

CS 378 Lecture 2

Classification 1: Features, Perceptron

Announcements

- AI released (due 9/8)
- Book notation diverges from lectures

Today

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

Classification Points $\bar{x} \in \mathbb{R}^n$

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

Classifier maps $\bar{x} \rightarrow y$

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

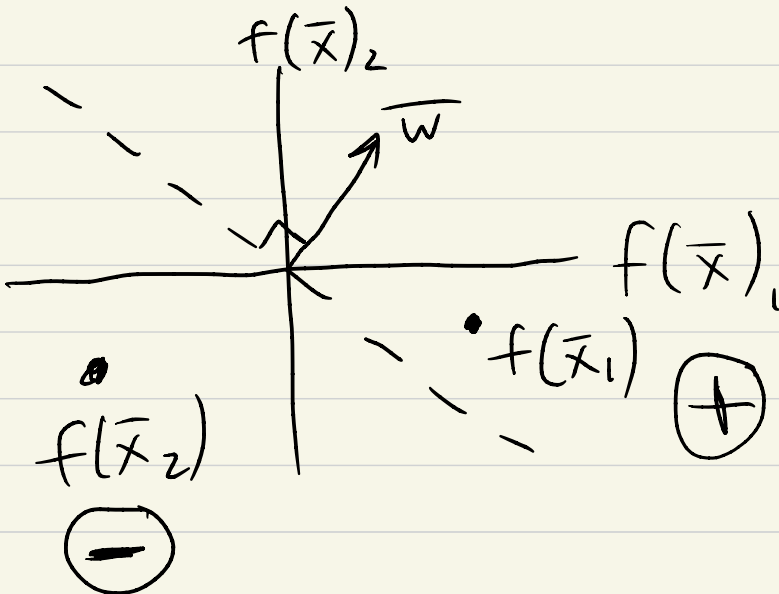


credit: Machine Learning Memes on Facebook

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

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

$\underbrace{\bar{w} \cdot f(\bar{x})}_{\text{a real number}}$ is greater than 0?



Sentiment Analysis

\bar{x} = the movie was great!
would watch again!

① Feature extraction

$$\bar{x} \Rightarrow f(\bar{x})$$

string \mathbb{R}^n

② Learning training set

$$\left\{ \left(f(\bar{x}^{(i)}) , y^{(i)} \right) \right\}_{i=1}^D \Rightarrow \bar{w} \in \mathbb{R}^n$$

D) labeled examples

Feature Extraction

\bar{X} = the movie was great

Bag-of-words featurization

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

Vocabulary of
 $\sim 10,000$ words

value is

the count of that word in \bar{X}

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

$\left[\begin{array}{cccc} -0.1 & +0.2 & \dots & +10 \\ \text{the} & \text{a} & & \text{great} \end{array} \right]$

Preprocessing

① Vocab selection: need a fixed set of words for the vector space

replace unseen words w/ UNK

② Tokenization

wasn't great! [| | ... 0]
wasn't great! great

typical tokenization:

- break out punctuation
- break out contractions

was n't great!

② Stopword filtering:

- prepositions

- a, the

- pronouns (maybe for debiasing)

③ Lowercasing / Stemming

Fix typos!

So far: unigram BoW

Bigram BoW

[1 1 0 ...]
the movie movie was not good

"Vocab" = vocab^2

Unigram: 10k

Bigram =

$(10k)^2$ in theory

1M in practice

Maintain an index

can combine =

[the: 0
a: 1
movie: 47]

[the: 0
:
was great: 1172
:
:]

Machine Learning

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

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

Find the best $\bar{w} \in \mathbb{R}^n$

Training objective = $\text{loss}(\text{dataset})$

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

"if we use \bar{w} as our weights, how badly do we screw up ex. (i) "

(sample i)

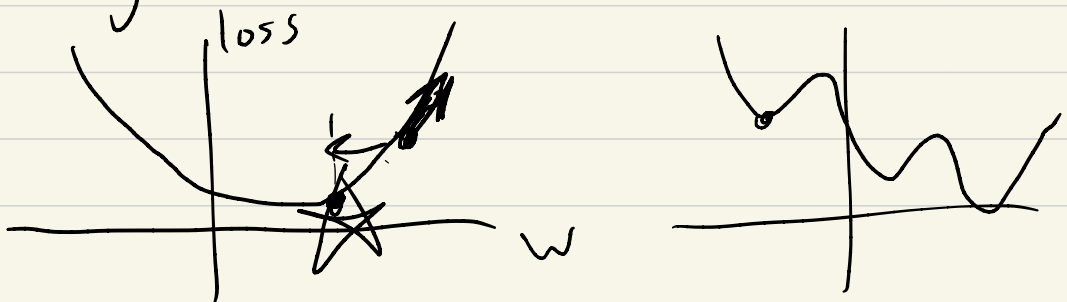
Stochastic gradient descent

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

for i in range $(0, D)$

$$\bar{w} \leftarrow \bar{w} - \underset{\substack{\uparrow \\ \text{step size} \\ \approx 1}}{\alpha} \cdot \frac{\partial}{\partial \bar{w}} \text{loss} \left(\bar{x}^{(i)}, y^{(i)}, \bar{w} \right)$$

Update \bar{w} by subtracting
gradient of the loss



Perceptron (instance of SGD)

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

for t in range(0, epochs)

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

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

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

Let $\alpha = 1$ for now

$y^{(i)} = -1$

\bar{w}

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

Suppose

$$\bar{w}^T f(\bar{x}^{(i)}) \Rightarrow -1.3$$

$$y^{(i)} = +1$$

After update:

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

$$\bar{w}^T f(\bar{x}^{(i)}) + \underbrace{f(\bar{x}^{(i)})^T f(\bar{x}^{(i)})}_{> 0}$$

larger than
-1.3