

CS 371N Lecture 2

Classification 1: Features, Perceptron



Announcements

- AI
- Reading notation != 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 maps $\bar{x} \rightarrow y$

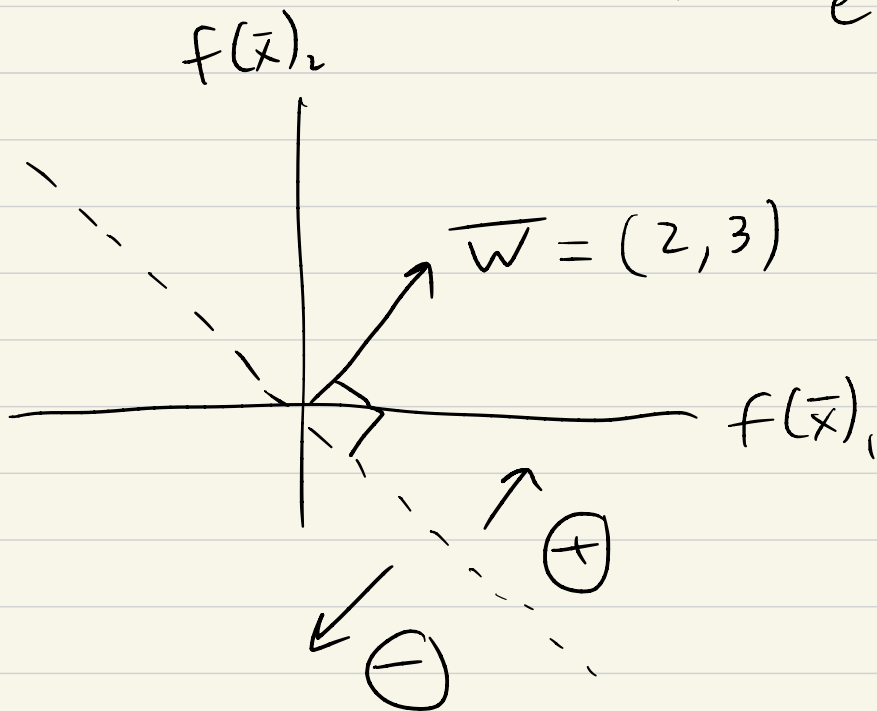
Linear classifier: represented by a weight vector $\bar{w} \in \mathbb{R}^n$

Decision rule: $\bar{w}^T f(\bar{x}) \stackrel{?}{>} 0$

$n=2$

$\bar{w} \cdot f(\bar{x})$

if $> 0 = +1$
else $= -1$



Sentiment Analysis

\bar{x} = the movie was great!

① Feature extraction

$$\begin{array}{l} \bar{x} \Rightarrow f(\bar{x}) \in \mathbb{R}^n \\ \text{string} \end{array}$$

② Learning algorithm

Training set $\left\{ (f(\bar{x}^{(i)}), y^{(i)}) \right\}_{i=1}^D$ $\begin{array}{l} D \text{ exs} \\ \downarrow \\ D \end{array}$

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

Today: cover #1

see Perceptron for #2

Feature Extraction

\bar{x} = the movie was great

- What do we want?
- reflect word order X
 - word frequency ✓
 - parts of speech X
 - word meaning / proximity X X
 - negation (not good) X

Bag-of-words featurization

$\left[\begin{array}{ccccccc} 1 & 0 & 0 & & 1 & & 1 & 1 \end{array} \right]$
the a of ... movie ... great ... was ...

vocabulary $\sim 10,000$ words

1 if present

OR count of the word

sparse
feat.
repr.

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

+7
awesome

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

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

"weighted voting" w_{great}

awesome and great have independent weights

Preprocessing ① Vocab selection:
vector space is a fixed set of words

replace unseen words w/ UNK
have a weight for UNK

- split(" ")

That movie... really, it wasn't great!
↓
was not
↓
great!

punc		+3	+7	
$\overline{w} =$	[great	... great!	18,000
split_punc		+3	+0	
$\overline{w} =$	[great	... ! ...	15,000

Typical tokenization

- break out punc.
 - break out contractions
 - ② Remove stop words (the, of, be, etc.)
 - ③ Lowercasing / stemming (arrived ⇒ arrive)
- Optional

Is lowercasing always good? No!

- text messages

- capture names

Revisit "not awesome"

{ ... not awesome ... movie was }

bigram bag of words

sentence \Rightarrow adjacent pairs

vocab is now huge (\sim old vocab²)

We can mix unigrams + bigrams in
our feature space

"Custom" feature space

movie	1
was	1
awesome	0
great	1
=	3
not awesome	0
movie was	1
⋮	

Machine Learning

Optimize parameters \bar{w} to fit some training data (labeled)

We want \bar{w} to make good predictions

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

"if we use \bar{w} as our weights, how badly do we mess up?"

Stochastic Gradient Descent

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

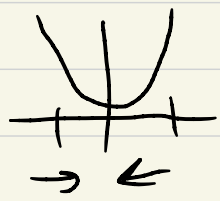
for i in range $(0, D)$

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

step size
 ≈ 1

Subtracting gradient of the loss \Rightarrow
finding a \bar{w} with lower loss

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



$w = 1$
gradient $= 2w = 2$

$w = -1$
gradient $= -2$

$$\bar{w}^T f(\bar{x}) y^{(i)}$$

Perceptron (instance of SGD)

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

for t in range (0, epochs)

for i in range (0, D) (shuffle exs each epoch)

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

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

Let $\alpha = 1$ for now

before

$$\bar{w}^T f(\bar{x}^{(i)})$$

after

$$(\bar{w} + f(\bar{x}^{(i)})^T) f(\bar{x}^{(i)})$$

||

$$\bar{w}^T f(\bar{x}^{(i)}) + \underbrace{f(\bar{x}^{(i)})^T f(\bar{x}^{(i)})}_{\|f(\bar{x}^{(i)})\|^2}$$

$$\|f(\bar{x}^{(i)})\|^2 > 0$$

Our update rule is sparse

$|V|$ vocab = 10k words

$f(\bar{x}^{(i)})$ has 4 words

we only update 4 weights

Step size



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

$$w = 1$$

$$\Rightarrow \text{grad} = 2$$

$$\Rightarrow \text{update} = -2\alpha \quad \alpha \text{ step size}$$

$$\alpha = 1$$

$$\Rightarrow 1 - 2 = -1$$

$$w = -1$$

$$\Rightarrow \dots 1$$

$$\alpha = 0.5 \quad \checkmark$$

For AI:

Schedule

Start with

$$\alpha = 1$$

then reduce

$$\alpha = \frac{1}{t} \text{ } t \text{ epoch}$$