# Lecture 6 : 2/6/17

## Naive Bayes Model

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1}^{N}$$

### Bernoulli Naive Model :

$$p(\vec{x}, y \mid \emptyset, \Theta) = p(x_1, \ldots, x_M, y \mid \emptyset, \Theta)$$

$$= p(y \mid \emptyset) \prod_{m=1}^{M} p(x_m \mid y, \Theta)$$

$$\Theta = \begin{bmatrix} \Theta_{10} & \Theta_{11} \\ \vdots & \vdots \\ \Theta_{M0} & \Theta_{M1} \end{bmatrix}$$

$$= \left[(\emptyset)^y (1-\emptyset)^{(1-y)}\right] \prod_{m=1}^{M} (\Theta_{my})^{x_m} (1-\Theta_{my})^{(1-x_m)}$$

### Naives Bayes Assumption :

$$\boxed{\text{Recall}: \quad p(\vec{x}, y) = p(\vec{x} \mid y)\, p(y)}$$

$$p(\vec{x} \mid y) = \prod_{m=1}^{M} p(x_m \mid y)$$

each $x_q$ is <u>conditionally independent</u> of $x_r$ given $y$ $\forall_{q,r}$

$$\boxed{\begin{array}{l} \text{Def: two r.vs } X, Y \text{ are } \underline{\text{cond. indep.}} \text{ given r.v. } Z \\ \text{written } X \perp Y \mid Z \\ \text{iff } P(X, Y \mid Z) = P(X \mid Z)\, P(Y \mid Z) \end{array}}$$

Q: Why is this "naive"?

A: in real data.

Q: Why is it useful?

A: Count parameters:

Case #1 : w/o NB assumption:

$$p(x_1, \ldots, x_M \mid y) =$$

| $x_1$ | $x_2$ | $\cdots$ | $x_M$ | $y$ | $p(\vec{x}\mid y)$ |
|---|---|---|---|---|---|
| 0 | 0 | $\cdots$ | 0 | 0 | . |
| 0 | 0 | $\cdots$ | 0 | 1 | . |
| 0 | 0 | $\cdots$ | 1 | 0 | . |
| 0 | 0 | $\cdots$ | 1 | 1 | . |

$2^{M+1}$ rows

$2^{M+1} - 2$ params

**Case #2 w/ NBA**

$$p(x_1 = 1 \mid y = 0) = \Theta_{10}$$
$$p(x_1 = 0 \mid y = 0) = 1 - \Theta_{10}$$

4M rows

2M params

$$p(x_1 \mid y) = \begin{array}{|c|c|c|} \hline x_1 & y & p(x_1 \mid y) \\ \hline 0 & 0 & . \\ \hline 0 & 1 & . \\ \hline 1 & 0 & . \\ \hline 1 & 1 & . \\ \hline \end{array}$$

$$p(x_2 \mid y) = \boxed{\phantom{xx}} \quad \cdots \quad p(x_M \mid y) = \boxed{\phantom{xx}}$$
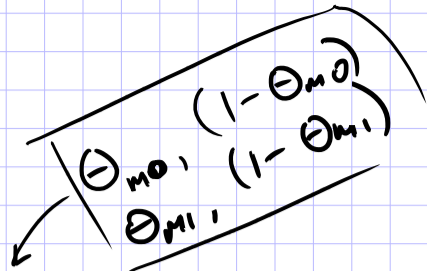
## MLE for Naive Bayes

① Data Likelihood

$$\ell(\phi, \Theta) = \log \prod_{i=1}^{N} p(x^{(i)}, y^{(i)} | \phi, \Theta)$$

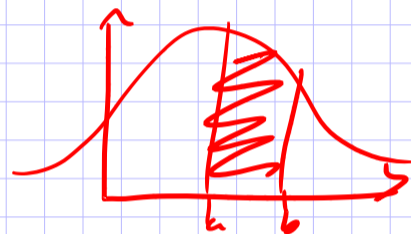$$= \sum_{i=1}^{N} \left[ \log p(y^{(i)}|\phi) + \sum_{m=1}^{M} \log \underline{p(x_m^{(i)} | y^{(i)}, \Theta)} \right]$$

$$\Theta_{m0}, (1-\Theta_{m0})$$
$$\Theta_{m1}, (1-\Theta_{m1})$$

# times $y^{(i)}$ is 1 in D

# times that $x_m^{(i)}=1$ and $y^{(i)}=0$

$$= \boxed{N_{y=1}} \log \phi + N_{y=0} \log(1-\phi)$$

$$+ \sum_{m=1}^{M} \left[ N_{x_m=1, y=1} \log \Theta_{m1} + N_{x_m=0, y=1} \log(1-\Theta_{m1}) \right]$$

$$+ \sum_{m=1}^{M} \left[ N_{x_m=1, y=0} \log \Theta_{m0} + N_{x_m=0, y=0} \log(1-\Theta_{m0}) \right]$$

