

CS371N Lecture 3

Classification 2: Logistic Regression and Optimization

Announcements - AI due in 9 days

Recap Linear binary classifier: $\bar{w}^T f(\bar{x}) \stackrel{?}{>} 0$

Bag-of-words featurization:

\bar{x} = the movie was great

$f(\bar{x}) = \begin{bmatrix} 0 & 1 & 1 & 0 & 0 & \dots & 1 & \dots \end{bmatrix}$
a the was of in movie
(4 1s)

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

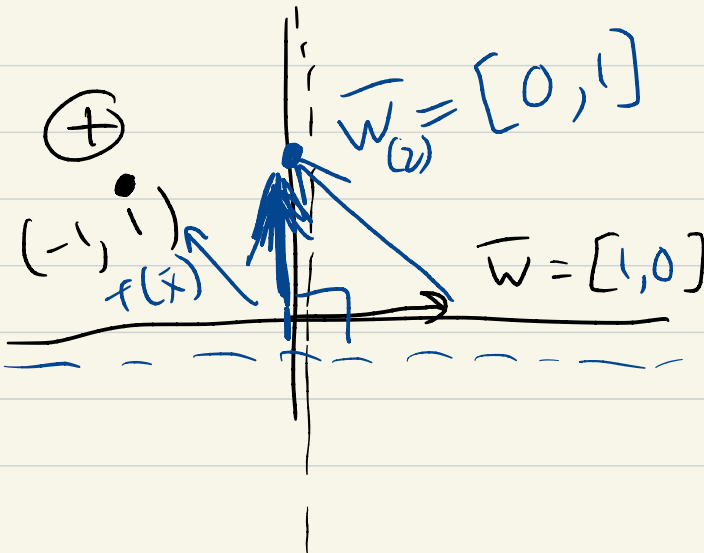
init $\bar{w} = 0$

for t in range(0, epochs)

for i in range(0, D)

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

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



Example

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

if $0 \Rightarrow -1$

\bar{x} : good

$y = +1$

not good

$y = -1$

bad

$y = -1$

$\alpha = 1$

① Write the feature vectors (x3)

② Execute one epoch of perceptron

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

Give final weight vector

	y	feats	n_b	n_g
g	+1	$\begin{bmatrix} 1 & 0 & 0 \end{bmatrix}$	0	0
ng	-1	$\begin{bmatrix} 1 & 1 & 0 \end{bmatrix}$	0	1
b	-1	$\begin{bmatrix} 0 & 0 & 1 \end{bmatrix}$	0	0

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

Ex 1 : $y_{\text{pred}} = -1$

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

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

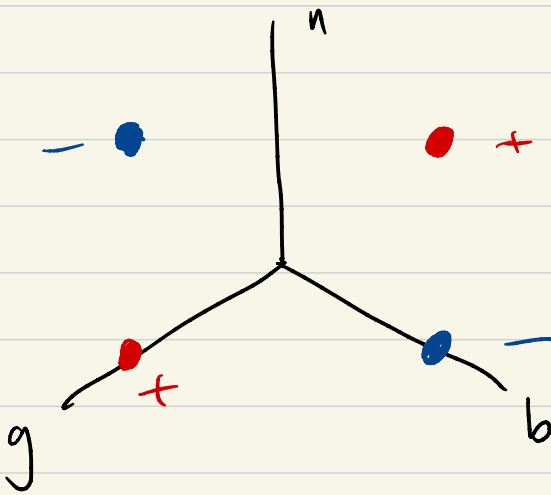
Ex 3 : $[0 \ -1 \ 0] \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} = 0$ $y_{\text{pred}} = -1$ no change

If we start epoch 2:

Update $\Rightarrow [1 \ -1 \ 0]$ converged

Ex Add the example "not bad"

nb +1 [0 1 1] $\begin{matrix} \text{nb} & \text{ng} \\ 1 & 0 \end{matrix}$



Logistic Regression

Discriminative probabilistic model

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

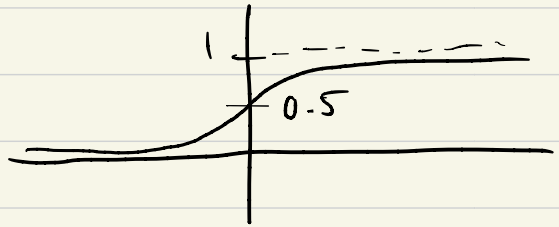
label features/instance

(generative: $P(\bar{x}, y)$)
Naive Bayes

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

$$\frac{e^z}{1 + e^z} \quad \left(\frac{1}{1 + e^{-z}} \right)$$

logistic
fcn



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

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

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

equivalent to
our earlier
decision
rule

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

Learning Maximize the data likelihood

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

$$\text{We want } \bar{w}^* = \underset{\bar{w}}{\text{argmax}} \underline{L(\bar{w})}$$

① Likelihood \Rightarrow log likelihood (LL)

$$\underset{\bar{w}}{\text{argmax}} \sum_{i=1}^D \log P(y=y^{(i)} | \bar{x}^{(i)})$$

log is monotonic



(NLL)
② Minimize the negative LL

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

$$\text{loss}(\bar{x}^{(i)}, y^{(i)}, \bar{w})$$

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

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

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

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

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

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

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

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

Scaling factor

What if $P(y=+1 | \bar{x}) \approx 1$ ("ypred")

Little update (but still nonzero)

What if $P(y=+1 | \bar{x}) \approx 0$? Essentially perc-update

$P(y=+1 | \bar{x}) = 0.50001$?

Perc: no update

HPerc: update!

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

Update:

step size
↓

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

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

$\underbrace{\hspace{10em}}_{P(y = +1 | \bar{x}^{(i)})}$

Still want a step size.

Consider constant, $\frac{1}{t}$, $\frac{1}{\sqrt{t}}$, ... t epoch number

SGD is first-order optimization

Newton's method is second-order

second deriv: Hessian $n \times n$ matrix

$$\frac{\partial^2}{\partial w_i \partial w_j}$$