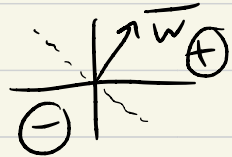


CS371N Lecture 3

Classification 2: Logistic Regression and Optimization

Announcements: AI due next Thurs

Recap Linear binary classification $y \in \{-1, +1\}$
 $\bar{w}^T f(\bar{x}) \stackrel{?}{>} 0$



features f : Bag of words

$f(\bar{x}) = [0 \quad 1 \quad 1 \quad 0 \quad 0 \quad 1]$
a the was of in -- movie--

\bar{x} = the movie was great

Perceptron: dataset $\{ (\bar{x}^{(i)}, y^{(i)}) \}_{i=1}^D$

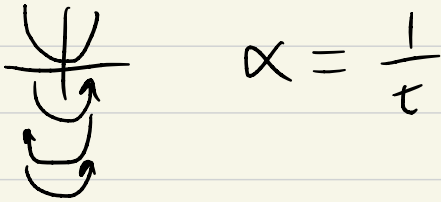
Init $\bar{w} = \bar{0}$

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

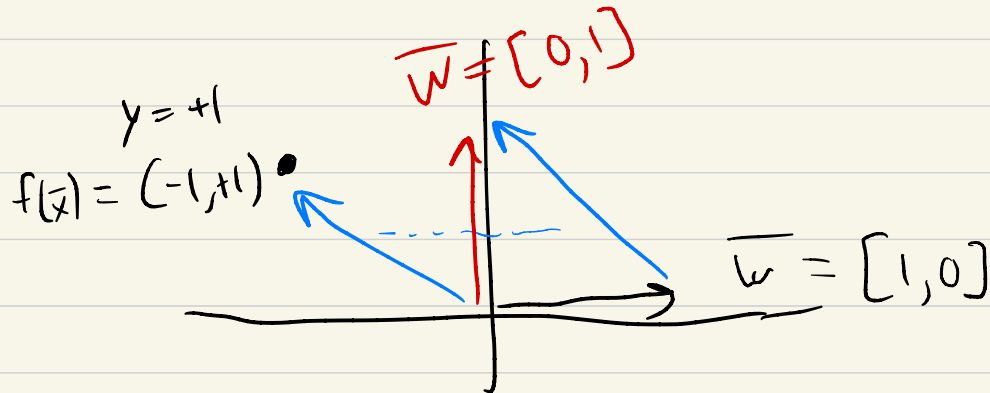
for i in range $(0, D)$

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

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



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



Ex $\bar{x}^{(1)}$: good

$$y = +1$$

$$\bar{w}^T f(L\bar{x}) > 0$$

if $0 \Rightarrow -1$

$\bar{x}^{(2)}$: not good

$$y = -1$$

$\bar{x}^{(3)}$: bad

$$y = -1$$

① Write the feature vectors

② Execute one epoch of perceptron

Start with $\bar{w} = 0$, go in order

Assume $x = 1$

\bar{x}	y	$f(\bar{x})$					
		g	n	b	ng	nb	
g	+1	[1	0	0]	0	0	
ng	-1	[1	1	0]	1	0	*
b	-1	[0	0	1]	0	0	

$$\bar{w} = [0 \ 0 \ 0]$$

Case 2, add $f(\bar{x})$

Ex 1 $y_{\text{pred}} = -1 \Rightarrow \bar{w} = [1 \ 0 \ 0] - \star \downarrow$

Ex 2 $[1 \ 0 \ 0] \begin{bmatrix} 1 \\ 1 \\ 0 \end{bmatrix} = 1, y_{\text{pred}} = 1 \Rightarrow \bar{w} = [0 \ -1 \ 0]$

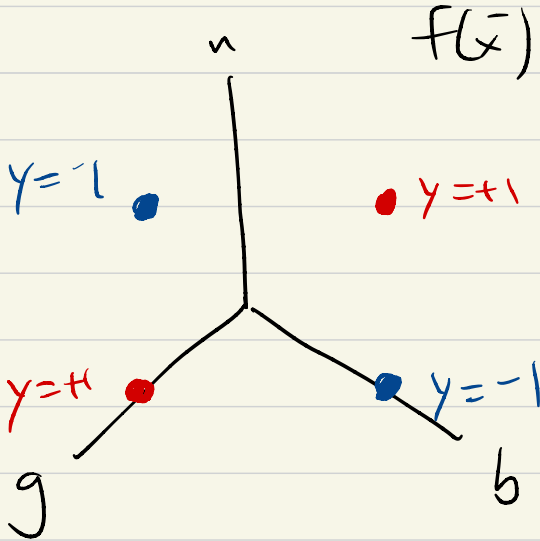
Ex 3 $[0 \ -1 \ 0] \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} = 0, y_{\text{pred}} = -1 \Rightarrow \text{no change}$

