

TA: Jiacheng Xu.

Syntactic semantic
propagation

eg.

The [city counsel] refused [the demonstrators] a permit
 because they violence.
 ↓ threatened (D)
 refer c/D prompted (D)
 used (D)
 dislike (C)
 caused (D)

* Difficulty in referring the pro-noun they

* Ambiguities

eg. Teacher Strikes Idle kids.
 v. adj.
 N. V.

eg. Iraq Head Seeks Arms
 of state weapons
 of body body part

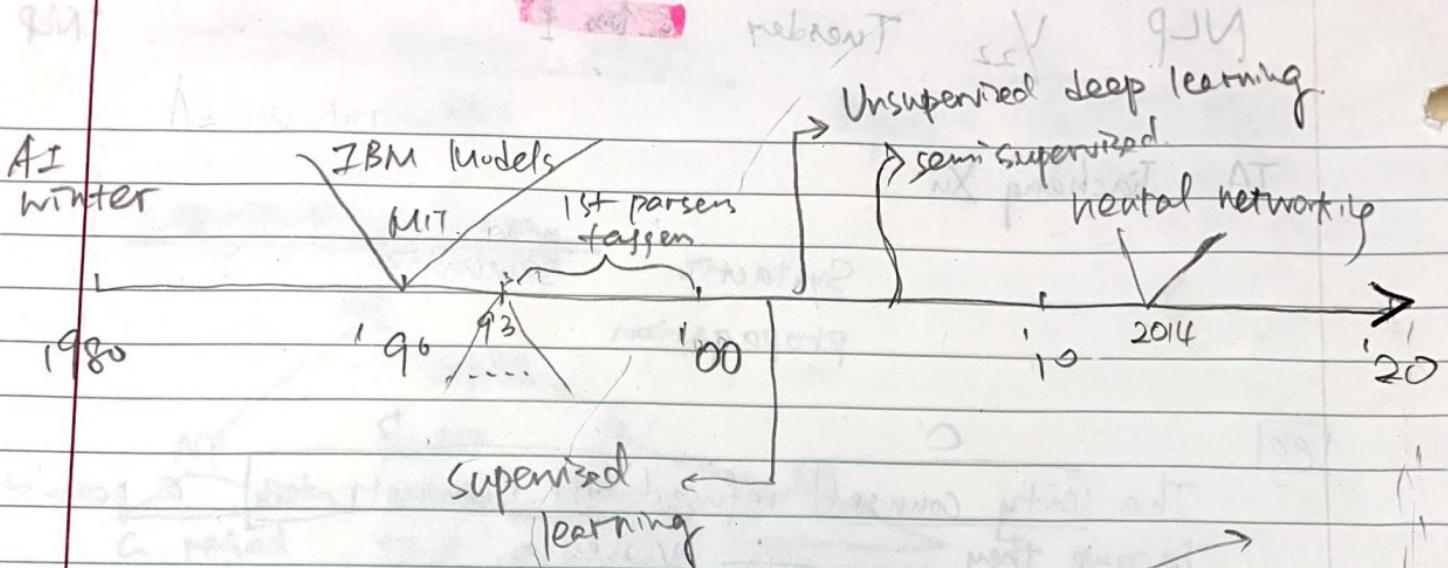
eg. Kids Make Nutritious Sharks
 create babies - sign taking hairy
 are/be

eg Translation

I fast vermont because → Today the weather is nice.

fast ↗ It's nice out.

He makes today beautiful.
 If fact actually handsome



NLP vs. Computational Linguistics.

words meaning

performant code

► Homework

► Midterm In-class

► Final Project : open-ended.

2.1, 4.1, 4.3

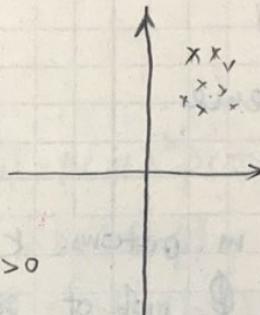
Outline

- Naive Bayes
- Feature Extraction
- Sentiment Analysis

ClassificationPoint $x \in \mathbb{R}^n$ Label $y \in \{0, 1\}$ Weights $w \in \mathbb{R}^n$

$$w^T x + b > 0$$

Transformation $x \in \mathbb{R}^n \rightarrow \mathbb{R}$
 $w^T x > 0$



Ex1 The movie was great! I would watch again

$$\bar{x} \Rightarrow f(x) \in \mathbb{R}^n$$

$$w^T f(x) > 0$$

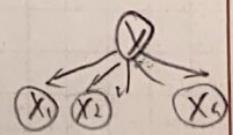
Ex2 the film was awful. I'll never watch it again

Bag of words: Say we have n wordsSentence $\bar{x} \Rightarrow [1.0 \dots 1.0]^n$
 the a ... great awful ...Value of i^{th} position = count(word i , \bar{x})Decisions: Use count or just 0-1?

tf - idf tf (term frequency) = Count(word) / Total words

idf = log $\frac{\# \text{ documents}}{\# \text{ docs containing word}}$ Index: Map from words to integersUnigram BOW [1 1 0] \rightarrow [the movie] (would watch) \rightarrow log(10)Naive Bayes

$$P(\bar{x}, y) = P(y) P(\bar{x}|y) \stackrel{\text{def}}{=} P(y) \prod_{i=1}^n P(x_i|y)$$



P

Maximum Likelihood

$$\prod P(x, y) = \prod P(y) P(\bar{x}|y) \cancel{\prod}$$

Weighted: HHTHT

$$P(H) = 0.75$$

$$P(H) = p \quad \prod P(y) = P(H)^3 P(T)$$

$$= p^3 (1-p)$$

$$3 \log p + \log(1-p)$$

$$\frac{3}{p} + \frac{1}{1-p} = 0$$

+ : It was good.

