

Problem 1

formula for the conditional likelihood

$$1. a: NB: P(Y=1 | X_1, \dots, X_p) = \frac{P(Y=1) \prod_{i=1}^p P(X_i | Y=1)}{\sum_{j=1}^2 P(Y=y_j) \prod_{i=1}^p P(X_i | Y=y_j)}$$

$$LR: P(Y=1 | X_1, \dots, X_p) = \frac{\exp(w_0 + \sum_{i=1}^p w_i X_i)}{1 + \exp(w_0 + \sum_{i=1}^p w_i X_i)}$$

b. classification rule:

NB: for $X^{new} = \langle X_1, \dots, X_n \rangle$ is: $Y^{new} \leftarrow \arg \max_{y_k} P(Y=y_k) \prod_i P(X_i^{new} | Y=y_k)$

LR: $w_0 + \sum_i w_i X_i \geq 0$

c. parameters we have to estimate:

For NB: NB for discrete-valued inputs: have to estimate two sets of parameters:

first: $\theta_{ijk} \equiv P(X_i = x_{ij} | Y=y_k)$ for each input features X_i , each of its possible values x_{ij} , and each of the possible values y_k of Y .

second: we have to estimate the parameter that define the prior probability over Y .

$$\pi_k \equiv P(Y=y_k)$$

NB for continuous inputs: must estimate the mean and standard deviation of each of these Gaussians (for example Gaussian NB)

μ_{ik} and σ_{ik}^2 for each feature X_i and each possible value y_k of Y

For LR: $W \leftarrow \arg \max_W \prod P(Y^i | X^i, W)$

where $W = \langle w_0, w_1, \dots, w_n \rangle$ is the vector of parameters to be estimated.

(d) methods:

For NB: We can use either MLE or MAP.

Choose parameters $W = \langle w_1, \dots, w_n \rangle$ to maximize conditional likelihood of training data, called MLE.

For LR: Since there is no closed form solution to maximizing $l(w)$.

Details are we use the gradient ascent in the following.

2. a) NB is classed a generative classifier, because we can view the distribution $P(X|Y)$ as describing how to generate random instances X conditioned on the target attribute Y .

LR is referred to as a discriminative classifier because we can view the distribution $P(Y|X)$ as directly discriminating the value of the target value Y for any given instance X .

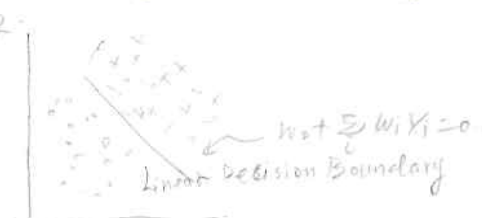
b)

NB



1. decide the position of the contour

LR



$W = \langle w_0, \dots, w_n \rangle$ is the parameters learned.

2. decide the shape of the contour.

(d) of 1.

d. methods to do MLE.

For NB:

$$\hat{\theta}_{MLE} = \arg \max_{\theta} P(D|\theta)$$
$$= \arg \max_{\theta} \prod_{i=1}^n P(Y_i|\theta)$$

Then take derivative and set it to 0, then we get $\hat{\theta}_{MLE}$.

$$\hat{\theta}_{MAP} = \arg \max_{\theta} P(\theta|D)$$
$$= \arg \max_{\theta} P(D|\theta) P(\theta)$$

Then take derivative and set it to 0, then

we get $\hat{\theta}_{MAP}$.

For LR: we solved for LR parameters with MLE.

$$l(w) = \log \prod_{i=1}^n P(Y^{(i)} = y | X^{(i)} = x; w)$$

Since there is no closed-form solution to maximizing $l(w)$, we use gradient ascent.

$$(c) P(Y=1|X) = \frac{P(Y=1)P(X|Y=1)}{P(Y=1)P(X|Y=1) + P(Y=0)P(X|Y=0)}$$

$$= \frac{P(Y=1)}{1 + \frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)}}$$

$$= \frac{1}{1 + \exp\left(\ln \frac{P(Y=0)P(X|Y=0)}{P(Y=1)P(X|Y=1)}\right)}$$

$$= \frac{1}{1 + \exp\left(\ln \frac{P(Y=0)}{P(Y=1)} + \sum_i \ln \frac{P(X_i|Y=0)}{P(X_i|Y=1)}\right)}$$

$$= \frac{1}{1 + \exp\left(\ln \frac{1-\pi}{\pi} + \sum_i \ln \frac{P(X_i|Y=0)}{P(X_i|Y=1)}\right)}$$

Suppose:

$$P(Y=1) = \pi$$

$$P(Y=0) = 1 - \pi$$

write $P(X_i=x|Y=1) = \theta_{i1}^x (1-\theta_{i1})^{1-x}$

$P(X_i=x|Y=0) = \theta_{i0}^x (1-\theta_{i0})^{1-x}$

where θ_{ii} = parameter of feature i ~~under class~~ ^{when $Y=1$}
 θ_{i0} : $Y=0$

