

Cs 3780/5780

Logistic Regression

Recall Naïve Bayes Model:

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

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

Multinomial NB:

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

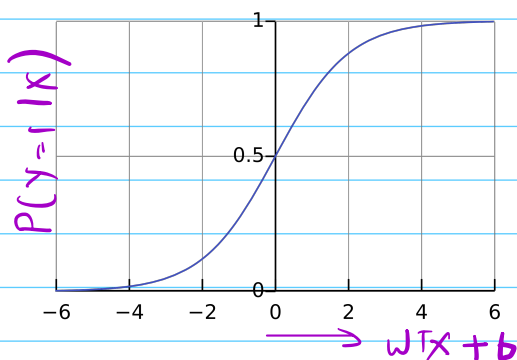
Gaussian NB:

Same variance across class
for each feature

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

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

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



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

Show for Multinomial NB case that:

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

=

=

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

=

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

$w =$

$b =$

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

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

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

NB is generative: we model $P(\mathbf{X}, Y)$

Discriminative model we only model $P(Y|\mathbf{x})$

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

Probabilistic model: (absorb bias into last dimension)

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

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

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

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

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

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

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

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

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

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

Find \vec{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

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

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

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

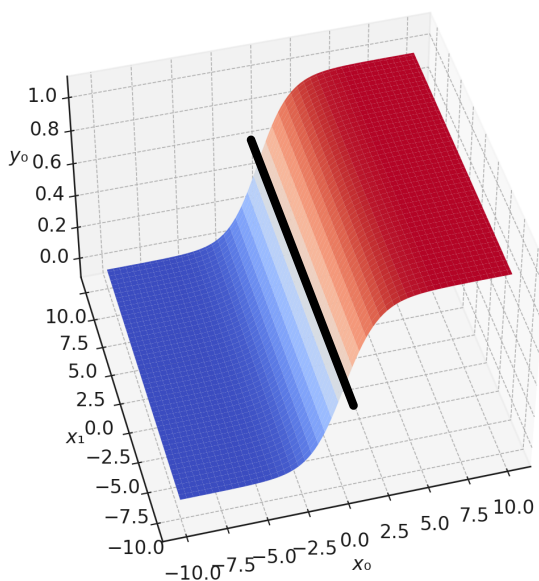
Multiclass version:

$$y = [K]$$

$$w_1, \dots, w_K$$

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

Sigmoid Output (y_0)



Sigmoid Output (y_1)

