Outline

Naive Bayes
— Laplace Smoothing

— Event Models

Kernel Methods

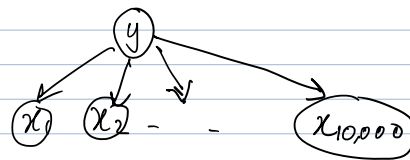Recap:

$$X = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 1 \end{bmatrix} \begin{matrix} a \\ \text{aardvark} \\ \vdots \\ \text{buy} \end{matrix} \quad \updownarrow d \qquad n \text{ examples}$$

$X_j = \mathbb{1}\{\text{word } j \text{ appears in email}\}$

Generative Model

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

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

$y = \begin{cases} 0 & \text{not spam} \\ 1 & \text{spam} \end{cases}$



Parameters

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

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

$\phi_y = P(y = 1)$

Joint Likelihood

$$\mathcal{L}(\phi_y, \phi_{j|y}) = \prod_{i=1}^{n} P(x^{(i)}, y^{(i)}; \phi_y, \phi_{j|y})$$

MLE

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

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

Prediction:

$$P(y=1|x) = \frac{\overbrace{P(x|y=1)}^{\phi_{j|y=1}} \cdot P(y=1)}{P(x|y=1)P(y=1) + \underbrace{P(x|y=0)}_{\phi_{j|y=0}} \cdot P(y=0)}$$

COVID    $j = 1273$

$$P(X_{1273}=1|y=1) = \frac{0}{\#\{y=1\}} = \phi_{1273|y=1}$$

$$P(X_{1273}=1|y=0) = \frac{0}{\#\{y=0\}} = \phi_{1273|y=0}$$

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

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

$\phi_{1273|y=1}$

$\phi_{1273|y=0}$

| | Won? |
|---|---|
| Wake forest | 0 |
| OSU | 0 |
| Arizona | 0 |
| Caltech | 0 |
| Oklahoma | ?? |

$$P(x=1) = \frac{\#\ "1"s \quad {\color{red}+1}}{\#\ "1"s + \#\ "0"s \quad {\color{red}+2}}$$

$$= \frac{0}{0+4} \quad {\color{red}\frac{+1}{+2}} \quad = {\color{red}\frac{1}{6}}$$

$$= 0$$

## Laplace Smoothing       $\#\ "1"s + 1$
$\#\ "0"s + 1$

$$X_i \in \{1 \cdots |V|\}$$

| size | $< 400$ feet$^2$ | 400–800 | 800–1200 | >1200 |
|---|---|---|---|---|
| X | 1 | 2 | 3 | 4 |

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

multinomial (vs. bernoulli)

$$X = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}$$

| a | 1 |
| aardvark | 2 |
| account | 800 |
| bank | 1600 |
| beneficiary | |

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

" bank -- account -- bank "

$$X \in \begin{bmatrix} 1600 \\ \vdots \\ 800 \\ \vdots \\ 1600 \\ 6200 \end{bmatrix} \in \mathbb{R}^{d_i}$$

$$x_j \in \{1 \cdots [V]\} \qquad [V] = 10,000$$

$d_i$ : length of email $i$

Multivariate Bernoulli event model

Multinomial event model

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

assume: $P(x \mid y) = \prod_{j=1}^{d} P(x_j \mid y)$

$$x_j \in \{1, \cdots [V]\}$$

Parameters

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

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

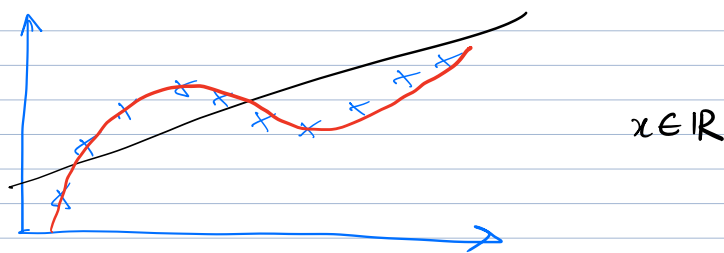Chance that word $j$ is $k^{th}$ word in dictionary
if $y = 0$

MLE $\quad \phi_{k|y=0} = \dfrac{\sum\limits_{i=1}^{n} \left( \mathbb{1}\{y^{(i)}=0\} \sum\limits_{j=1}^{d_i} \mathbb{1}\{x_j^{(i)} = k\} \right)}{\sum\limits_{i=1}^{n} \mathbb{1}\{y^{(i)}=0\} \cdot d_i}$

<span style="color:red">Laplace Smoothing :  $+1$   to numerator</span>
<span style="color:red">$+|V|$   to denominator</span>
<span style="color:red">10,000</span>

Map rare words to "UNK"

