

Outline

Naive Bayes

- Laplace Smoothing
- Event models

Comments on applied ML

Kernel Methods

Recap:

$$X = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ \text{aardvark} \\ \text{buy} \end{bmatrix} \quad X_j = \mathbb{1}_{\{\text{word } j \text{ appears in email}\}}$$

Generative Model

$$P(X|y) = \prod_{i=1}^d p(x_i|y)$$

$$\text{Parameters: } P(y=1) = \phi_y$$

$$P(X_j=1 | y=0) = \phi_{j|y=0}$$

$$P(X_j=1 | y=1) = \phi_{j|y=1}$$

Max Likelihood estimates

$$\phi_y = \frac{\sum_{i=1}^n \mathbb{1}_{\{y^{(i)}=1\}}}{n}$$

$$\phi_{j|y=0} = \frac{\sum_{i=1}^n \mathbb{1}_{\{x_j^{(i)}=1, y^{(i)}=0\}}}{\sum_{i=1}^n \mathbb{1}_{\{y^{(i)}=0\}}}$$

$$\phi_{j|y=1}$$

At prediction time

$$P(y=1|x) = \frac{P(x|y=1) \cdot P(y=1)}{P(x|y=1) \cdot P(y=1) + P(x|y=0) \cdot P(y=0)} = \frac{0}{0}$$

NeurIPS

$$j = 5500$$

$$P(X_{5500} = 1 | y=1) = \frac{0}{\#\{y=1\}} = 0 = \phi_{5500|y=1}$$

$$P(X_{5500} = 1 | y=0) = \frac{0}{\#\{y=0\}} = 0 = \phi_{5500|y=0}$$

$$P(x|y=1) = \prod_{j=1}^{10,000} P(x_j|y=1) \rightarrow \phi_{5500|y=1}$$

$$P(x|y=0) = \prod_{j=1}^{10,000} P(x_j|y=0) \rightarrow \phi_{5500|y=0}$$

Laplace Smoothing

2009

		Won?
9/12	Wakeforest	O
10/10	OSU	O
10/17	Arizona	O
11/21	Caltech	O
12/31	Oklahoma	O

$$\begin{aligned} P(x=1) &= \frac{\#1's}{\#1's + \#0's} & +1 \\ &= \frac{0}{0+4} & \frac{1}{6} \\ &= 0 & = \frac{1}{6} \end{aligned}$$

Laplace Smoothing. $\frac{\#1's + 1}{\#0's + 1}$

More generally $X \in \{1 \dots k\}$

$$\text{Estimate } P(X=j) = \frac{\sum_{i=1}^n \mathbb{1}_{\{X^{(i)}=j\}} + 1}{n+k}$$

$$\sum_j P(X=j) = 1$$

$$P_{j|y=0} = \frac{\sum_{i=1}^n \mathbb{1}_{\{X_j^{(i)}=1, y^{(i)}=0\}} + 1}{\sum_{i=1}^n \mathbb{1}_{\{y^{(i)}=0\}} + 2} = P(X_j=1|y=0)$$

$P(X_j=1|y=0) + P(X_j=0|y=0) = 1$

$$X_i \in \{1 \dots k\}$$

size	$< 400 \text{ feet}^2$	$400-800$	$800-1200$	> 1200
X	1	2	3	4

$$P(x|y) = \prod_{i=1}^d \underbrace{P(X_j|y)}_{\text{multinomial (vs. bernoulli)}}$$

$$\text{Is } P(X_j=1) + P(X_j=0) = 1$$

$$P(X_j=1|y=0) + P(X_j=0|y=0) = 1$$

$$X = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \begin{array}{l} \text{aardvark} \\ \text{account} \\ \text{bank} \\ \text{beneficiary} \end{array}$$

$$X_i \in \{0, 1\}$$

account bank account ...

$$800 \quad 1600 \quad 800$$

$$x = \begin{bmatrix} 800 \\ 1600 \\ 800 \\ \vdots \\ \vdots \end{bmatrix} \in \mathbb{R}^{d_i} \quad d_i = \text{length of email } i$$

$x_j \in \{1, \dots, 10,000\}$

So far: Multivariate Bernoulli event model

New: Multinomial event model

Generative model

$$P(x, y) = P(x|y) \cdot P(y)$$

$$P(x|y) = \prod_{j=1}^{d_i} P(x_j|y)$$

Parameters

$$\phi_y = P(y=1)$$

$$\phi_{k|y=0} = \underbrace{P(X_j=k|y=0)}_{\text{chance of word } j \text{ being } k \text{ if } y=0}$$