+ : It was not bad.

- : It was terrible

- : It was bad.

$$P(+)=P(-)=0.5$$

$$P(+|+) = \phi_{+,+} = \frac{2}{7}$$

$$P(\text{was}) = \frac{2}{7}$$

$$= \frac{1}{7}$$

$$= \frac{1}{7}$$

$$= \frac{1}{7}$$

It was good:

$$+ : \frac{1}{7} \cdot \frac{2}{7} \cdot \frac{2}{7} \cdot \frac{1}{7} = \dots$$

$$- : \frac{1}{7} \cdot \frac{2}{7} \cdot \frac{2}{7} \cdot 0 = 0$$

$$P(\text{negative given positive}) = 0$$

Smoothing

$$P(x|y) = \frac{\text{count}(x|y)}{\text{count}(\#, y)}$$

$$P(\bar{x}_i | y) =$$

Recap!

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

(\bar{x}, y) $\bar{x} \in \mathbb{R}^n$ $y \in \{0, 1\}$ $\bar{x}^T \bar{w}$ misaligned softmax

Bags of words \bar{x}_i = count of word i in the sentence

Naive Bayes : $P(\bar{x}, y) = P(y) P(\bar{x}|y)$ new planned

$$\propto \frac{P(y)}{\text{categorical}} \prod_{i=1}^n \varphi_{i,y}^{x_i} \quad \varphi_{i,y} = P(w_i|y) \text{ multinomial}$$

Estimation Data (\bar{x}_j, y_j) $j=1, \dots, D$

Maximum Likelihood : $P(y) = \frac{\text{Count}(y)}{\text{total}}$

$$\varphi_{i,y} = \frac{\text{Count of word } i \text{ in sents w/ label } y}{\text{Count of all words in sents w/ label } y}$$

Inference

$$P(y|\bar{x}) \propto P(y) \cdot P(\bar{x}|y)$$

this is computed these come from model params

Generative Vs Discriminative

$$P(\bar{x}, y)$$

$$P(y|\bar{x})$$

Directly optimizes what we care about (prediction)

- Can sample $(\bar{x}, y) \sim P$

(Logistic Regression)

Perceptron: error-driven discriminative method

$\times \times \times \times \times \times$

Sentiment Analysis

eg. On PowerPoint.

Perception

Decision boundary $\bar{w}^T \bar{x} \geq 0$

if ≥ 0 +
else -

Learning Given $(\bar{x}^i, y^i)_{i=1}^D \Rightarrow$ produce \bar{w}
for i in range (0, epochs)

for j in range (0, D)

$y_{pred} = \text{sign}(\bar{w}^T \bar{x}^j)$ +1 if ≥ 0
-1 else

w_0

$$\bar{w}_i = \bar{w}_0 + \bar{x}^i$$

Update \bar{w} : if $y_{pred} = y^j$, continue

$$w_0 \cdot \bar{x}^i = c$$

$$\bar{w}_i \cdot \bar{x}^i = c + \bar{x}^i \cdot \bar{x}^i > c$$

if $y_{pred} = -1$, $y^j = \pm 1$, $\bar{w} \leftarrow \bar{w} + \bar{x}^j$

$y_{pred} = +1$, $y^j = -1$, $\bar{w} \leftarrow \bar{w} - \bar{x}^j$

Ex $\bar{x}^1 = \text{good}$ $y^1 = +$

$$\begin{bmatrix} \text{good} \\ 1 \\ 0 \end{bmatrix} \circ \bar{w}_0 \leftarrow [0, 0]$$

$\bar{x}^2 = \text{bad}$ $y^2 = -$

$$\begin{bmatrix} \text{bad} \\ 0 \\ 1 \end{bmatrix} \circ \bar{w}_0 \leftarrow [0, 0]$$

not good

$$\begin{bmatrix} \text{good} \\ 1 \\ 0 \\ 1 \end{bmatrix}$$

not bad

$$\begin{bmatrix} \text{bad} \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

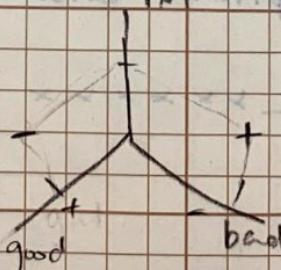
$$\bar{w}_0 \leftarrow \bar{w}_0 + [1, 0]$$

$$\bar{w}_1 = [1, 0]$$

$$\bar{w} = [1, -1]$$

$$w = [0, -1, -1]$$

$$w = [0, 0, 0]$$



separable

Lecture 3 (2)