### MLE for $\phi$ and $\Theta$:

$$(\hat{\phi}, \hat{\Theta}) = \underset{\phi, \Theta}{\arg\max} \; \ell(\phi, \Theta)$$

② Take partials wrt $\phi$:

$$\frac{d\ell(\phi, \Theta)}{d\phi} = \frac{N_{y=1}}{\phi} + \frac{N_{y=0}}{\phi - 1}$$

we already know the MLE!
set equal to zero, solve for $\phi$...

$$\phi^{MLE} = \frac{N_{y=1}}{N_{y=1} + N_{y=0}} = \frac{N_{y=1}}{N}$$

③ Take partials wrt $\Theta_{my}$:
Case where $y=0$, $\Theta_{m0}$

$$\frac{d\ell(\phi, \Theta)}{d\Theta_{m0}} = \frac{N_{x_m=1, y=0}}{\Theta_{m0}} + \frac{N_{x_m=0, y=0}}{\Theta_{m0} - 1}$$

$$\Rightarrow \Theta_{m0}^{MLE} = \frac{N_{x_m=1, y=0}}{N_{x_m=1, y=0} + N_{x_m=0, y=0}} = \frac{N_{x_m=1, y=0}}{N_{y=0}}$$

# MAP Estimation for NB

### The Problem w/ MLE:

Suppose we never observe "Brexit" in Onion article

$\forall i$ where $y^{(i)} = $ Onion, $X_{Brexit}^{(i)} = 0$

Q: What is the MLE $\Theta_{Brexit, Onion} = ?\,0$

$$p(y = Onion \mid \vec{x}^{(new)}) = 0$$

↑ contains Brexit

### Beta Priors:

$$\emptyset \sim Beta(\alpha, \beta) \quad \longleftarrow \text{not in HW}$$

$$\Theta_{m,y} \sim Beta(\alpha, \beta) \quad \forall y \in \{0,1\} \quad \# m \in \{1, \ldots, M\}$$

for $i$ in $1 \ldots N$

$$y^{(i)} \sim Bern(\emptyset)$$
$$x_1^{(i)} \sim Bern(\quad)$$
$$\vdots$$
$$x_m^{(i)} \sim Bern(\quad)$$

$$\ell_{MAP}(\emptyset, \Theta) = \log\left[p(\emptyset, \Theta \mid \alpha, \beta) \, p(D \mid \emptyset, \Theta)\right]$$

$$= \log\left(\left[p(\emptyset \mid \alpha, \beta) \prod_{m=1}^{M} p(\Theta_{m0} \mid \alpha, \beta) p(\Theta_{m1} \mid \alpha, \beta)\right] \left[\prod_{i=1}^{N} p(x^{(i)}, y^{(i)} \mid \emptyset, \Theta)\right]\right)$$

### MAP Estimates:

$$(\hat{\emptyset}, \hat{\Theta}) = \underset{\emptyset, \Theta}{argmax} \; \ell_{MAP}(\emptyset, \Theta)$$

Take partials, set equal to zero, and solve.

$$\Theta_{m0} = \frac{N_{x_m=1, y=0} + (\alpha - 1)}{N_{y=0} + (\alpha-1) + (\beta-1)}$$

$$1 - \Theta_{m0} = \frac{N_{x_m=0, y=0} + (\beta - 1)}{N_{y=0} + (\alpha-1) + (\beta-1)}$$

## Gaussian Naive Bayes

Gaussian NB Model :

$y \sim Bern(\emptyset)$     $= p(y)$
$x_1 \sim Gaussian(\mu_{1y}, \sigma^2_{1y}) = p(x_1|y)$
$\vdots$
$x_M \sim Gaussian(\mu_{My}, \sigma^2_{My}) = p(x_M|y)$

Data: $y \in \{0,1\}$
$\vec{x} \in \mathbb{R}^M$

$$p(\vec{x}, y) = p(y) \prod_{m=1}^{M} p(x_m|y)$$
$\underbrace{\qquad}_{\text{event model}}$

Gaussian NB Parameters

$\emptyset \in \mathbb{R}$

$\mu = \begin{bmatrix} \mu_{10} & \mu_{11} \\ \vdots & \vdots \\ \mu_{M0} & \mu_{M1} \end{bmatrix}$  ← for $x_1$

  ← for $x_M$

$\uparrow$ if $y=0$   $\uparrow$ if $y=1$

$\mu_{my} \in \mathbb{R} \quad \forall y, m$

$\sigma^2 = \begin{bmatrix} \sigma^2_{10} & \sigma^2_{11} \\ \vdots & \vdots \\ \sigma^2_{M0} & \sigma^2_{M1} \end{bmatrix}$

$\sigma^2_{my} > 0 \quad \forall y, m$

# Details:

① Data likelihood:

② Learning (MLE or MAP)

③ Prediction:

$$\hat{y} = h(\vec{x}) = \underset{y \in \{0,1\}}{\arg\max}\ p(y \mid \vec{x}, \phi, \mu, \sigma^2)$$

Everything just works like "Bernoulli" NB!