linear.

## for Glove

- Good dimension is  $\sim 300$
- Asymmetric context rate not as good
- This might be different for downstream tasks
- Window size of 8 around each center word is good for glove vectors

## Recent works on dimensionality of Word Embedding [NeurIPS 2018]

- Wikipedia is better than news text

- Word vector distances and their correlation with human judgment.

Example Wordsim353 (check out)

Cosine distance for similarity judging

Word senses and word sense ambiguity

→ Most words have lots of meanings

Does one vector capture all these meanings?

Improving Word Rep via Global Context (Huang et al 2012)

Ideal Cluster ~~around~~ windows around words, word retain with each word assigned to multiple different clusters bank<sub>1</sub>, bank<sub>2</sub>

Linear Algebraic Structure of Word Senses with Applications to Polysemy (Arora, TACL 2018)

Different senses of a word reside in a linear superposition (weighted sum). In standard word embedding a like wordvec

$$V_{\text{pike}} = d_1 V_{\text{pike}_1} + d_2 V_{\text{pike}_2} + d_3 V_{\text{pike}_3}$$

$$d_1 = \frac{f_1}{f_1 + f_2 + f_3} \quad f = \text{frequency}$$

- Because of ideas from sparse coding we can actually separate out the senses.

### Extrinsic Word Vector Evaluation:

All subsequent tasks in this class

- Good word vectors would help in named entity recognition:

### Assignment: Exploring Word Vectors

- Word vectors and word embeddings are often used interchangeably. The term "embedding" refers to the fact that we are encoding aspects of a word's meaning in a lower dimensional space

## Problems with thesaurus-based meaning:

- They have problems with recall
- Many words are missing
- Most phrases are missing
- Some connection between senses are missing
- Doesn't work well for verbs and adjectives

business and financial news corpus  
Reuters

## Singular Value Decomposition:

PCA

## Lecture 3 - Word Window Classification, Neural Nets

### Classification setup:

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

- $x_i$  are inputs, e.g. words, sentences, documents (Dimension)
- $y_i$  are labels (one of  $C$  classes) we try to predict, for example: → classes: sentiment, named entities, buy/sell
  - other words
  - multi-word sentences

### Traditional ML/Stats approach:

Assume  $x_i$  are fixed, train & softmax/logistic regression weights  $w \in \mathbb{R}^C$  to determine a decision boundary as in the picture

$$\exp(w_y \cdot x)$$

$$p(y/x) = \frac{\exp(w_y \cdot x)}{\sum_c \exp(w_c \cdot x)}$$

### Softmax in two steps

$$1) \text{ Find all } f_c \text{ for } c=1 \dots C$$

$$w_y \cdot x = \sum_{i=1}^d w_{y,i} x_i = f_y$$

$$2) P(y/x) = \frac{\exp(f_y)}{\sum_c \exp(f_c)} = \text{softmax}(f_y)$$

→ Our objective is to maximize the probability of the correct class  $y$ .

→ Or we can minimize the negative log probability of that class.

$$-\log p(y|f) = -\log \left( \frac{\exp(p_y)}{\sum_c \exp(f_c)} \right)$$

• NLL loss (Negative log-likelihood loss)

→ Cross entropy loss

True probability distribution  $p$

computed probability distribution  $q$

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

• Assuming a ground truth (one-hot or gold or target) probability distribution that is 1 at the right class and 0 everywhere else

• Because of one-hot  $p$ , the only term left is the negative log prob of the true class.

Classification over a full dataset:

• Cross entropy loss function over full dataset

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

$$f_y = f_y(x) = w_y x = \sum_j w_{yj} x_j$$

$$f = w x$$

Traditional ML Optimizations

$$\theta = \begin{bmatrix} w_1 \\ \vdots \\ w_d \end{bmatrix} = w(:) \in \mathbb{R}^d$$