Ex It was good

the movie was not bad.

not so good

$$w_{xit} = 1000$$

$$w_{movie} = 1000$$

$$w_{so} = -1000$$

negative word "stop"

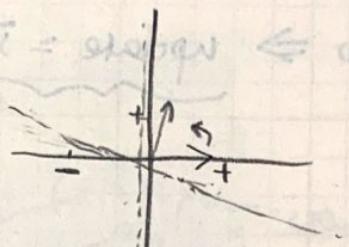
"stop" "good" "stop" "stop"

$$((\bar{w}^T \bar{x}) + 9 - 1) \bar{w}^T \bar{x} =$$

$$8 \text{ words} \rightarrow \text{it's good}$$

$$\text{"stop" or } \Leftrightarrow 1 \approx (+) 9$$

$$\bar{w}^T \bar{x} = \text{stop} \Leftrightarrow w_{so} (-) 9$$

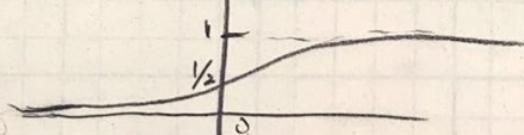


Logistic Regression

$$\text{Discriminative : } P(y|\bar{x}) \triangleq \frac{e^{\bar{w}^T \bar{x}}}{1 + e^{\bar{w}^T \bar{x}}} = \text{logistic}(\bar{w}^T \bar{x})$$

$$\text{logistic}(x) = \frac{e^x}{1 + e^x} = \frac{1}{1 + e^{-x}}$$

$$P(y|\bar{x}) > \frac{1}{2} \Leftrightarrow \bar{w}^T \bar{x} > 0$$



Training

Maximizing discriminative likelihood

$$L = \sum_{j:y_j=+} \log P(y=+|\bar{x}^j) + \sum_{j:y_j=-} \log P(y=-|\bar{x}^j) \approx \log \prod_j P(y \text{ is correct} | \bar{x}^j)$$

Train with stochastic gradient ascent

\bar{w}

Pick up (\bar{x}^j, y) from data

$$\bar{w} \leftarrow \bar{w} + \alpha \cdot \frac{\partial}{\partial \bar{w}} L((\bar{x}^j, y))$$

Stop

$$\alpha = 1$$

example

$$L((\bar{x}^j, y^j=+)) = \log P(y=+|\bar{x}^j) \\ = \log \frac{e^{\bar{w}^T \bar{x}^j}}{1 + e^{\bar{w}^T \bar{x}^j}}$$

$$P(y=+|\bar{x}) \\ = \bar{x}_i (1 - P(y=+|\bar{x}))$$

$$L = \bar{w}^T \bar{x} - \log(1 + e^{\bar{w}^T \bar{x}}) \\ \frac{\partial L}{\partial w_i} = \bar{x}_i - \frac{1}{1 + e^{\bar{w}^T \bar{x}}} \cdot e^{\bar{w}^T \bar{x}} \cdot \bar{x}_i \\ = \bar{x}_i \left(1 - \frac{e^{\bar{w}^T \bar{x}}}{1 + e^{\bar{w}^T \bar{x}}}\right)$$

$$\frac{\partial}{\partial w} L = \bar{x}^j (1 - P(+ | \bar{x}^j))$$

Recall $y^j = +$

$P(+) \approx 1 \Rightarrow$ "no update"

$P(-) \approx 0 \Rightarrow$ update = \bar{x}^j

↳ Looks like "soft" perception

$$\text{Loss} = -\log P(+ | \bar{x})$$

Maximize

[+ example]

$$\log \frac{1}{2} = -0.691$$

$$1 - \frac{1}{2} = \frac{1}{2}$$

$$\bar{w}^T \bar{x}$$

$$0 < \bar{x}^T \bar{w} < \bar{x}^T \bar{y}$$

$$\text{softmax } \frac{e^{x_i}}{\sum e^{x_j}}$$

$$(x_i, y_i) \in \{(x_j, y_j)\}_{j=1}^n$$

Recap Perceptron & Logistic Regression

weight vector \bar{w} , decision rule: $\bar{w}^T \bar{x} > 0$ LR: $P(y=+|\bar{x}) = \frac{e^{\bar{w}^T \bar{x}}}{1+e^{\bar{w}^T \bar{x}}}$

Learning

Perce for i in range (0, epochs)

for j in range (0, D)

$$y_{pred} = \text{sign}(\bar{w}^T \bar{x})$$

if $y_{pred} \neq y^j$

$$\text{if } y^j = + : \bar{w} \leftarrow \bar{w} + \alpha \bar{x}^j$$

$$y^j = - : \bar{w} \leftarrow \bar{w} - \alpha \bar{x}^j$$

LP

$$P(y=+|\bar{x}) = \frac{e^{\bar{w}^T \bar{x}}}{1+e^{\bar{w}^T \bar{x}}}$$

$$\text{if } y^j = + : \bar{w} = \bar{w} + \alpha \bar{x} (1 - P(y=+|\bar{x}))$$

$$y^j = - : \bar{w} = \bar{w} - \alpha \bar{x} (1 - P(y=-|\bar{x}))$$

Today: Optimization, Multiclass, Application examples

Optimization

LR had an objective function

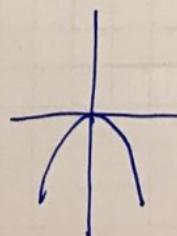
$$L(\bar{w}) = \sum_{j: y^j=+} \log P(y=+|\bar{x}^j) + \sum_{j: y^j=-} \log P(y=-|\bar{x}^j) \quad (\|w\|^2)$$