Now consider just

$$\sum_i \ln \frac{P(X_i|Y=0)}{P(X_i|Y=1)} = \sum_i \ln \left(\frac{\theta_{i0}}{\theta_{i1}} \right)^{X_i} \cdot \left(\frac{1-\theta_{i0}}{1-\theta_{i1}} \right)^{1-X_i}$$

$$= \sum_i \left[\left(\ln \frac{\theta_{i0}}{\theta_{i1}} - \ln \frac{1-\theta_{i0}}{1-\theta_{i1}} \right) X_i + \ln \frac{1-\theta_{i0}}{1-\theta_{i1}} \right]$$

Let $w_i = \ln \frac{\theta_{i0}}{\theta_{i1}} - \ln \frac{1-\theta_{i0}}{1-\theta_{i1}}$

then $P(Y=1|X) = \frac{1}{1 + \exp\left[\ln \frac{1-\pi}{\pi} + \sum_i (w_i X_i + \ln \frac{1-\theta_{i0}}{1-\theta_{i1}})\right]}$

Let $w_0 = \ln \frac{1-\pi}{\pi} + \sum_i \ln \frac{1-\theta_{i0}}{1-\theta_{i1}}$

then $P(Y=1|X) = \frac{1}{1 + \exp(w_0 + \sum_{i=1}^p w_i X_i)}$

So we can write $P(Y=1|X)$ in a form that matches the Logistic class distribution

(d) From (c), we show that the value of the weights w_i of LR can be ~~be~~ provided in terms of the parameters estimated by the NB classifier. They have the same form. And when we optimize the conditional likelihood, ~~they~~ we get the same classifier.

(2) The Naive Bayes assumption that assumes the attributes X_1, \dots, X_d are all conditionally independent of one another, given Y makes NB less generic than LR.

(f) LR and NB are identical in the limit as the number of training examples approaches infinity, provided the Naive Bayes assumptions hold

3 (a) LR will generally have a lower asymptotic error rate.

(b) $n = \Omega(p)$

(c) $n = O(\log p)$

(d) When the training data is less than $\Omega(p)$ but greater than $O(\log p)$
number of

Problem 2

1. (a) 1. Maximize the conditional log-likelihood is easier.
2. Maximize the log of a function and the function itself is the same

(b) Because log function is a monotonically increasing function

$$2. (a) \ell(w) = \log \left[\prod_{j=1}^n \left(\sum_{c=1}^C I(y_j=c) \frac{e^{w_c^T x_j}}{\sum_{c=1}^C e^{w_c^T x_j}} \right) \right]$$

$$= \sum_{j=1}^n \sum_{c=1}^C \left[I(y_j=c) w_c^T x_j - \log \sum_{c=1}^C e^{w_c^T x_j} \right]$$

$$(b) \nabla_{w_c} \ell(w) = \sum_{j=1}^n \left(\frac{I(y_j=c)}{\sum_{c=1}^C e^{w_c^T x_j}} w_c^T x_j - \nabla_{w_c} \log \sum_{c=1}^C e^{w_c^T x_j} \right)$$

$$= \sum_{j=1}^n \left[I(y_j=c) x_j - \frac{e^{w_c^T x_j} x_j}{\sum_{c=1}^C e^{w_c^T x_j}} \right]$$

$$= \sum_{j=1}^n x_j [I(y_j=c) - P(y_j=c | x_j, w)]$$

(c)

$$\text{When } c=c, \frac{\partial \ell(w)}{\partial w_c} = \sum_{j=1}^n x_j [I(y_j=c) - P(y_j=c | x_j, w)]$$

$$= \sum_{j=1}^n x_j \left[I(y_j=c) - \frac{e^{w_c^T x_j}}{\sum_{c=1}^C e^{w_c^T x_j}} \right]$$

$$= \sum_{j=1}^n x_j (x_j)^T \left[\left(\frac{e^{w_c^T x_j}}{\sum_{c=1}^C e^{w_c^T x_j}} \right)^2 - \left(\frac{e^{w_c^T x_j}}{\sum_{c=1}^C e^{w_c^T x_j}} \right) \right]$$

$$= \sum_{j=1}^n P(y_j=c | x_j, w) [P(y_j=c | x_j, w) - 1] x_j (w_j)^T$$

when $c' \neq c$

$$\begin{aligned}\frac{\partial^2(\omega)}{\partial w_c \partial w_{c'}} &= \left\{ \sum_{j=1}^n x_j \left[I(y_j = c) - P(y_j = c | x_j, w) \right] \right\}' \\&= \sum_{j=1}^n x_j \left[I(y_j = c) - \frac{e^{w_c^T x_j}}{\sum_{c'=1}^C e^{w_{c'}^T x_j}} \right]' \\&= \sum_{j=1}^n x_j (x_j)^T \cdot \frac{e^{w_c^T x_j} \cdot e^{w_{c'}^T x_j}}{\left(\sum_{c'=1}^C e^{w_{c'}^T x_j} \right)^2} \\&= \sum_{j=1}^n (P(y_j = c | x_j, w) \cdot P(y_j = c' | x_j, w) \cdot x_j (x_j)^T)\end{aligned}$$