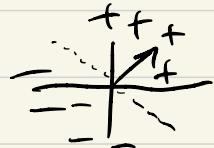


CS 378 Lecture 3

Classification 2: Logistic Regression, Optimization

Announcements

- AI Seating chart
- Video on course website



Recap Linear binary classifier: $\bar{w}^T f(\bar{x}) \geq 0$
Bag-of-words featurization

\bar{x} = movie was great

$$\Rightarrow f(\bar{x}) = [0 \ 1 \ 0 \ 0 \dots \ 1 \ \dots \ 1]$$

the was a of movie great

Perceptron: initialize $\bar{w} = \bar{0}$

for + in range(0, epochs)

for i in range(0, D)

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

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

Example

	y	g	b	n	feats
good	+1	[1	0	0]	$= f(\bar{x}^{(1)})$
not good	-1	[1	0	1]	
bad	-1	[0	1	0]	

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

$$Ex \ 1: \ \bar{w}^\top f(\bar{x}^{(1)}) = 0 > 0 \stackrel{?}{\Rightarrow} y_{pred} = -1$$

$$\bar{w} \leftarrow \bar{w} + \alpha \cdot [1 \ 0 \ 0] \quad \alpha = 1$$

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

$$Ex \ 2: \ \bar{w}^\top f(\bar{x}^{(2)}) = 1 > 0 \stackrel{?}{\Rightarrow} y_{pred} = +1$$

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

$$\bar{w} = [0 \ 0 \ -1] \quad \text{et again:}$$

$$Ex \ 3: \text{correct} \quad \text{After} \quad [1 \ 0 \ -1]$$

Today

- Logistic regression
- Optimization
- Sentiment systems (if time)

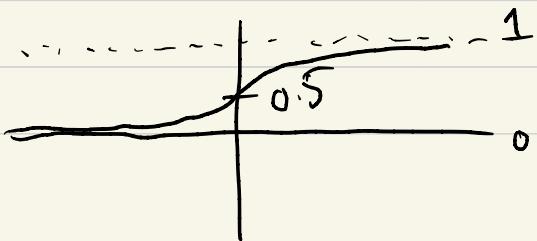
Logistic Regression

Discriminative probabilistic model

Models a distribution $P(y|\bar{x})$
generative $P(\bar{x}, y)$ label instance
(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 function



Decision boundary: predict +1 if

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

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

equivalent to

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

in order for it to be
a dist. over y

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

Learning

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

maximize the data likelihood

$$\max_{\bar{w}} \prod_{i=1}^D P(y^{(i)} | \bar{x}^{(i)})$$

logistic reg. w/weights \bar{w}

$$\Rightarrow \max_{\bar{w}} \log \prod_{i=1}^D P(y^{(i)} | \bar{x}^{(i)})$$

~~\log~~

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

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

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

For stochastic gradient descent:

$$\text{need } \frac{\partial}{\partial \bar{w}} \log (\bar{x}, y, \bar{w})$$

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

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

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

CALCULUS

Logistic regression update

$$\bar{w} \leftarrow \bar{w} - \alpha \frac{\partial}{\partial \bar{w}} \text{loss}$$

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

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

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

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

① Think about when $y^{(i)} = +1$

What happens when:

not much update!

$P(y^{(i)} = +1 | \bar{x})$ is close to 1?

close to 0? perc update

close to 0.5?

Optimization

$$\sum_{i=1}^D \text{loss}(\bar{x}^{(i)}, y^{(i)}, \bar{w}) \triangleq \mathcal{L}\left((\bar{x}^{(i)}, y^{(i)})_{i=1}^D | \bar{w}\right)$$

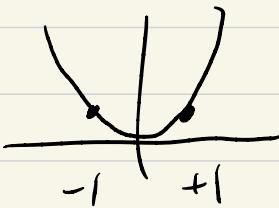
function of \bar{w}

$$\mathcal{L}(\bar{w})$$

$$\text{SGD: } \bar{w} \leftarrow \bar{w} - \alpha \frac{\partial}{\partial \bar{w}} \mathcal{L}(\bar{w})$$

$$L(x) = x^2$$

$$x = -1$$



$$\frac{\partial}{\partial x} L = 2x$$

Suppose $\alpha = 1$

SGD: start at $x = -1$

$$\Rightarrow x \leftarrow x - 1 \cdot \frac{\partial}{\partial x} L(x)$$
$$\leftarrow -1 - 1 \cdot (-2)$$

$$x = 1$$

Oscillates with $\alpha = 1$

If $\alpha = 1/2$, converges instantly

How to choose step size?

- Different constants
- Large \rightarrow small, e.g. $\frac{1}{t}$ where
 $t = \text{epoch number}$
- $\frac{1}{\sqrt{f}}$, ... other options

e^{-t} : decreases too fast

How to do something smarter?

Newton: $\bar{w} \leftarrow \bar{w} - \underbrace{\left(\frac{\partial^2 L}{\partial w^2} \right)^{-1} \frac{\partial L}{\partial w}}$

inverse Hessian

Approximate 2nd-order
methods: Adagrad,

Adadelta, Adam, Adam^W ...

$n \times n$ matrix
 $n = \# \text{ of features}$
parameters

Regularization: not used