Maximize log likelihood \Leftrightarrow find \bar{w} to maximize $L(\bar{w})$

Alg! Stochastic gradient descent

$$\text{Repeat } \bar{w} \leftarrow \bar{w} + \alpha \frac{\partial}{\partial \bar{w}} L(\bar{w})$$



$$L(x) = -x^2 \quad \frac{\partial}{\partial x} -x^2 = -2x$$

$$x_0 = -2$$

$$\frac{\partial}{\partial x} -x^2 \Big|_{x=-2} = 4$$

$$x_1 = x_0 + \alpha \cdot g = -2 + 4 = 2$$

(1) Newton's method vs SGD approximation

Alg 2

Newton's method

$$\bar{w} \leftarrow \bar{w} + \left(\frac{\partial^2}{\partial w^2} L \right)^{-1} \frac{\partial}{\partial w} L$$

Hessian
matrix nxn
 $\frac{\partial^2}{\partial w^2}$
inverse

gradient
vector len n

good Newton

conjugate gradient

L-BFGS
Ade
Adam

SGD

fast

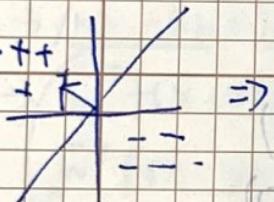
SGD momentum

Objective

$L = \text{maximize } \log \text{likelihood}$

Regularization:

Multiclass



$$\Rightarrow \begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 1 & 2 & 3 \\ \hline 2 & 3 & 1 \\ \hline 3 & 1 & 2 \\ \hline \end{array}$$

One vs-all

BC₁: (1 vs 2, 3)

BC₂: (2 vs 1, 3)

BC₃: (3 vs 1, 2)

$\bar{w}_{BC} \vec{x} > \bar{w}_{BC'}$?

$$\begin{array}{|c|c|c|} \hline 1 & 2 & 3 \\ \hline 2 & 1 & 3 \\ \hline 3 & 3 & 1 \\ \hline \end{array} \quad \text{all vs all}$$

BC₁₂: (1 vs 2)

key idea

Two ways to do this

① Different weights (DW):

w_y for each label y

$\underset{y}{\operatorname{argmax}} \bar{w}_y^T \vec{x} = \text{prediction}$

② Different features (DF): ~~fix y~~

Eisenstein

$f(\vec{x}, y)$ changes feature representation depending on y

$\underset{y}{\operatorname{argmax}} \bar{w}^T f(\vec{x}, y)$

"hypothesized" y

$$\begin{aligned} y &= \vec{w}^T \vec{x} + b \\ &= \vec{w}^T (\vec{x} + b) \\ &= \vec{w}^T \vec{x} + \vec{w}^T b \\ &= \vec{w}^T \vec{x} + w_0 \end{aligned}$$

NLP / Lecture 4 (2)

Ex Topic classification

\bar{x} = too many drug trials, too few patients

Bow: [dmg patients basebal]

DW $\bar{w}_{\text{health}} = [2.0, 5.6, -3] \Rightarrow 7.6$

$$\bar{w}_{\text{sports}} = [1.2, -3.1, 5.7] \Rightarrow -1.9$$

DF $f(\bar{x}, y) := \bar{x}$ replicated $|y|$ times (Cartesian product)

$$f(\bar{x}, y=H) = [1, 1, 0, 0, 0, 0, 0, 0]$$

$$f(\bar{x}, y=S) = [0, 0, 0, 1, 1, 0, 0, 0, 0]$$

$$\bar{w} = [2.0, 5.6, -3 / 1.2, -3.1, 5.7, \dots] \stackrel{\uparrow}{I} [\text{count of word, 1 label = sports}]$$

$\bar{w}^T f(\bar{x}, y=H) = 7.6$

$\bar{w}^T f(\bar{x}, y=S) = -1.9$

MC - Perception

for i range (0, epochs)

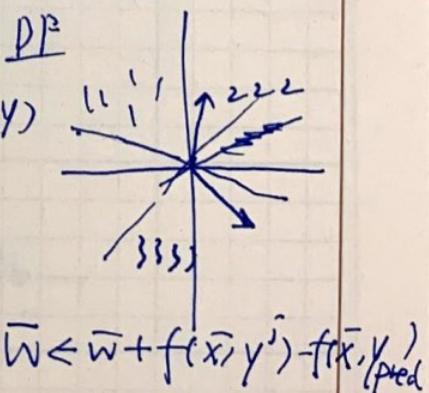
for j in range (0, D)

$$y_{\text{pred}} = \arg \max_y \bar{w}_y^T \bar{x} / \bar{w}^T f(\bar{x}, y)$$

if $y_{\text{pred}} \neq y_j$

$$\bar{w}_{\text{pred}} \leftarrow \bar{w}_{\text{pred}} - \bar{x}$$

$$\bar{w}_j \leftarrow \bar{w}_{\text{pred}} + \bar{x}$$



NLP 2/5 Tuesday Lecture 5

Recap (some subset) words still not written: UN

Optimization maximize (minimize) objective function $f(\bar{w})$

Stochastic gradient descent (SGD for minimizing)

