

# Cs 3780/5780 Logistic Regression

Recall Naïve Bayes Model:

Assumption:  $P(\vec{X} = \vec{x} | Y = y) = \prod_{\alpha=1}^d P(x^{[\alpha]} = x^{[\alpha]} | Y = y)$

$$P(Y = y | X = x) = \frac{\prod_{\alpha=1}^d P(x^{[\alpha]} = x^{[\alpha]} | Y = y) P(Y = y)}{\sum_{c \in \mathcal{Y}} \prod_{\alpha=1}^d P(x^{[\alpha]} = x^{[\alpha]} | Y = c) P(Y = c)}$$

Multinomial NB:

$$P(X^{[\omega]} = x^{[\omega]} | Y = y) \propto \theta_{y, [\omega]}^{x^{[\omega]}}$$

Eg  $\mathcal{Y} = \{\text{spam}, \text{Not spam}\}$   
 $= \{+1, -1\}$

$\theta_{+1}$  = distribution of words in spam emails

Gaussian NB:

Same variance across class for each feature

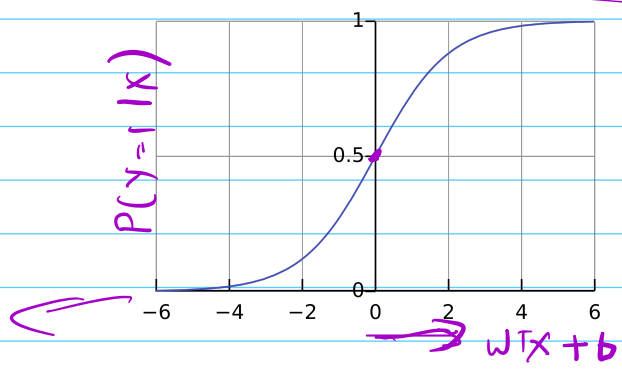
$y = \text{Adult}$

$$P(X^{[\omega]} = x^{[\omega]} | Y = y) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x^{[\omega]} - \mu_y^{[\omega]})^2}{2\sigma^2}}$$

Eg: Take  $\mathcal{Y} = \{+1, -1\}$ , in both the above cases:  
 show that  $P(Y | X = x)$  has the following form:

Logit

$$P(Y = +1 | X = x) = \frac{1}{1 + e^{-(w^T x + b)}}$$



In each of multinomial (and gaussian NB) cases, what are  $w$  and  $b$ ?

Show for Multinomial NB case that:

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

NB Assumption = 
$$\prod_{\alpha=1}^d P(x_{(\alpha)} = s_{(\alpha)} | y=1) P(y=1)$$

$$\rightarrow = \frac{\prod_{\alpha=1}^d (\theta_{+1}(\alpha))^{x_{(\alpha)}} P(y=1)}{\prod_{\alpha=1}^d (\theta_{+1}(\alpha))^{x_{(\alpha)}} P(y=1) + \prod_{\alpha=1}^d (\theta_{-1}(\alpha))^{x_{(\alpha)}} P(y=-1)}$$

$$= \frac{\prod_{\alpha=1}^d e^{w_{+1}(\alpha) \cdot x_{(\alpha)}} \times e^{b_{+}}}{\prod_{\alpha=1}^d e^{w_{+1}(\alpha) \cdot x_{(\alpha)}} e^{b_{+}} + \prod_{\alpha=1}^d e^{w_{-1}(\alpha) \cdot x_{(\alpha)}} e^{b_{-}}}$$

set  $b_{+} = \log(P(y=1))$   
 $w_{+1}(\alpha) = \log(\theta_{+1}(\alpha))$

$$\rightarrow = \frac{e^{w_{+1}^T x + b_{+}}}{e^{w_{+1}^T x + b_{+}} + e^{w_{-1}^T x + b_{-}}}$$

$$= \frac{1}{1 + e^{w_{-1}^T x - w_{+1}^T x + b_{-} - b_{+}}}$$

$$= \frac{1}{1 + e^{-(w^T x + b)}}$$

$$\underline{w} = \begin{bmatrix} w_{+} & -w_{-} \end{bmatrix}$$

$$\underline{b} = \begin{bmatrix} b_{+} & -b_{-} \end{bmatrix}$$

(in terms of  $w_{+}$  and  $w_{-}$ )  
 (in terms of  $b_{+}$  and  $b_{-}$ )

Since  $Y = \{+1, -1\}$   $P(Y=y | \vec{x}=\vec{x}) = \frac{1}{1 + e^{-y(\mathbf{w}^T \mathbf{x} + b)}}$

NB is generative: we model  $P(X, Y)$

Discriminative model we only model  $P(Y|X)$

Discriminative counterpart of Multinomial NB (and Gaussian NB) is Logistic Regression.

