

Logistic Regression.

We are trying to directly approximate the posterior using linear discriminant functions.

We have,

$$Q_0(x) = \frac{f_0(x) P_0}{f_0(x) P_0 + f_1(x) P_1}$$
$$= \frac{1}{1 + \exp(-\xi)}$$

$$\xi = -\ln \left(\frac{f_1(x) P_1}{f_0(x) P_0} \right) = \log \left(\frac{f_0(x) P_0}{f_1(x) P_1} \right)$$

Now if $\xi = w^T x + w_0$ [a linear function of x]

then $Q_0()$ can be approximated by

$$\hat{y} = h(w^T x + w_0), \quad h(a) = \frac{1}{1 + e^{-a}}$$

The assumption that $E = w^T x + w_0$ is valid if f_0, f_1 are Gaussians.

[also called
Gaussian discriminant
analysis]

Learning of w is done via LMS algorithm

$$J(w) = \sum_{i=1}^n \left[h(w^T x_i + w_0) - y_i \right]^2$$

Another motivation for LR is

Among $P_0(x)$, $P_1(x)$, Bayes' Classifier chooses the maximum.

This is equivalent to looking at the sign of $\log \left[\frac{P_0 f_0(x)}{P_1 f_1(x)} \right]$.

A linear dis fn chooses the classes looking at the sign of $w^T x$. Thus if E is approximated (154) using linear fn of x , it is close to logit.

One can also formulate Logistic regression from an ML perspective.

suppose we model the Posterior using a logit Model

$$\text{i.e., } p(y_i | x_i) = \left(\frac{e^{w^T x_i}}{1 + e^{w^T x_i}} \right)^{y_i} \left(\frac{1}{1 + e^{w^T x_i}} \right)^{1 - y_i}$$

$$\text{Log likelihood} = \sum_{i=1}^n \log p(y_i | x_i)$$

$$= \sum_{i=1}^n y_i w^T x_i - y_i \log(1 + e^{w^T x_i})$$

$$- \log(1 + e^{w^T x_i}) + y_i \log(1 + e^{w^T x_i})$$

$$= \sum_{i=1}^n y_i w^T x_i - \log(1 + e^{w^T x_i})$$

$$\frac{\partial L}{\partial w} = \sum_{i=1}^n y_i x_i^T - \left(\frac{e^{w^T x_i}}{1 + e^{w^T x_i}} \right) x_i^T$$

(155)

To find out the optimal W^* , one can use numerical techniques such as gradient descent.

Multi-class discriminat functions.

We saw LDA for binary problems.

Now suppose we have k classes:

C_1, C_2, \dots, C_k .

Question: How to solve this problem?

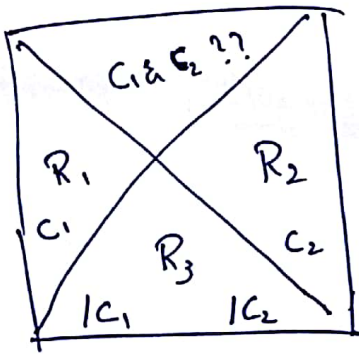
One possible soln: Learn multiple 2-class classifiers & use one-vs-rest approach.

⇐.

Learn k - binary classifiers & perform
a C_i vs not- C_i classifications.

Or learn $\frac{k(k-1)}{2}$ binary classifiers
to perform C_i vs C_j [Computationally expensive]

Both of these approaches however leaves
ambiguous 'gaps' in the feature space.



similarly for C_i vs C_j

A better approach is to learn k disc
fns as follows.

Let us define k fns, $g_s, s=1, \dots, k$

$$g_s(x) = w_s^T x + b_s.$$

Classifier: Assign class C_j to x if

$$g_j(x) \geq g_s(x) \forall s.$$

To learn these, one can convert the class-labels into vectors of k -components as follows.

If $x_i \in C_j$, then y_i would be a k -dim vector with j^{th} component 1 & all other zero.

Learn k fn with vector valued targets.

How to generalize this idea with logistic regression?

Recall that we are trying to approximate the posterior with a composite fn of logits & linear fns

In other words, we were seeking which of $P_1(x)P_1$ & $P_0(x)P_0$ were greater.

For a multi-class Bayes classifier,

the class was chosen as following,

if $Q_i(x_i) > Q_j(x_i) \forall j$, then class

of x_i was taken to be i .

or we need to find the max of

$$Q_i(x) = f_i(x) p_i, i = 1, \dots, K$$

From Bayes rule,

$$Q_j(x) = \frac{f_j(x) p_j}{\sum_s f_s(x) p_s} = \frac{\exp(a_j)}{\sum_s \exp(a_s)}$$

$$a_s = \ln [f_s(x) p_s]$$

Now, the idea is to approximate

$$a_s = W_s^T x + W_{s0}.$$

[As before, the above will be true if

all f_i are Gaussian with same Covariance

Essentially, we are seeking the $\arg\max_a(a_s)$.

Now, Define a fn $g: \mathbb{R}^k \rightarrow \mathbb{R}^k$ with

$$g(a) = [g_1(a) \cdots g_k(a)]^T \quad \forall j=1, \dots, k.$$

$$g_j(a) = \frac{\exp(a_j)}{\sum_s \exp(a_s)}, \quad a = (a_1, \dots, a_k)^T \in \mathbb{R}^k.$$

$g(a)$ is called the Soft-max fn.

[The softened version of the argmax fn]

Now, $a_s = w_s^T x \quad \forall s$.

Let W be a matrix with columns as w_s .

After learning, the final classifier will be

$x \in C_j$ if j^{th} component of $g(W^T x)$ is the highest.
(161).

Ideally, when $x \in C_j$, we need

$$g_j(w^T x) = 1 \quad \& \quad g_i(w^T x) = 0 \\ \forall i \neq j$$

Question: How to compute w .

use gradient descent.

$$L(w) = \sum_{i=1}^n \|g(w^T x_i) - y_i\|^2.$$

$$x_i \in \mathbb{R}^d, \quad y_i = (y_i^1, \dots, y_i^k), \quad y_i^s \in \{0, 1\} \forall i.$$

$$\sum_s y_i^s = 1 \quad \forall i.$$

$$L(w) = \sum_{i=1}^n \left[\sum_{s=1}^k \left(\frac{\exp(w_s^T x_i)}{\sum_j \exp(w_j^T x_i)} - y_i^s \right)^2 \right]$$

↓

different w_s are not decoupled.

(162).

use an approximate version of L :

$$L(w) = \sum_{i=1}^n \left[\sum_{s=1}^k \left(\frac{\exp(w_s^T x_i)}{z_i} - y_i^s \right)^2 \right]$$

$$z_i = \sum_j \exp(w_j^T x_i) \quad \text{computed using}$$

the previous values of w_s .

End of linear families.