Now we could update the decision boundary

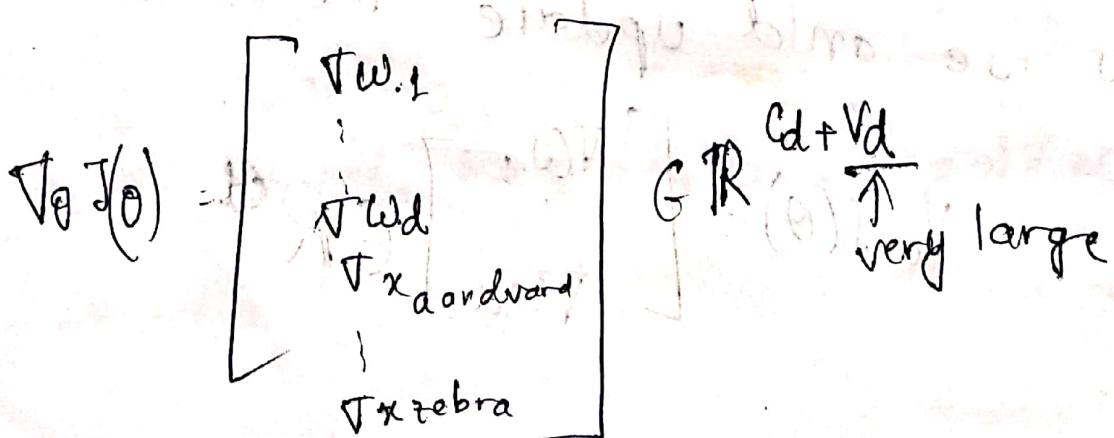
$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \nabla w_1 \\ \vdots \\ \nabla w_d \end{bmatrix} \in \mathbb{R}^d$$

## Neural Networks Classifiers (Non-linearity)

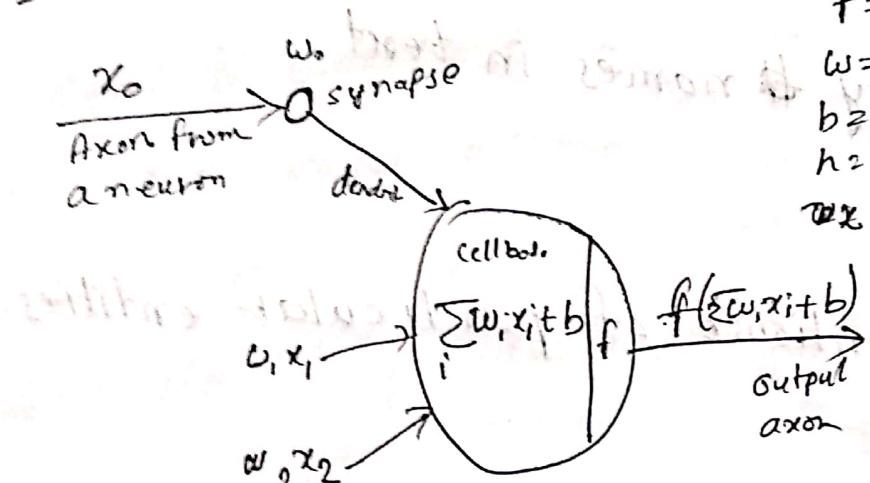
- can learn more complex functions and non-linear decision boundaries.

Classification difference with word vectors.

- Commonly in NLP deep learning:
  - We learn both  $w$  and word vectors  $x$
  - We learn both conventional params and representation
  - The word vectors re-represent one-hot vectors - move them around in an intermediate layer vector space - for easy classification with a (linear) softmax classifier via layer  $x \in \mathbb{R}^d$



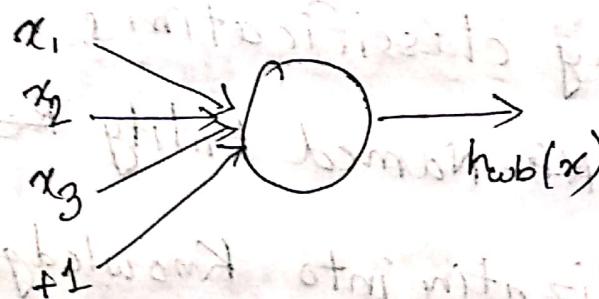
## An artificial neuron



$f$  = nonlinear activation  
 $w$  = weights  
 $b$  = bias  
 $h$  = hidden  
 $x$  = input

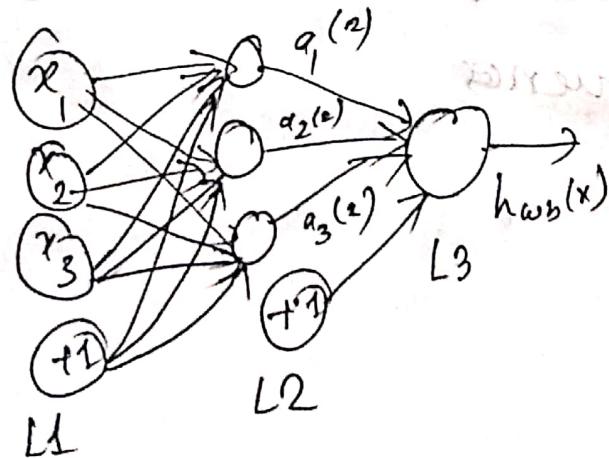
$$f(z) = \frac{1}{1+e^{-z}}$$

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



A neuron can be a binary logistic regression unit

A neural net = running several logistic regressions at the same time



It is the loss function that will decide what the intermediate hidden variables should be

# Named Entity Recognition (NER)

- Find and classify names in text

## Possible purposes:

- Tracking mentions of particular entities in document
- For question answering. Answers are usually named entities
- The same technique can be extended to other slot-filling classifications.