Probabilistic model: (absorb bias into last dimension)

$$P(Y=y | \vec{x}=\vec{x}) = \frac{1}{1 + e^{-y\mathbf{w}^T \mathbf{x}}}$$

$$\hat{\mathbf{w}}_{MLE} = \underset{\mathbf{w}}{\operatorname{argmax}} P(D|\mathbf{w}) \quad (\text{Definition of MLE})$$

$$= \underset{\mathbf{w}}{\operatorname{argmax}} P((y_1, \mathbf{x}_1), \dots, (y_n, \mathbf{x}_n) | \mathbf{w}) \quad (\text{Substituting in D.})$$

$$= \underset{\mathbf{w}}{\operatorname{argmax}} \prod_{i=1}^n P(y_i, \mathbf{x}_i | \mathbf{w}) \quad (\text{Data is i.i.d.})$$

$$= \underset{\mathbf{w}}{\operatorname{argmax}} \prod_{i=1}^n P(y_i | \mathbf{x}_i, \mathbf{w}) P(\mathbf{x}_i | \mathbf{w}) \quad (\text{Chain Rule of Statistics})$$

$$= \underset{\mathbf{w}}{\operatorname{argmax}} \prod_{i=1}^n P(y_i | \mathbf{x}_i, \mathbf{w}) P(\mathbf{x}_i) \quad (\mathbf{x}_i \text{ does not depend on } \mathbf{w})$$

$$= \underset{\mathbf{w}}{\operatorname{argmax}} \prod_{i=1}^n P(y_i | \mathbf{x}_i, \mathbf{w}) \quad (P(\mathbf{x}_i) \text{ does not affect } \mathbf{w})$$

$$= \underset{\mathbf{w}}{\operatorname{argmax}} \sum_{i=1}^n \log [P(y_i | \mathbf{x}_i, \mathbf{w})]. \quad (\text{Taking the log})$$

$$= \underset{\mathbf{w}}{\operatorname{argmax}} - \sum_{i=1}^n \log(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i}) \quad (\text{Substituting in } P(y_i | \mathbf{x}_i, \mathbf{w}))$$

$$= \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=1}^n \log(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i}) \quad (\text{We prefer minimization.})$$

Find  $\mathbf{w}$  st.

$$\nabla \sum_{i=1}^n \log(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i}) = 0$$

No closed form, use GD to optimize

loss  $(\mathbf{w}, \mathbf{x}_i, y_i)$   
 " "  
 $\log(1 + e^{-y_i \mathbf{w}^T \mathbf{x}_i})$   
 Logistic loss

Maximum a posterior: Prior  $w \sim N(0, \sigma^2 I)$

$$\begin{aligned}
 \hat{w}_{MAP} &= \underset{w}{\operatorname{argmax}} P(D|w) P(w) \\
 &= \underset{w}{\operatorname{argmax}} P((y_1, \mathbf{x}_1), \dots, (y_n, \mathbf{x}_n) | w) P(w) \\
 &= \underset{w}{\operatorname{argmax}} \left( \prod_{i=1}^n P(y_i | \mathbf{x}_i, w) \right) P(w) \\
 &= \underset{w}{\operatorname{argmax}} \sum_{i=1}^n \log [P(y_i | \mathbf{x}_i, w)] + \log P(w) \\
 &= \underset{w}{\operatorname{argmin}} - \sum_{i=1}^n \log [P(y_i | \mathbf{x}_i, w)] - \log P(w) \\
 &= \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \log [1 + e^{-y_i w^T \mathbf{x}_i}] + \frac{1}{2\sigma^2} w^T w \\
 &= \underset{w}{\operatorname{argmin}} \sum_{i=1}^n \log [1 + e^{-y_i w^T \mathbf{x}_i}] + \lambda w^T w
 \end{aligned}$$

$$\lambda = 1/2\sigma^2$$

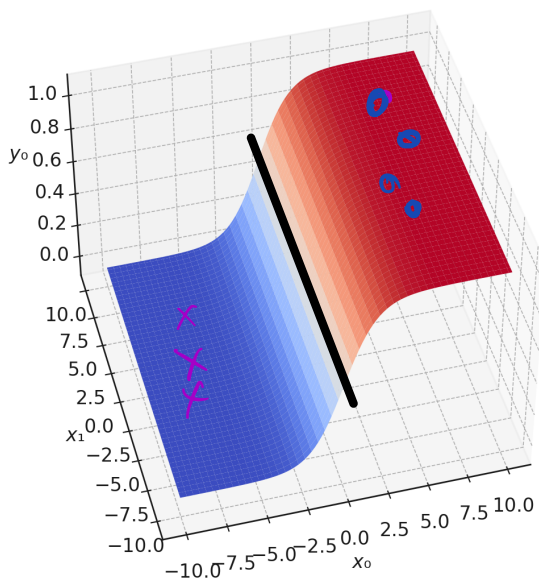
Multiclass version:

$$y = [K]$$

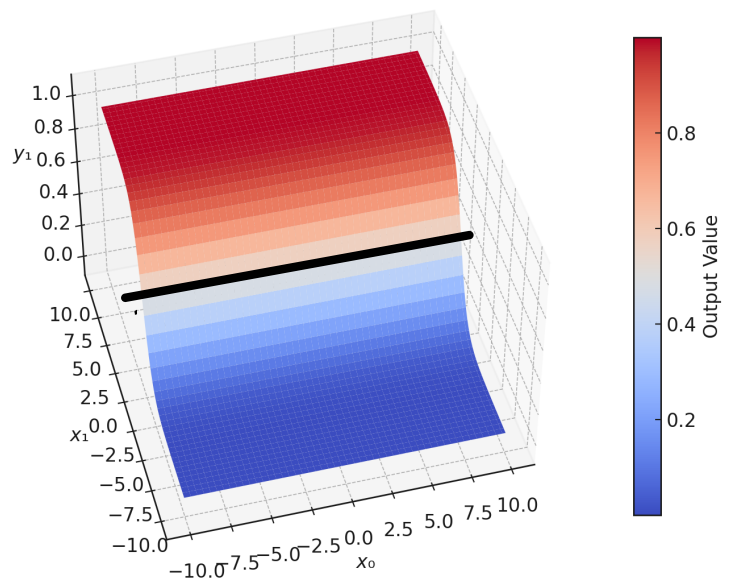
$$w_1, \dots, w_K$$

$$P(y=y | x=x) = \frac{e^{w_y^T x}}{\sum_{k=1}^K e^{w_k^T x}}$$

Sigmoid Output ( $y_0$ )



Sigmoid Output ( $y_1$ )



Multinomial NB:

$$\begin{aligned} P(Y=+1 | \vec{X}=\vec{x}) &= \frac{\prod_{\alpha=1}^d P(X^{(\alpha)} = x^{(\alpha)} | Y=+1) P(Y=+1)}{\sum_{c \in \mathcal{Y}} \prod_{\alpha=1}^d P(X^{(\alpha)} = x^{(\alpha)} | Y=c) P(Y=c)} \\ &= \frac{\prod_{\alpha=1}^d \theta_{\alpha,+1}^{x^{(\alpha)}} P(Y=+1)}{\prod_{\alpha=1}^d \theta_{\alpha,+1}^{x^{(\alpha)}} P(Y=+1) + \prod_{\alpha=1}^d \theta_{\alpha,-1}^{x^{(\alpha)}} P(Y=-1)} \end{aligned}$$

pick  $\vec{w}_y^{(\alpha)} = \log \theta_{\alpha,y}$  and  $b_y = \log P(Y=y)$

$$\begin{aligned} &= \frac{\left( \prod_{\alpha=1}^d e^{\vec{w}_{+1}^{(\alpha)} x^{(\alpha)}} \right) e^{b_{+1}}}{\left( \prod_{\alpha=1}^d e^{\vec{w}_{+1}^{(\alpha)} x^{(\alpha)}} \right) e^{b_{+1}} + \prod_{\alpha=1}^d e^{\vec{w}_{-1}^{(\alpha)} x^{(\alpha)}} e^{b_{-1}}} \end{aligned}$$

$$\begin{aligned} &= \frac{e^{\vec{w}_{+1}^T \vec{x} + b_{+1}}}{e^{\vec{w}_{+1}^T \vec{x} + b_{+1}} + e^{\vec{w}_{-1}^T \vec{x} + b_{-1}}} \end{aligned}$$

$$= \frac{1}{1 + e^{(\vec{w}_{-1} - \vec{w}_{+1})^T \vec{x} + (b_{-1} - b_{+1})}}$$

$$\vec{w} = \vec{w}_{+1} - \vec{w}_{-1} \quad b = b_{+1} - b_{-1}$$

$$\begin{aligned} &= \frac{1}{1 + e^{-(\vec{w}^T \vec{x} + b)}} \end{aligned}$$

$$P(Y = \text{spam})$$

$$P(Y = \text{N-spam})$$

$\theta_{\text{spam}}$  = distribution over words  
when mail is spam

$\theta_{\text{NS}}$  = distribution over words  
when mail is not spam.

$$P(X=x | Y=y) = \prod_{\alpha} P(X[\alpha] = x[\alpha] | Y = \text{spam}) \\ \propto (\theta_{\text{spam}}[\alpha])^{x[\alpha]}$$

$d$  = # words in dictionary.

$\theta_{\text{spam}}$  =  $d$  dim vector  
represents probability of each  
word in spam emails.

$\theta_{\text{NS}}$  =  $d$  dim vector -

$$P(\vec{X} = \vec{x} | Y = \text{spam}) \\ \propto \prod_{\alpha=1}^d (\theta_{\text{spam}}[\alpha])^{x[\alpha]}$$

NB Assumption for multinomial:

Given Email is spam (or not spam)

There is a fixed distribution over words and each word is drawn independently of others.

$x$  is a  $d$  dimensional vector  
 $d =$  size of lexicon

$x[i]$  = # occurrences of  $i$ th word in dictionary.

$$P(x | y = \text{spam}) = \prod_{\alpha=1}^d (\theta_{\text{spam}}[\alpha])^{x[\alpha]}$$

" Multinomial distribution "

$$= \frac{m!}{x[1]! \dots x[d]!} \prod_{\alpha=1}^d (\theta_{\text{spam}}[\alpha])^{x[\alpha]}$$