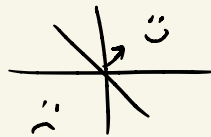


CS 371N Lecture 2



Classification 1: Features, Perceptron

Announcements

- AI released, due in 2 weeks
- Reading notation \neq lecture notation

- Today
- Classification (linear, binary)
 - Feature extraction
 - ML basics + perceptron

Classification

Points \bar{x}

for us:
strings

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

f : feature
extractor

Label $y \in \{-1, +1\}$

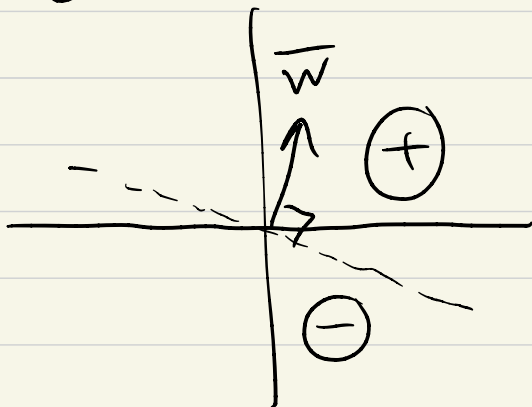
Classifier $\bar{x} \rightarrow y$

Linear classifier: weight vector $\bar{w} \in \mathbb{R}^n$

Decision rule: $\bar{w}^T f(\bar{x}) \stackrel{+b}{\geq} 0$

if $> 0 \Rightarrow +1$ else $\Rightarrow -1$

$n=2$



Sentiment Analysis

\bar{x} = the movie was great

wow, that was sooooo bad!

① Feature extraction:

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

string

② Learning algorithm D exs

Training set $\left\{ \left(f(\bar{x}^{(i)}), y^{(i)} \right) \right\}_{i=1}^D$

$= ??? \Rightarrow \bar{w}$ learned weight vector

Feature Extraction

Combos

- what words are there

~~X~~ - are the words "positive" or "negative"

~~X~~ - order \approx - punc. (! : ())

\approx "intensity" (good vs. great)

- context of words ("not great")

~~X~~

\bar{x} = the movie was great

Our basic tool: bag-of-words features

$[1 \ 0 \ 1 \ 0 \dots 0 \ 1 \dots 0 \dots 1]$
the a... movie... good great... was...

Vocabulary of n words $n \approx 10,000$

1 if present (or count)

10K-dim vector, 4 1s 9996 0s
¹
sparse

weight vector $\bar{w} \in \mathbb{R}^{10,000}$

$[-0.1 \ +0.2 \ \dots \ +0.3 \ \dots \ +10 \ \dots \ -0.1]$
the a movie great was

$$\bar{w}^T f(\bar{x}) = 10.1 = w_{\text{the}} \cdot 1 + w_{\text{movie}} \cdot 1 \\ + w_{\text{was}} \cdot 1 + w_{\text{great}} \cdot 1$$

Problems ① "not great" $w_{\text{not}} + w_{\text{great}}$

↳ also "really great"

↳ no order of words

② Static weights

word senses (awesome, good)

③ Different weights for related words

good, great
+10 "never seen"