Assume that this does not depend on j

$$\phi_{k|y=1} = P(X_j=k|y=1)$$

MLE

$$\hat{\phi}_{k|y=0} = \frac{\sum_{i=1}^n \mathbb{1}_{\{y^{(i)}=0\}} \sum_{j=1}^{d_i} \mathbb{1}_{\{x_j^{(i)}=k\}}}{\sum_{i=1}^n \mathbb{1}_{\{y^{(i)}=0\}} \cdot d_i} \quad \begin{array}{l} +1 \\ \hline +10,000 \end{array}$$

Laplace Smoothing

$x_j = \frac{\text{index of } j^{\text{th}} \text{ word in email}}{\text{position in the dictionary}}$

map rare words to UNK

mortgage mφrtgʌnge
 ↓
 UNK

$$P_{k|y=1} = P(j^{\text{th}} \text{ word} = k | y = 1)$$

↓
independent of j

$$P(x|y) = \prod_{i=1}^d P(x_i|y)$$

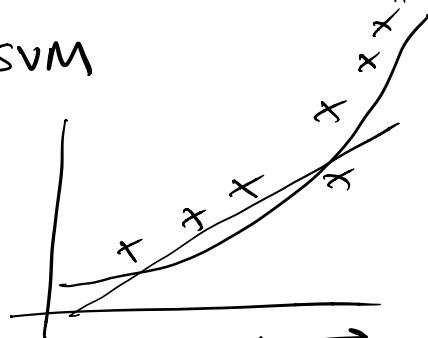
$$= P(x_1=1|y) \cdot P(x_2=0|y) \cdot P(x_3=1|y) \cdots$$

$$x \begin{bmatrix} 1 \\ 0 \\ 1 \\ \vdots \\ 1 \end{bmatrix}$$

$$P(x|y) = \prod_{i=1}^d P(x_i|y)$$

Kernel Methods → SVM

linear f^n
 quadratic f^n
 cubic f^n of data



$$h_\theta(x) = \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ x^3 \end{bmatrix} \quad \phi: \mathbb{R} \rightarrow \mathbb{R}^4$$

$$h_\theta(x) = [\theta_0 \ \theta_1 \ \theta_2 \ \theta_3] \begin{bmatrix} 1 \\ x \\ x^2 \\ x^3 \end{bmatrix}$$

$$= \theta^T \phi(x)$$

$$(x^{(1)}, y^{(1)}) \dots (x^{(n)}, y^{(n)})$$

↓

$$(\phi(x^{(1)}), y^{(1)}) \dots (\phi(x^{(n)}), y^{(n)})$$

Recall: linear regression

$$\theta = 0$$

$$\text{Loop } \theta := \theta + \alpha \sum_{i=1}^n \underbrace{(y^{(i)} - \theta^T x^{(i)})}_{\text{scalar}} x^{(i)}$$

$\uparrow \mathbb{R}^d \quad \uparrow \mathbb{R}^d$

New data set

$$\theta = 0$$

$$\text{Loop } \theta := \theta + \alpha \sum_{i=1}^n \underbrace{(y^{(i)} - \underbrace{\theta^T \phi(x^{(i)})}_{\text{scalar}})}_{\mathbb{R}^P} \underbrace{\phi(x^{(i)})}_{\mathbb{R}^P}$$

$$\phi: \mathbb{R}^d \rightarrow \mathbb{R}^P \quad (d=1, P=4 \text{ in example})$$

Terminology

ϕ : feature map

$\phi(x)$: (new) feature

x : attributes

Kernel methods:

$$d > 1 \quad x = (x_1 \dots x_d)$$

$$\phi(x) = \begin{bmatrix} 1 \\ x_1 \\ x_d \\ x_1 x_1 \\ x_1 x_2 \\ \vdots \\ x_d x_d \end{bmatrix} \quad \left| \begin{array}{c} 1 \\ \} d \\ \} d^2 \end{array} \right.$$

$\theta^T \phi(x)$ can represent any degree 3 poly in $x_1 \dots x_d$

$P = 1 + d + d^2 + d^3$
 Suppose $d = 1000$
 $P \approx 10^9$

$$\begin{bmatrix} x_1 x_1 x_1 \\ x_1 x_1 x_2 \\ \vdots \\ x_d x_d x_d \end{bmatrix} \quad \left\} d^3 \right.$$

time per iteration $O(n^p)$ $p \sim d^3$

Improve this to $O(n^2)$ per iteration

Even storing Θ takes time $O(p)$
 $p \gg n$