- Mortgage
  mørtgↄge
  $\hookleftarrow$ UNK

- spoofed headers

- fetching URL

<span style="color:blue">Kernel Methods</span>



$x \in \mathbb{R}$

Linear models : $\Theta^T x$

$h_\theta(x) = \Theta_3 x^3 + \Theta_2 x^2 + \Theta_1 x + \Theta_0$

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

$h_\theta(x) = [\Theta_0, \Theta_1, \Theta_2, \Theta_3] \begin{bmatrix} 1 \\ x \\ x^2 \\ x^3 \end{bmatrix} = \Theta^T \phi(x)$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{bmatrix}$$

$h_\theta(x) =$ linear in $\theta$, $\phi(x)$

$$\{ (x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots \quad (x^{(n)}, y^{(n)}) \}$$

$$\Downarrow$$

$$\{ (\phi(x^{(1)}), y^{(1)}), (\phi(x^{(2)}), y^{(2)}) \dots \quad (\phi(x^{(n)}), y^{(n)}) \}$$

cubic polynomial for old dataset

$\Longleftrightarrow$ linear on new dataset

## LMS on new dataset

$$\min_\theta \frac{1}{2} \sum_{i=1}^{n} (y^{(i)} - \theta^T \phi(x^{(i)}))^2$$

Gradient Descent:

Loop $\{ \quad \theta := \theta + \alpha \sum_{i=1}^{n} (y^{(i)} - \theta^T \phi(x^{(i)})) \phi(x^{(i)})$

$\}$

$\in \mathbb{R}^p \qquad \in \mathbb{R}^p$

$O(np)$

## Terminology

$\phi : \mathbb{R}^d \longrightarrow \mathbb{R}^p \qquad$ feature map

attributes $\qquad$ features

$x :$ attributes

$\phi(x) :$ "features"

What to do if $p$ is very large?

$d > 1$ for cubic polynomial

$$\phi(x) = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_d \\ x_1^2 \\ \vdots \\ x_i x_j \\ x_d^2 \\ x_1^3 \\ \vdots \\ x_i x_j x_k \\ \vdots \\ x_d^3 \end{bmatrix} \begin{array}{l} \left.\vphantom{\begin{matrix}1\\x\\x\end{matrix}}\right\} d \\ \left.\vphantom{\begin{matrix}x\\x\\x\end{matrix}}\right\} d^2 \\ \left.\vphantom{\begin{matrix}x\\x\\x\\x\end{matrix}}\right\} d^3 \end{array}$$

$$\boxed{\theta^T \phi(x)}$$

$$= \_\_ \cdot 1 + \_\_ x_1 + \_\_ x_2$$
$$+ \quad \_\_ x_i x_j$$
$$+ \quad \_\_ x_i x_j x_k$$

Problem: $\phi(x)$ is high dimensional!

$$p = 1 + d + d^2 + d^3 \qquad O(d^3)$$
$$d = 10^3 \qquad p \sim 10^9$$

$$\theta := \theta + \alpha \sum_{i=1}^{n} \left( y^{(i)} - \theta^T \phi(x^{(i)}) \right) \phi(x^{(i)})$$

Runtime for 1 iteration of GD is $O(np)$

Key observation

If $\theta$ initialized at 0,
then at any time, $\theta$ can be written as
$$\theta = \sum_{i=1}^{n} \beta_i \, \phi(x^{(i)}) \quad \text{for some } \beta_1 \dots \beta_n \in \mathbb{R}$$
$$\in \mathbb{R}^p \hspace{10cm} \in \mathbb{R}^n$$

Proof of observation:
By induction on # iterations

Base Case: iteration 0
$$\theta = 0 = \sum_{i=1}^{n} 0 \cdot \phi(x^{(i)})$$
$$\beta_i$$

Assume at iteration $t$, $\theta = \sum_{i=1}^{n} \beta_i \phi(x^{(i)})$

Next iteration:

$$\theta := \theta + \alpha \sum_{i=1}^{n} \left( y^{(i)} - \theta^T \phi(x^{(i)}) \right) \phi(x^{(i)})$$

$$= \sum_{i=1}^{n} \left( \underbrace{\beta_i + \alpha (\underbrace{y^{(i)} - \theta^T \phi(x^{(i)})}_{\text{scalar}})}_{\text{new } \beta_i} \right) \phi(x^{(i)})$$

$$\underbrace{\qquad\qquad\qquad}_{\text{new } \beta_i}$$

New Algo: represent $\theta \in \mathbb{R}^p$ by $\beta \in \mathbb{R}^n$

$p$ param $\longrightarrow$ $n$ param

$$\beta_i := \beta_i + \alpha \left( y^{(i)} - \theta^T \phi(x^{(i)}) \right)$$

$$= \beta_i + \alpha \left( y^{(i)} - \left( \sum_{j=1}^{n} \beta_j \phi(x^{(j)}) \right)^T \phi(x^{(i)}) \right)$$

$$= \beta_i + \alpha \left( y^{(i)} - \sum_{j=1}^{n} \beta_j \underbrace{\langle \phi(x^{(j)}), \phi(x^{(i)}) \rangle}_{p} \right)$$

$$\underbrace{\qquad}_{n}$$

① $\langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$ can be precomputed

② $\langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle$ can often be computed much faster without explicitly computing $\phi()$