→ often followed by ~~NE~~ Named Entity Recognition

Linking / Canonicalization into Knowledge Base

→ We predict Entities by classifying words in context and then extracting entities as word subsequences

Why is it hard?

- Hard to work out boundaries of entity
- Hard to know if something is an entity
- Hard to know the class of unknown/novel entity
- Entity class is ambiguous and depends on context.

## Binary word Window classification

- In general, classifying single words is rarely done
- Interesting problems like ambiguity arise in context.

Example: Auto antonyms

- "To sanction" cannot mean 'to permit' or 'to punish'
- "To seed" can mean 'to place seeds' or 'to remove seeds'

Example: Resolving of ambiguous named entities

- Paris → Paris, France vs. Paris Hinton vs. Paris, Texas
- Hathaway → Berkshire Hathaway vs Anne Hathaway

## Window classification

- Idea: classify a word in its context of neighboring words.
  - For example, Named Entity Classification of a word in context
- A simple way to classify a word in context might be to average the word vectors in a window and to classify the average vector.
- Problem: that would lose positional info.

## Softmax idea:

Train softmax classifier to classify a center word by taking concatenation of word vectors surrounding it in a window.

$$\text{Resulting vector } \mathbf{x}_{\text{window}} = \mathbf{x} \in \mathbb{R}^{\text{wd}}$$

$\text{wd} = \text{windowsize} + 1$

## Simplest window classifier

with  $x = X_{\text{window}}$  we can use the same software

classifier as before

predicted  
model  
output  
probability

$$\hat{y}_g$$

$$P(y/x) =$$

$$\frac{\exp(w_y x)}{\sum_{c=1}^C \exp(w_c x)}$$

$$\exp(w_y x)$$

$$\sum_{c=1}^C \exp(w_c x)$$

$$f_{y_i} = w_y x$$

with cross entropy error as before

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

To update the word vectors we need to calculate grads like we did previously

## Binary classification with unnormalized scores

Collabor & Weston (2008, 2011) ICM Test of time 2018

- Assume we want to classify whether the center word is a location
- 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 centers are "true" positions and get a high score.

## NN Feed forward Computation:

$$\text{score}(x) = \mathbf{v}^T \mathbf{d}(\mathbf{f}(\mathbf{R}))$$

Compute a window's score with a 3-layer neural net:

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

$$s = \mathbf{v}^T f(\mathbf{w}x + b)$$

$$\mathbf{x} \in \mathbb{R}^{20 \times 1}, \mathbf{w} \in \mathbb{R}^{8 \times 20}, \mathbf{v} \in \mathbb{R}^{8 \times 1}$$

Main intuition before extra layer:

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

Alternative: Max-margin loss (no Softmax)

- Idea for training objective: Make true window's score larger and corrupt window's score lower (until they are good enough)

$$\text{Minimize } J = \max(0, 1 - s + s_c)$$

- This is not differentiable but continuous  
→ we can use SGD.

$$\theta^{\text{new}} = \theta^{\text{old}} - \alpha \nabla J(\theta)$$

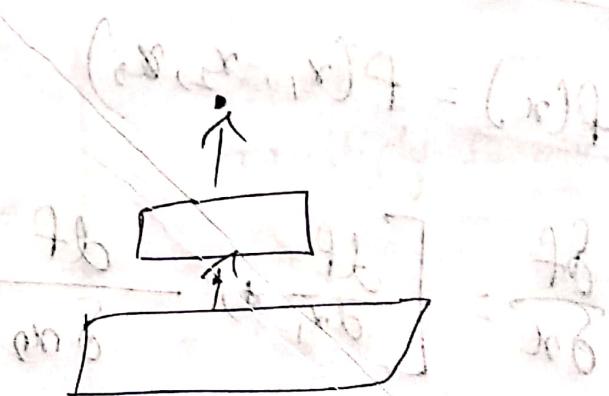
Simple net for score

$$s = u^T w$$

$$f = f(wx + b)$$

$x$  = input

let's find out  $\frac{\partial s}{\partial b}$



- In practice we care about gradient of the loss.