If we start epoch 2:

update $\Rightarrow [1 \ -1 \ 0]$ no more updates

Ex Add the example "not bad"

$n_b + 1$ $[0 \ 1 \ 1]$ $n_g \ n_b$
 $0 \ 1$



Counter and Indexer

$$\text{Indexer} = \left\{ \begin{array}{l} 0 \Leftrightarrow \text{"the"} \\ 1 \Leftrightarrow \text{"a"} \\ \vdots \\ 147 \Leftrightarrow \text{"not good"} \end{array} \right\}$$

Counter: sparse vector

$$\left\{ \begin{array}{l} 147 \Rightarrow 1 \\ 192 \Rightarrow 2 \end{array} \right\} \text{ feature vector}$$

Logistic Regression

Discriminative probabilistic model

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

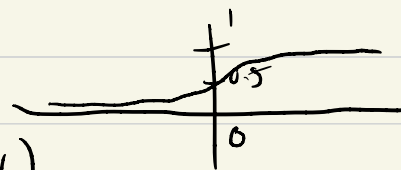
$$P(y=+1 | \bar{x}) = \frac{e^{\bar{w}^T f(\bar{x})}}{1 + e^{\bar{w}^T f(\bar{x})}}$$

$$\frac{e^z}{1 + e^z}$$

$$\frac{1}{1 + e^{-z}}$$

logistic function

maps $z \in \mathbb{R} \Rightarrow (0, 1)$



$$\bar{w}^T f(\bar{x}) \stackrel{?}{>} 0 \Leftrightarrow y = +1$$

$$P(y = +1 | \bar{x}) \stackrel{?}{>} 0.5$$

Learning Maximize data likelihood

$$\text{Likelihood } \mathcal{L} = \prod_{i=1}^D P(y = y^{(i)} | \bar{x}^{(i)})$$

Transform: likelihood \Rightarrow log likelihood

$$LL = \sum_{i=1}^D \log P(y = y^{(i)} | \bar{x}^{(i)})$$

$$\operatorname{argmax}_{\bar{w}} \mathcal{L}(\bar{w}) = \operatorname{argmax}_{\bar{w}} LL(\bar{w})$$

Actually: minimize negative LL

$$\operatorname{argmin}_{\bar{w}} \sum_{i=1}^D \underbrace{-\log (P(y = y^{(i)} | \bar{x}^{(i)}; \bar{w}))}_{\text{log loss}}$$

$$\text{SGD} = \frac{\partial}{\partial \bar{w}} \text{loss}(\bar{x}^{(i)}, y^{(i)}, \bar{w})$$

This gradient \approx perceptron update

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

$$\frac{\partial}{\partial \bar{w}} -\log P(y = +1 | \bar{x})$$

$$= \frac{\partial}{\partial \bar{w}} -\log \left[\frac{e^{\bar{w}^T f(\bar{x})}}{1 + e^{\bar{w}^T f(\bar{x})}} \right]$$

$$= \frac{\partial}{\partial \bar{w}} \left[-\bar{w}^T f(\bar{x}) + \log(1 + e^{\bar{w}^T f(\bar{x})}) \right]$$

$$= -f(\bar{x}) + \frac{1}{1 + e^{\bar{w}^T f(\bar{x})}} \cdot e^{\bar{w}^T f(\bar{x})} \cdot f(\bar{x})$$

$$= f(\bar{x}) \left(-1 + \frac{e^{\bar{w}^T f(\bar{x})}}{1 + e^{\bar{w}^T f(\bar{x})}} \right)$$

$$= f(\bar{x}) \left(-1 + P(y = +1 | \bar{x}) \right)$$

$$-\frac{\partial}{\partial w} = f(\bar{x})(1 - P(y=+1|\bar{x}))$$

In SGD: we add $-\frac{\partial}{\partial w}$ to \bar{w}

If $P(y=+1|\bar{x}) \approx 1$, what happens?
 \approx no change

$$\underline{P(y=+1|\bar{x}) \approx 0?} + f(\bar{x})$$

If $P(y=+1|\bar{x}) = 0.5$ "half an update"



differs from perceptron

Update

$$y^{(i)} = +1: \bar{w} \leftarrow \bar{w} + \alpha f(\bar{x}^{(i)}) (1 - P(y=+1|\bar{x}^{(i)}))$$

$$y^{(i)} = -1: \bar{w} \leftarrow \bar{w} - \alpha f(\bar{x}^{(i)}) \underbrace{(1 - P(y=-1|\bar{x}^{(i)}))}_{\substack{1 - P(y=-1) \\ = P(y=+1)}}$$

Step size:

$\frac{1}{t}$, $\frac{1}{\sqrt{t}}$... others possible

SGD = first-order