10-301/601: Introduction to Machine Learning Lecture 17.5 - Naïve Bayes Predictions

Henry Chai 11/1/22

## Bernoulli Naïve Bayes

- Binary label
  - $Y \sim \text{Bernoulli}(\pi)$

• 
$$\hat{\pi} = N_{Y=1}/N$$

- $\sim$  N = # of data points
  - $N_{Y=1}$  = # of data points with label 1
- Binary features
  - $X_d | Y = y \sim \text{Bernoulli}(\theta_{d,y})$

$$\bullet \ \widehat{\theta}_{d,y} = \frac{N_{Y=y,X_d=1}}{N_{Y=y}}$$

- $\vec{N}_{Y=y}$  = # of data points with label y
  - $N_{Y=y, X_d=1}$  = # of data points with label y and feature  $X_d=1$

What if some
Wendelabel
Phair reever
Appears in our
thanking data?
Predictions

• Given a test data point  $\mathbf{x}' = [x_1', ..., x_D']^T$ P(Y=1|X') or P(X'|Y=1) P(Y=1) $= \left(\frac{1}{\pi} P(x_d'|Y=1)\right) P(Y=1)$  $P(Y=0|X') \propto (\frac{\pi}{1-\theta_{1,1}})^{1-X_{1}} (1-\frac{\pi}{1-x_{1}})^{1-X_{1}} (1-\frac{\pi}{1-x_{1}}$  $\hat{\gamma} = \begin{cases} 1 & \text{if } P_1 \geq P_0 \\ 0 & \text{otherwise} \end{cases}$ 

## What if some Word-Label pair never appears in our training data?

<i>x</i> <sub>1</sub> ("hat")	x <sub>2</sub> ("cat")	x <sub>3</sub> ("dog")	x <sub>4</sub> ("fish")	x <sub>5</sub> ("mom")	<i>x</i> <sub>6</sub> ("dad")	<i>y</i> (Dr. Seuss)
1	1	0	0	0	0	1
0	0	1	0	0	0	0
0	0	0	1	0	0	1
0	0	0	0	1	0	0

The Cat in the Hat gets a Dog (by ???)

- If some  $\hat{\theta}_{d,y} = 0$  and that word appears in our test data x', then P(Y = y | x') = 0 even if all the other features in x' point to the label being y!
- The model has been overfit to the training data
- We can address this with a prior over the parameters!

## Setting the Parameters via MAP

- Binary label
  - $Y \sim \text{Bernoulli}(\pi)$

• 
$$\hat{\pi} = \frac{N_{Y=1}}{N}$$

- N = # of data points
- $N_{Y=1}$  = # of data points with label 1
- Binary features

• 
$$X_d | Y = y \sim \text{Bernoulli}(\theta_{d,y}) \text{ and } \theta_{d,y} \sim \text{Beta}(\alpha, \beta)$$

$$\hat{\theta}_{d,y} = \frac{N_{Y=y,X_{d=1}} + (\alpha - 1)}{N_{Y=y} + (\alpha - 1) + (\beta - 1)}$$

- $N_{Y=y}$  = # of data points with label y
- $N_{Y=y, X_d=1}$  = # of data points with label y and feature  $X_d=1$
- Common choice:  $\alpha = 2$ ,  $\beta = 2$