Adam/Adagrad/Adadelta

$$\bar{w} \leftarrow \bar{w} + \alpha \frac{\partial f(w)}{\partial w}$$

$\bar{x}^T (\bar{w} - w)$ evaluated on one instance
decrease during training

Multiclass

$$\operatorname{argmax}_y \bar{w}^T \bar{x}$$

different weights

e.g. $\bar{x} = [1, 1, 1, 0]$ drug patient baseball

$$w = [x \ x \ x \ x] w_y$$

$$\operatorname{argmax}_y \bar{w}^T f(\bar{x}, y)$$

different features

$$f(\bar{x}, y=y_1) = [1, 1, 0, 1, 0, 0, 0, 0, 0]$$

$$f(\bar{x}, y=y_2) = [0, 0, 0, 1, 1, 0, 0, 0, 0]$$

$$w = [x \ x \ x \ x \dots x]$$

MC Perceptron Update: $\bar{w}_{y_{\text{pred}}} \leftarrow \bar{w}_{y_{\text{pred}}} - \alpha \bar{x}$

$$\bar{w}_{y_i} \leftarrow \bar{w}_{y_i} + \alpha \bar{x}$$

Neutral Nets

Linear classifier $\bar{w}^T \bar{x} \geq 0$

good

$$[g \ b \ n]$$

bad

$$[0 \ 1 \ 0]$$

not good

$$[1 \ 0 \ 1]$$

not bad

$$[0 \ 0 \ 1]$$

neutral

$$[0 \ 0 \ 0]$$

add

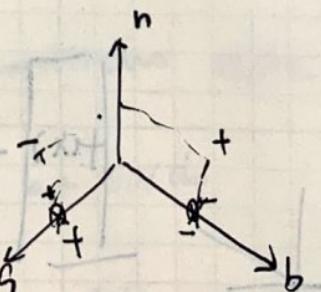
$$[g \ b \ n \ ng \ nb]$$

$$[1 \ 1 \ 0 \ 0 \ 0]$$

$$[1 \ 1 \ 1 \ 0 \ 0]$$

$$[1 \ 1 \ 1 \ 1 \ 0]$$

$$[1 \ 1 \ 1 \ 1 \ 1]$$



$$z = 0 \quad [0 \ 0 \ 0]$$

$$z = 2.0 \quad [2.0 \ 0 \ 0]$$

$$z = 4.0 \quad [0 \ 2.0 \ 0]$$

- Bigrams require crafting. What if we need trigrams / 4-gram?
- Kernels are expensive
- We probably don't need all pairs

NN: transform data into latent feature space

$$\hat{z} = g(W\bar{x}) \quad \bar{x} \in \mathbb{R}^3$$

$$W^T z > 0$$

g = nonlinearity

\tanh



$$\Rightarrow W^T (W\bar{x})$$

$$+ \tanh(1) \approx 1$$

$$\Rightarrow (W^T W)^T \bar{x} \geq 0$$

$$+ \tanh(2) \approx 1$$

\bar{W}^T

$$W = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 1 \end{bmatrix}$$

$$x = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$$

$$g(W\bar{x}) = \begin{bmatrix} 1 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix}$$

