- Naive Bayes
  - Laplace Smoothing
  - Event Models
- Comment on applied ML
- Kernel Methods

Recap: Spam filter Naive Bayes

$$X = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ \vdots \end{bmatrix} \begin{matrix} a \\ add \\ bee \end{matrix} \qquad X_j = \mathbb{1}\{\text{work } j \text{ appears in email}\}$$

### Generative Model

$$P(x|y) \qquad P(y)$$

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

### Parameters:

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

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

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

### Max Likelihood estimates

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

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

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

### Prediction time

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

If word $x_{5500}$ never occurs.

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

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

$$P(x | y = 1) = \prod_{j=1}^{10000} P(x_j | y = 1) = 0 \quad \Bigg\} \Rightarrow \quad P(y=1 | x) = \frac{0}{0+0}$$

$$P(x | y = 0) = \prod P(x_j | y=0) = 0$$

## Laplace   Smoothing

$$\# 1's \quad +1$$
$$\# 0's \quad +1$$

$$x \in \{1, ..., k\}$$

Estimate $P(x = j) = \dfrac{\sum_{i=1}^{n} \mathbb{1}\{x^{(i)} = j\} + 1}{n + k}$

~ Space size

$$\sum_{j=1}^{n} P(x = j)$$

· Back  to  Naïve  Bayes

$$\phi_{j | y=0} = \frac{\sum_{i=1}^{n} \mathbb{1}\{x_j^{(i)} = 1, \; y = 0\} + 1}{\sum_{i=1}^{n} \mathbb{1}\{y^{(i)} = 0\} + 2}$$

$$x_i \in \{1, ..., k\}$$

size $X$ $\left\langle \begin{array}{cccc} 400 \text{ feet}^2 & 400\text{-}800 & 800\text{-}1200 & >1200 \\ 1 & 2 & 3 & 4 \end{array} \right\rangle$  Discrete variable.

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

Multinomial (vcs. bernoulli)

$$X = \begin{bmatrix} 800 \\ 1600 \\ 800 \end{bmatrix}$$

account   bank   account...
800    1600    800

$\in \mathbb{R}^{di}$, $d_i$ = length of email $\underline{i}$,   $x_j \in \{1, ..., 10,000\}$   ← dict size

So far:  Multivarite  Bernoulli  event  model

New:    Multinomial    event  model

Generative    model

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

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

$P(x|y), P(y)$  generative    $P(y|x)$ .  discriminative

Parameters

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

$$\phi_{k|y=0} = P(x_j = k | y = 0)$$

chance of word $j$ being $k$ if $y=0$

Assume that does not depend on $j$ (independent of $j$)

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

MLE

# of $\{x_j = k$ when $y=0\}$

$$\phi_{k|y=0} = \frac{\left(\sum_{i=1}^{n} \mathbb{1}\{y^{(i)} = 0\} \sum_{j=1}^{d_i} \mathbb{1}\{x_j^{(i)} = k\}\right) + 1}{\left(\sum_{i=1}^{n} \mathbb{1}\{y^{(i)} = 0\} \cdot d_i\right) + 10,000}$$

(Laplace Smoothing)

# of $\{$words, when $y=0\}$

$x_j$ : index of the $j$-th word in the email

Original.

$$P(x|y) = \prod_{i=1}^{d} P(x_j | y) = P(x_1 = 1 | y) P(x_2 = 0 | y) P(x_3 = 1 | y=0) \cdots \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \end{bmatrix} \overset{x}{}$$

New model

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

Kernel Methods → SVM

linear fn.
quadratic fn.
cubic fn.



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

$$\phi(x) = \begin{bmatrix} 1 \\ x \\ x^2 \\ x^3 \end{bmatrix} \quad \text{(feature mapping)}$$

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

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

$\Downarrow$

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

Recall: linear regression

$\theta : 0$

Loop:
$$\theta := \underset{\mathbb{R}^d}{\theta} + \alpha \sum \underbrace{(y^{(i)} - \theta^T x^{(i)})}_{\text{scalar}} \underset{\mathbb{R}^d}{x^{(i)}}$$

New data set (after feature map)

$\theta : 0$

Loop:
$$\theta := \theta + \sum_{i=1}^{n} \underset{\mathbb{R}^p}{\underbrace{(y^{(i)} - \theta^T \phi(x^{(i)}))}_{\text{scalar}}} \underset{\mathbb{R}^p}{\underbrace{\phi(x^{(i)})}_{\mathbb{R}^p}} \qquad O(np)$$

$\phi : \mathbb{R}^d \to \mathbb{R}^p$

**Terminology**

$\phi$ : feature map

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

$x$ : attributes

**Kernel Method**

$d > 1. \qquad x = (x_1, \dots x_d)$

$$\phi(x) = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_1 x_1 \\ x_1 y_2 \\ x_1 y_1 x_1 \\ \vdots \\ x_d x_d x_d \end{bmatrix} \begin{matrix} \}1 \\ \}d \\ \\ \}d^2 \\ \\ \}d^3 \end{matrix}$$

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

$p = 1 + d + d^2 + d^3$

Suppose $d = 1000, \qquad p \cong 10^9$ Bad!

Time per iteration
$O(np) \qquad p \sim d^3$

$\xrightarrow[\text{Method}]{\text{Kernel}}$

Improve to $O(n^2)$
Even $\theta$ takes $O(p)$
$p \gg n$