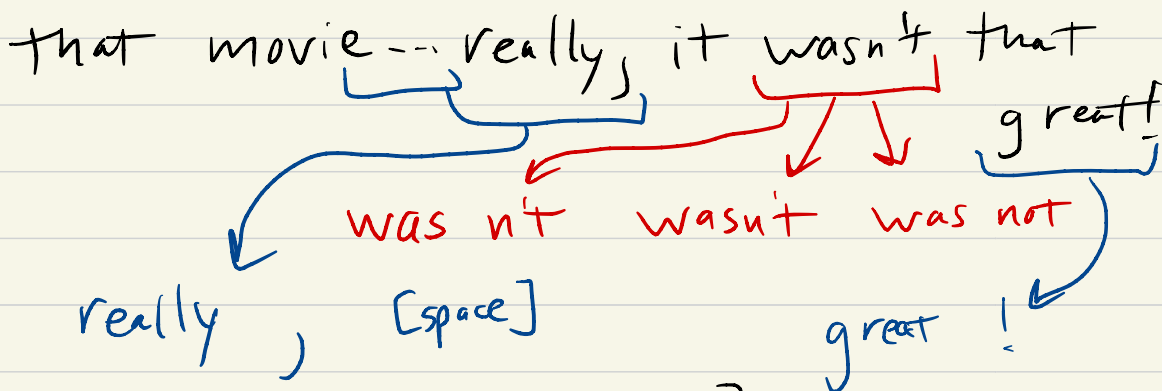
Preprocessing ① Vocab selection

vector space is fixed

maybe we look at our training data

$\bar{w} \in \mathbb{R}^n$, n doesn't change at
test time

replace rare words in train w/UNK
learn a weight for it



$\bar{w} = \{ \overset{+3}{\text{great}} \dots \overset{+7}{\text{great!}} \dots \} \dots$???

Tokenization - break out punc.
- break out contractions

② Remove stopwords (the, of, a, ...) optional

③ Lowercasing / Stemming → arrives
→ arrive

Not "not great" should be different from "great"

[the -- great -- ¹ not_great ...]

Bigram bag-of-words

Unigrams (each word)

Bigrams (each adjacent pair) } → V

{ the a not_good movie horror_movie ... }

U U B U B

We can select what bigrams go in the vocab! Did we solve "not"?

- ① No: not + X for all X
- ② "not very good"

Machine Learning

Optimize parameters \bar{w} to fit training data

$$\text{loss} = \sum_{i=1}^D \text{loss}(\bar{x}^{(i)}, y^{(i)}, \bar{w})$$

→ if we use \bar{w} to predict on $x^{(i)}$, how badly do we mess up w.r.t. $y^{(i)}$?

(Stochastic) Gradient Descent

for t in range(0, epochs):

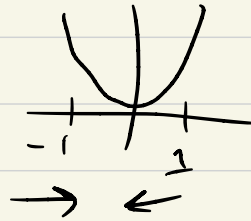
for i in range(0, D):

$$\bar{w} \leftarrow \bar{w} - \alpha \frac{\partial}{\partial \bar{w}} \text{loss}(\bar{x}^{(i)}, y^{(i)}, \bar{w})$$

α step size for now = 1

Subtract gradient of loss \Rightarrow
find \bar{w} with lower loss

$$\text{loss}(w) = w^2$$



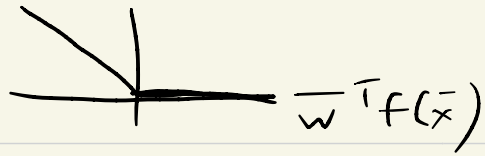
$$w=1$$

$$\frac{\partial}{\partial w} \text{loss} = 2$$

$$-\frac{\partial}{\partial w} \text{loss} = -2$$

Perceptron

$$y^{(i)} = +1$$
$$\text{loss:}$$



Init $\bar{w} = \bar{0}$

for t in range $(0, \text{epochs})$

for i in range $(0, D)$

$$y_{\text{pred}} \leftarrow \begin{cases} 1 & \text{if } \bar{w}^T f(\bar{x}^{(i)}) > 0 \\ -1 & \text{else} \end{cases}$$

$$\bar{w} \leftarrow \begin{cases} \bar{w} & \text{if } y_{\text{pred}} = y^{(i)} \\ \bar{w} + \alpha f(\bar{x}^{(i)}) & \text{if } y^{(i)} = +1 \\ \bar{w} - \alpha f(\bar{x}^{(i)}) & \text{if } y^{(i)} = -1 \end{cases}$$

Set $\alpha = 1$ for now

Suppose $\bar{w}^T f(\bar{x}^{(i)})$ was < 0

$$\text{Let } \bar{w}' = \bar{w} + \alpha f(\bar{x}^{(i)})$$

$$\bar{w}'^T f(\bar{x}) = \bar{w}^T f(\bar{x}^{(i)}) + \underbrace{\alpha f(\bar{x}^{(i)})^T f(\bar{x}^{(i)})}_{> 0}$$

Sparsity

If $f(\bar{x}^{(i)})$ only involves 4 features w/nonzero values, Computing y_{pred} and the new \bar{w} only involves those 4 features

Step size For w^2 case
 ~~Ψ~~ need $\alpha < 1$

In general: decrease α over training

One possibility: $\alpha = \frac{1}{t}$ + epochs

$\alpha = e^{-t}$ ← drops too fast

Do not randomly init. \bar{w} on A_1