$$\begin{aligned} y &= g(Wx + b) \\ z &= g(Vy + c) \end{aligned}$$

$$z = g(Vg(Wx + b) + c)$$

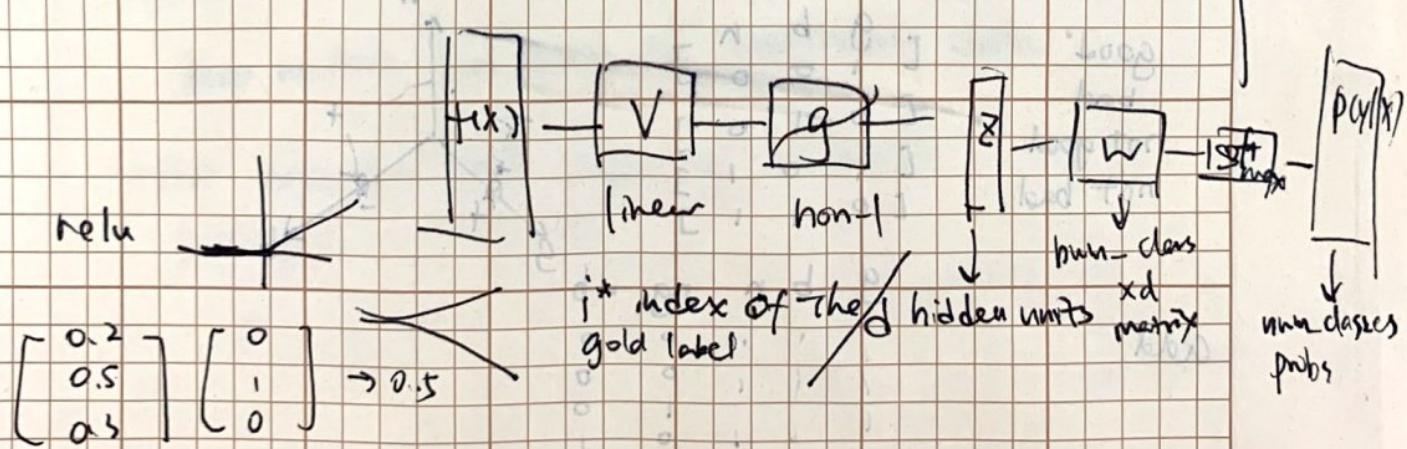
output of 1st layer

Deep Neural Network

Feed forward computation (not recurrent)

On PowerPoint
Vectorization and softmax.

exponentiate and normalize



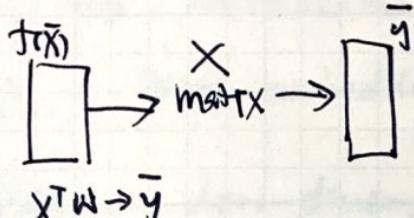
$$L(x, i^*) = Wz \cdot e_{i^*} - \log \sum_j (Wz) \cdot e_j$$

$$\frac{\partial}{\partial W_{ij}} L(x, i^*) = \begin{cases} z_j - p(y=i|X)z_j & \text{if } i = i^* \\ -p(y=i|X)z_j & \text{otherwise} \end{cases}$$

Recap $P(y|x) = \text{softmax}(Wg(Vfx))$

--- slides ---

Batching



$$[n] \quad [n \times m] \quad [m]$$

$$(\bar{x}^j, \bar{y}^j)_{j=1}^D \quad \text{if } \bar{x}^j = x^T W.$$

for ~~for $j = 1$ to D~~ $\bar{x}^j = f^j x = W$

$$x^T W \geq y$$

Restoring competence to handle multiple examples all at once.

2

$$\begin{bmatrix} \vdots \\ \vdots \end{bmatrix} \cdot \begin{bmatrix} \vdots & \vdots \\ 100 & 100 \\ \vdots & \vdots \end{bmatrix}$$

Slides

Word Embeddings : mapping word $n \in \mathbb{R}^d$: $n = 50 \sim 300$

. Deep models are good at continuous data

• What does the first layer of an NN do

$$\beta \cdot W = \begin{bmatrix} \text{the} & \text{good} & \text{matrix} \\ 1 & 0 & 0 \\ 1 & 0 & 0 \end{bmatrix} = X$$

$$Wx = \sum_{x_i} w$$

The president of
the
of
decreed.

NLP 2/12 Tuesday Lecture 7 (1)

Recap

- Pytorch

- Training NNs: Initialization, right optimizer

Today's Outline: Word embedding, evaluation, NNs for NLP

- Word embeddings low-dimensional (50-1000) vector representation of word types
(high-dim: \mathbb{R}^n each word has its own axis)

High-low dim: "You shall know a word by the company it keeps"
- JR Firth 1957

e.g. the president signed
 { the president announced
 { the governor signed
 { the governor announced

Mikolov predict a word from its context (and vice versa)

Each word w has a vector \bar{u}_w and a context vector \bar{c}_{cont}

Continuous bag-of-words (CBOW)

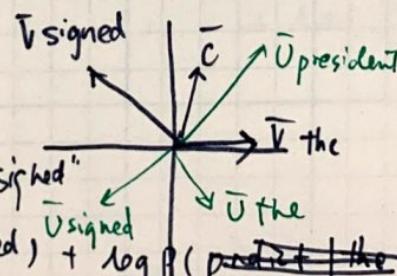
$$\boxed{\bar{u}_{\text{the}}} \oplus \boxed{\bar{u}_{\text{signed}}} = \bar{c} = \frac{1}{n} \sum_{j \in \text{context}} \bar{v}_j \quad p(w=w_i | \bar{c}) = \frac{\exp(\bar{u}_i \cdot \bar{c})}{\sum_{j \in v} \exp(\bar{u}_j \cdot \bar{c})}$$

Predict "president" from context "the - signed"

Maximize $\log p(\text{predict } (\text{the} - \text{signed})) + \log p(\text{president} | \text{the signed} | \text{president} - \text{e})$

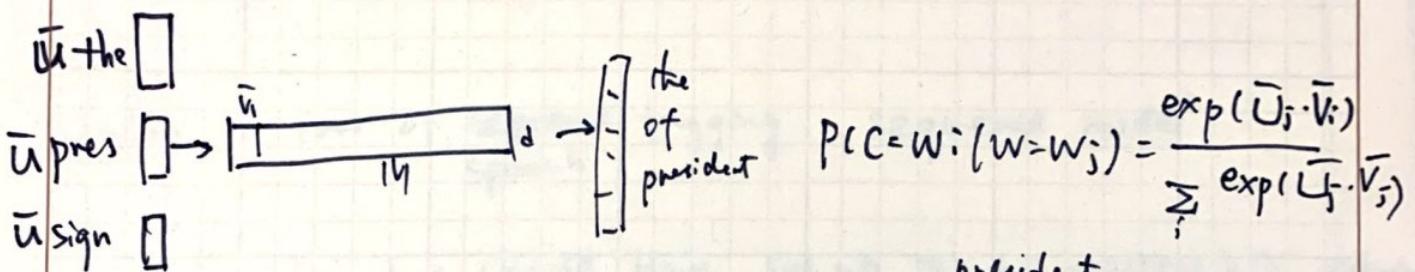
Optimize \bar{v}_s and \bar{u}_s both.