# Computing Gradients by Hand

- Review of multivariable derivatives
- Matrix calculus; fully vectorized gradients
  - Much faster and more useful than non-vectorized gradients
- But doing a non-vectorized gradient can be good practice

Gradients:

$$f(x) = f(x_1, x_2, \dots)$$

$$\frac{\delta f}{\delta x} = \left[ \frac{df}{dx_1}, \frac{df}{dx_2}, \dots, \frac{df}{dx_n} \right]$$

# Jacobian Matrix: Generalization of the Gradient

Given a function with  $m$  outputs and  $n$  inputs

$$f(x) = \begin{bmatrix} f_1(x_1, x_2, \dots, x_n), \dots, f_m(x_1, x_2, \dots, x_n) \end{bmatrix}$$

It's Jacobian is an  $m \times n$  matrix of partial derivatives

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \dots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

## Chain Rule

single variable

multiple derivatives

$$z = 3y$$

$$y = x^2$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx} = (3)(2x) = 6x$$

multivariable  
multiple Jacobian

$$h = f(x)$$

$$z = Wx + b$$

$$\frac{dh}{dm} = \frac{dh}{dx} \frac{dx}{dm}$$

$$= (AP^{-1}) \frac{dx}{dm}$$

~~gradient with respect to weights~~  
Jacobian: Elementwise activation function

~~what is~~

$h = f(z)$  what is  $\frac{dh}{dz} \quad h, z \in \mathbb{R}^n$

then  $h_i = f(z_i)$

$$\left( \frac{dh}{dz} \right)_{ij} = \frac{dh_i}{dz_j} = \frac{\delta_0}{\delta z_j} \circ f(z_i)$$

$$= \begin{cases} f'(z_i) & \text{if } i=j \\ 0 & \text{otherwise} \end{cases}$$

regular 1-variable derivative

$$\frac{\delta h}{\delta z} = \begin{pmatrix} f'(z_1) & & & \\ & \ddots & & \\ & & 0 & \\ & & & f'(z_n) \end{pmatrix} = \text{diag}(f'(z))$$

$$\frac{\delta}{\delta x} (w^T x + b) = w$$

$$\frac{\delta}{\delta b} (w^T x + b) = I$$

$$\frac{\delta}{\delta u} f(u^T h) = \star h^T$$

(identity matrix)

$$s = u^T h$$

$$h = f(Wx + b) \Rightarrow h = f(z)$$

$x$  (input)



(2)  $x$  (input)

$$\frac{\delta s}{\delta b} = \frac{\delta s}{\delta h} \frac{\delta h}{\delta z} \frac{\delta z}{\delta b}$$

$$= h^T \cdot \text{diag}(f'(z)) I$$

$$= h^T \circ f'(z)$$

$$\frac{\delta s}{\delta w} = \frac{\delta s}{\delta h} \frac{\delta h}{\delta z} \frac{\delta z}{\delta w} = \cancel{\delta h} \frac{\delta z}{\delta w} = \cancel{\delta h} \delta \frac{\partial z}{\partial w}$$

$$= h^T \cdot \text{diag}(f'(z))$$

$$\delta = \frac{\delta s}{\delta h} \frac{\delta h}{\delta z} = h^T \circ f'(z)$$



local error signal

$W \in \mathbb{R}^{n \times m}$

- L output nm inputs: 1 by nn Jacobian
- Inconveniant to do  $\theta^{\text{new}} = \theta^{\text{old}} + d\theta J(\theta)$

Instead, follow convention: shape of the gradient is shape of params

$\frac{\partial s}{\partial w}$  is nby m:

$$\left[ \begin{array}{c} \frac{\partial s}{\partial w_{11}} \\ \vdots \\ \frac{\partial s}{\partial w_{n1}} \end{array} \right] \quad \left[ \begin{array}{c} \frac{\partial s}{\partial w_{1m}} \\ \vdots \\ \frac{\partial s}{\partial w_{nm}} \end{array} \right]$$

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

↑      ↑  
local   local  
error   input  
signal   signal

[to get to appropriate shape]

$\frac{\partial s}{\partial b} = h^T f'(x)$  is a row vector

but according to shape convention in this  
should be column vector since be is a  
column vector.

Disagreement between Jacobian form (which makes chain rule easy) and the shape convention (which makes implementing SGD easy)

- We expect answers to follow the shape convention
- But Jacobian form is useful for computing answers

So,

Two options:

1) Use Jacobian form as much as possible, reshape to follow the convention at the end.

2) Always follow the convention.

[Look at dimensions to figure out when to transpose and/or reorder them]