

CS 378 Lecture 3

Classification 2: Logistic Regression, Optimization

Announcements

- AI
- Videos/readings

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

Bag-of-words featurization:

\bar{x} = the movie was great

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

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

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

Example

\bar{x} : good $y: +1$
not good $y: -1$
bad $y: -1$

$$\bar{w}^T f(\bar{x}) > 0 \\ \text{if } = 0 \Rightarrow -1$$

① Write the feature vector for each

② Execute 1 epoch of perceptron

	y	feats		
		g	b	n
g	+1	[1	0	0]
ng	-1	[1	0	1]
b	-1	[0	1	0]

↓

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

$$\text{Ex 1: } \bar{w}^T \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \Rightarrow 0 \quad (-1)$$

$$\bar{w} \leftarrow \bar{w} + 1 [1 \ 0 \ 0] \Rightarrow [1 \ 0 \ 0]$$

$$\text{Ex 2: } [1 \ 0 \ 0] \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} \Rightarrow 1 \quad (+1)$$

$$\bar{w} \leftarrow \bar{w} - 1 [1 \ 0 \ 1] \Rightarrow [0 \ 0 \ -1]$$

Ex 3: Correct, no change

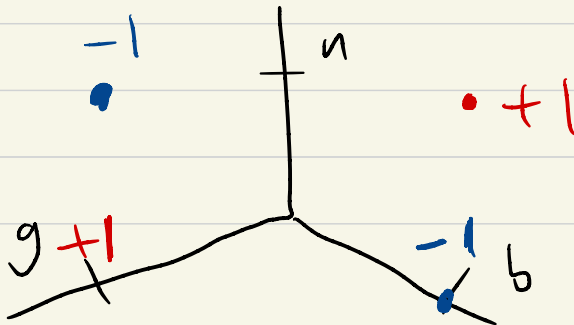
If we do ex 1 again:

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

Example 2

		g	b	n	ng
good	+1	1	0	0	0
not good	-1	1	0	1	1
bad	-1	0	1	0	0
not bad	+1	0	1	1	0

$[1 \ 0 \ -1]$ $\leftarrow -1 \Rightarrow \textcircled{-1} \ y_{\text{pred}}$



Logistic Regression

Discriminative probabilistic model

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

label features / instance

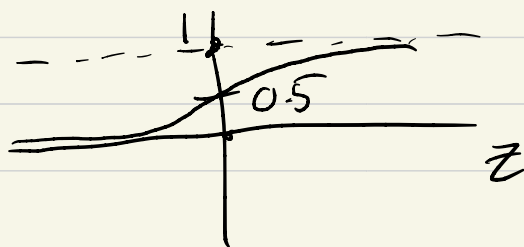
(generative: $P(\bar{x}, y)$)

Naive Bayes

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

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

logistic
fcn



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

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

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

Learning

Maximize the data
likelihood

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

We want $\bar{w} = \underset{\bar{w}}{\operatorname{argmax}}$ data likel.

✂

$\Rightarrow \operatorname{argmax}_{\bar{w}} \log \text{ data likelihood}$

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

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

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

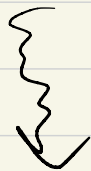
$$\text{SGD: } \frac{\partial}{\partial \bar{w}} \text{loss}(\bar{x}^{(i)}, y^{(i)}, \bar{w})$$

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

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

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

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



Update:

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

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

Update:

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

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

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

What happens if $P(y = +1 | \bar{x})$
on this ex. is 1?

\bar{w} Doesn't change

$$P(y = +1 | \bar{x}) \text{ is } 0?$$

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

LR never
"converges"
exactly

$$P(y = +1 | \bar{x}) \text{ is } 0.7? \quad + 0.3\alpha f(\bar{x}^{(i)})$$

Optimization

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

$\mathcal{L}(\bar{w})$ loss on the whole dataset

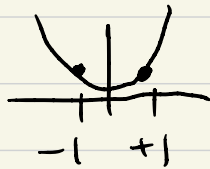
$$\left(\sum \mathcal{L}(\bar{x}^{(i)}, y^{(i)}, \bar{w}) \right)$$

$$\text{SGD: } \bar{w} \leftarrow \bar{w} - \alpha \frac{\partial}{\partial \bar{w}} \mathcal{L}(\bar{x}^{(i)}, y^{(i)}, \bar{w})$$

step size

$$\mathcal{L}(w) = w^2$$

$$w = -1$$



SGD: loops infinitely

$$w = -1$$

$$w \leftarrow w - \alpha \cdot 2w$$

$$\frac{\partial}{\partial w} \mathcal{L} = 2w$$

$$\alpha = 1: w = -1 \rightarrow +1$$

$$\alpha = 0.5: w = -1 \rightarrow 0 \quad \checkmark$$

How to choose step size?

— Start with 1, decrease it
 $1/t$ $t = \text{epoch}$

(may decrease too fast)

$\frac{1}{\sqrt{t}}$, e^{-t} is even faster

2nd deriv: Newton

quasi-Newton

Adagrad, Adam, AdamW

approximate 2nd-order