Tractable parameter: window size



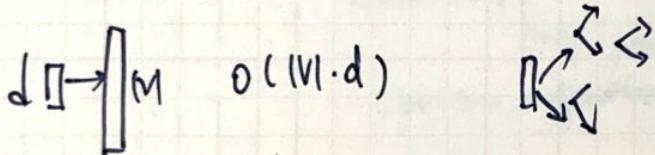
$p(w=\text{president} | \text{the} - \text{signed})$
 $p(\text{signed})$
 $\Leftrightarrow \underbrace{\text{the president signed ability}}_{W_{S1}} \underbrace{\text{in window}}_{W_{S2}}$

Skip-gram "Inverse of CBOW"



$$\text{Max : } \log P(\text{the} | \text{president}) + \log p(\text{signed} | \cancel{\text{president}}) \xrightarrow{ws=1} d=100 \\ + \log (\text{president} | \text{signed}) + \log (\text{a} | \text{signed}) \quad d=100 \\ M=10,000 \quad dW \sim 1000, 10^4$$

Runtime Solution 1 hierarchical softmax



Solu 2 skip-gram w/ negative sampling

Det a gives ws (word, context) pairs

Make synthetic (fake word, center)
 fake word freq

$$P(\text{real} | w=w_i, c=c_j) = \frac{\exp(\vec{U}_i \cdot \vec{V}_j)}{1 + (\vec{U}_i \cdot \vec{V}_j)}$$

~~$$\text{Max} \times \log P_{\cancel{\text{real}}}(\cancel{\text{real}} | \text{real words})$$~~
~~$$+ \log P_{\cancel{\text{real}}}(\text{real} | \cancel{\text{real}})$$~~

$$\text{Max} \times \log P(\text{real} | \text{real data}) \\ \rightarrow \log P(\text{real} - \text{fake data})$$

NLP 7/4 Thursday. Lecture 8

Recap

word embedding

Today Part of sequence tagging, sequence model
speech

Where's ambassador should have set up the big meeting in DC yesterday.

NNP NN MD VB (model) RB JJ NN. NNP NV
VBZ (verb) DT (determiners) IN

Closed class, pronouns, determiners, conjunctions,
prep. particles
Open class, adjectives, adverbs.

Words Ambiguity: Fed raises the interest by 5% rates.

NN VB

NN VB

NN VB

Input = $x_1 x_2$

Output = $y_1 y_2$

A multi class

Indicator features.

"Different Features"

$I_{[prev=rares \wedge tag=NN]}$

$I_{[prev=raises \wedge tag=NN]}$

Taggy

$$\text{Input} = \overline{x}(x_1, x_2)$$

$$\text{Output } \overline{y} = (y_1, y_2)$$

Reasons Match out reason about a Sentence's bags collectively.

① Generative Model = Hidden Markov Model

② discriminative model = conditional random fields (CRFs)

HMMs

$$P(\overline{y}, \overline{x}) = P(y_1) P(x_1 | y_1) \cdot P(y_2 | y_1) P(x_2 | y_2) \underbrace{P(y_3 | y_2)}_{P(y_3 | y_1)}$$
$$P(y_1) P(y_2 | y_1) P(y_3 | y_2).$$

Why Markov.

Parameter Parameter (initial and $P(y_i)$)

$$P(y_{i+1} | y_i) = \boxed{\quad}$$

$$P(x_i | y_i) = \boxed{\quad}$$

cols
starts 2.

new row

Today Estimation, Inference (NER, CRFs)

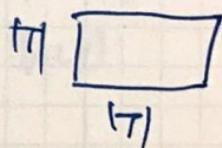
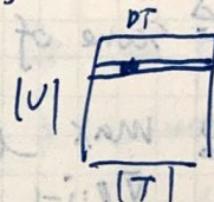
Parameter Estimation

Assume we have D examples (\bar{x}^i, \bar{y}^i) for $i=1 \dots D$

Find params to maximize $\sum \log P(\bar{x}^i, \bar{y}^i)$

$$= \sum_{i=1}^D [\log p(y_i) + \sum_{j=1}^n \log P(x_j^i | y_j^i) + \sum_{j=2}^n \log P(y_j^i | y_{j-1}^i)]$$

The $p(\text{word}|\text{tag}) = \frac{p(\text{tag}(\text{word}))}{p(\text{tag})}$



