# Naive Bayes

Feature vector $X$     (ex. spam classification)

$$X = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ \\ \end{bmatrix} \begin{matrix} a \\ aardvark \\ \vdots \\ \\ zymurgy \end{matrix}$$

$\uparrow$ top-d words $\downarrow$

$X \in \{0,1\}^d$

$x_i : \mathbb{1}_{\{word\ i\ in\ email\}}$

Want to model $p(x|y)\, p(y)$

$2^d$ possible values of $x$

Assume $x_i$'s are conditionally independent given $y$

$\Rightarrow p(x_1 \ldots x_d | y) = p(x_1|y)\, p(x_2|x_1,y) \cdots p(x_d|x_{d-1} \ldots x_1, y)$

$\underset{assume}{=} p(x_1|y)\, p(x_2|y) \cdots p(x_d|y)$

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

$\left.\begin{array}{c}\\\\\end{array}\right\}$ Naive Bayes assumption

## Parameters

$\phi_{j|y=1} = p(x_j=1|y=1)$    if it is a spam

$\phi_{j|y=0} = p(x_j=1|y=0)$    if it is not spam

$\phi_y = p(y)$         Pr(spam)

## Joint Likelihood

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

MLE: $\phi_y = \sum_{i=1}^{n} \dfrac{\mathbb{1}\{y^{(i)}=1\}}{n}$

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

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

## Prediction

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

with annotations: $\phi_{j|y=1}$ (pointing to $p(x|y=1)$ in numerator), $\phi_{j|y=1}$, $\phi_{j|y=0}$

With new word (i.e. COVID) word $k$

$$P(x_k=1|y=1) = \frac{0}{\#\{y=1\}} = \phi_{k|y=1}$$

$$P(x_k=1|y=0) = \frac{0}{\#\{y=0\}} = \phi_{k|y=0}$$

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

0: never seen

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

with annotations: $0$, $0$

## Laplace Smoothing

Add 1 to count of 1s

Add 1 to count of 0s

## Multivariate Bernoulli Event Model
## (Multinomial Event Model)

$$p(x,y) = p(x|y)p(y)$$

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

↑ no longer binary

$x_j \in \{1 \dots |v|\}$

ex.

$$V = \begin{bmatrix} a & 1 \\ aardvark & 2 \\ \vdots \\ account & 800 \\ \vdots \\ bank & 1600 \\ \vdots \end{bmatrix}$$

$$x = \begin{bmatrix} 1600 \\ \vdots \\ 800 \\ 1600 \\ \vdots \\ 6200 \end{bmatrix} \in \mathbb{N}^d$$

### Parameters

$\phi_y = p(y=1)$

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

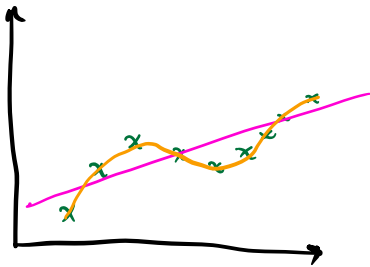↑ Chance that word $j$ is $k^{th}$ word in dictionary if $y=0$

MLE

$$\phi_{k|y=0} = \frac{\sum\limits_{i=1}^{\hat{n}} \mathbb{1}\{y^{(i)}=0\} \sum\limits_{j=1}^{d_i} \mathbb{1}\{x_j^{(i)}=k\}}{\sum\limits_{i=1}^{\hat{n}} \mathbb{1}\{y^{(i)}=0\}d_i}$$

Laplace Smoothing : +1 to numerator
+|V| to denominator

Map rare words to UNK

## Kernel Methods



Linear methods : $\Theta^T x$

What we want : $h_\Theta(x) = \Theta_3 x^3 + \Theta_2 x^2 + \Theta_1 x + \Theta_0$

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

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

$$\Theta = \begin{bmatrix} \Theta_0 \\ \Theta_1 \\ \Theta_2 \\ \Theta_3 \end{bmatrix}$$

$h_\Theta(x)$ : linear in $\Theta$, $\phi(x)$

Dataset

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

$$\Downarrow$$

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

Cubic polynomial on old dataset,
linear on new dataset

## LMS on new dataset

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

Gradient Descent Loop:

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

## Terminology

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

attributes $\uparrow$ $x$

features $\uparrow$ $\phi(x)$

What to do if $p$ is very large?

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

$$\Theta^T \phi(x) = \underline{\quad} \cdot 1 + \underline{\quad} x_1 + \underline{\quad} x_2$$
$$+ \underline{\quad} x_i x_j$$
$$+ \underline{\quad} x_i x_j x_k$$

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

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

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

## Key observation:

If $\Theta$ initialized as $0$, then at any time,

$$\Theta = \sum_{i=1}^{n} \beta_i \phi(x^{(i)}) \quad \text{for } \beta_1 \cdots \beta_n \in \mathbb{R}$$

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