Estimation Count and normalize

Smoothing Laplace add α

Inference $\arg \max_{\bar{y}} P(\bar{y} | \bar{x}) = \arg \max_{\bar{y}} \left(\frac{P(\bar{x}, \bar{y})}{P(\bar{x})} \right) = \arg \max_{\bar{y}} P(\bar{x}, \bar{y})$

HMM $P(\bar{x}, \bar{y}) = \prod_{i=1}^n P(x_i | y_i) = \arg \max_{\bar{y}} \log P(\bar{x}, \bar{y})$

$N \xrightarrow{V_2(N)} V_1(N) \xrightarrow{V_1(N)} \dots$

$= \arg \max_{\bar{y}} \left[\sum_i \log P(x_i | y_i) + \sum_i \log P(y_i | y_{i-1}) \right]$

Viterbi Algorithm compute most likely tag seq from HMM

	they	can	fish	
N	-1	-3	-1	$P(w t)$
V	-3	-1	-1	

(rows sum to 1)

Transition

$$\begin{aligned} \text{START} &\rightarrow N_1 \rightarrow V_1 \rightarrow \text{STOP} \\ N_1 &\rightarrow N_2 \rightarrow V_2 \rightarrow \text{STOP} \\ V_1 &\rightarrow N_2 \rightarrow V_2 \rightarrow \text{STOP} \end{aligned}$$

$$\max_{\bar{y}_{1:n}} \left[\sum_{i=1}^n \log P(x_i | y_i) + \sum_{i=1}^n \log P(y_i | y_{i-1}) + \underbrace{\log P(\text{stop} | y_n)}_{\text{drop}} \right]$$

$$= \max_{y_n} \left[\log P(x_n | y_n) + \max_{y_{n-1}} \left[\log P(y_n | y_{n-1}) + \max_{y_{n-2}} \left[\dots \sum_{i=1}^{n-1} \dots \right] \right] \right]$$

Define matrix $V_i(y) = n \times |\mathcal{T}|$ matrix

$V_i(y)$ \triangleq score of best path ending in y at time i

$$\triangleq \max_{\bar{y}_{1:i-1}} \log P([\bar{y}_{1:i-1}, y], \bar{x}_{1:i})$$

Compute $V_i(y) = \log P(x_i | y) + \max_{y_{\text{prev}}} \log P(y | y_{\text{prev}}) + V_{i-1}(y_{\text{prev}})$

$$V_i(y) = \log P(x_i | y) + \log P(y | \text{START})$$

Viterbi for $i=1 \dots n$ n
 for $y \in \mathcal{T}$ $|\mathcal{T}|$
 $V_i(y) \leftarrow \text{compute } -V(i, y) \quad |\mathcal{T}|$

for $y_{\text{last}} \in \mathcal{T}$ $|\mathcal{T}|$
 final-score = $V_n(y_{\text{last}}) + \log P(\text{stop} | y_{\text{last}})$

Result = max final-score

Inference

$$\begin{array}{r} N \xrightarrow{-2} -2-3=\boxed{-1} \\ V \xrightarrow{-4} -1-3=-8 \\ \quad \quad \quad \xrightarrow{-4} \text{(targ)} \end{array}$$

Viterbi w/ log probs. max, + (scoring)

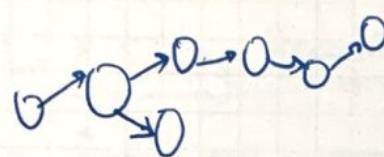
Viterbi w/ probs. max, *

$$\sum_i * (\log - \Sigma, \Sigma)$$

MP | Lecture 9 (2)

$\sum_{\bar{Y}} P(\bar{X}, \bar{Y})$

Algorithm: Sum-product : $\Sigma, * (\log - \Sigma, \Sigma) \quad \sum_{\bar{Y} \in \mathcal{Y}^{1:n}} \prod_i P(x_i | y_i) \dots = p(\bar{x})$



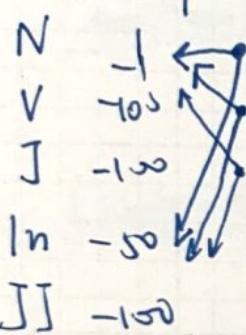
n words

$|T|$ tags

Time complexity $(n|T|^n)$

$V_i(y)$

they can fish



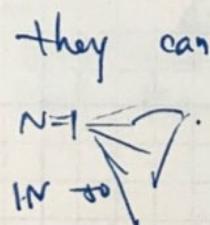
Beam Search: at i th timestep, only maintain

top k values of $V_i(y)$

$k=1$: greedy

$k=2$

Maintain priority queue over $V_i(y)$ for each i
beam



$$\underbrace{n \text{ words}}_{O(n|T|k \log(|T|k))} \underbrace{|T| \text{ tags}}_{L_s=k}$$

Today

Named Entity Recognition

- Frame as a sequence problem with a BIO tagset: Begin Inside Outside
- Why might HMM not do so well now?

what if unknown words :

Conditional Random Fields (CRF)

$$P(y|x) = \frac{\prod_k \exp(\phi(x_i, y_i))}{\sum \prod_k \exp(\phi(x_i, y_i))}$$

Constituency Parsing

