

MAST 90138: MULTIVARIATE STATISTICAL TECHNIQUES

See Hastie and Tibshirani, chapters 2, 4.

Classifier $\begin{cases} \text{Linear} \\ \text{logistic} \end{cases}$

7 LINEAR AND QUADRATIC CLASSIFICATION

7.1 INTRODUCTION TO CLASSIFICATION

➡ In the classification problem, we know that individuals from the population come from several, K say, different classes/groups

➡ Example when $K = 2$: healthy patients, unhealthy patients.

➡ We have at our disposal a training sample of individuals from the population for which we observe $(\mathbf{X}_1, G_1), \dots, (\mathbf{X}_n, G_n)$, where

- $G_i = 1, 2, \dots$ or K is the class/group label of the i th individual (an individual belong to exactly one group).

- $\mathbf{X}_i = (X_{i1}, \dots, X_{ip})^T$ is vector of explanatory variables (e.g. age, blood pressure, etc).

- TX file

definition: By comparing a new individual with the training sample, we decide which class it belongs to.

- prediction: by comparing that new individual with the ~~same~~ sample of individuals that I given to study.*
- ☞ A new individual from the same population comes in.
 - ☞ For that individual we only observe the value $\mathbf{x} = (x_1, \dots, x_p)^T$ of \mathbf{X} but we do not know which group the individual comes from (we don't observe his/her G).
 - ☞ Goal: classify the new individual in the correct group, i.e. find the group of that individual. Equivalently, find the class label $G \in \{1, \dots, K\}$ of that new individual.
 - ☞ How can we do that?
 - ☞ **Notation:** throughout chapter we use the notation (\mathbf{X}, G) for a generic individual from the population. The training data $(\mathbf{X}_1, G_1), \dots, (\mathbf{X}_n, G_n)$ have the same distribution as a generic (\mathbf{X}, G) , and so do observed or unobserved data from new individuals.

7.2 MAIN IDEAS OF CLASSIFICATION TECHNIQUES

☞ Main idea: compare the new value \mathbf{x} with the \mathbf{X}_i 's from the training individuals.

☞ If \mathbf{x} looks more similar to \mathbf{X}_i 's from group k than those from the other groups, classify new indiv in group k .

☞ How to decide if \mathbf{x} is more similar to the \mathbf{X}_i 's from group k than from other groups?

☞ For $k = 1, \dots, K$, use training data $(\mathbf{X}_1, G_1), \dots, (\mathbf{X}_n, G_n)$ to estimate

compare

$P(G = k | \mathbf{X} = \mathbf{x})$, 属于 G 里某类别的概率。

the probability that an individual comes from group k given that his/her \mathbf{X} is equal to \mathbf{x} , by

$$\hat{P}(G = k | \mathbf{X} = \mathbf{x}).$$

(We will see later how to compute the estimated probabilities).

- Then, classify new individual in group k if

$$\hat{P}(G = k | \mathbf{X} = \mathbf{x}) > \max_{j=1, \dots, K, j \neq k} \hat{P}(G = j | \mathbf{X} = \mathbf{x}).$$

有给概率最高的 计算每个类别的概率 \rightarrow 最大的那个就是 k 的估计概率。

- Many different methods of classification exist. They essentially correspond to various ways of estimating the above probabilities.

不同分类器差异：概率计算方式不同。

- Two important classes of methods are based on regression estimation techniques and on so-called Bayes methods.

- We start by introducing methods in the case where we have only $K = 2$ classes. We will see later how to extend them to $K > 2$ classes.

7.3 SIMPLE REGRESSION APPROACHES FOR $K = 2$ CLASSES

- Suppose the individuals come from only $K = 2$ groups. Then we classify a new individual in group 1 if (值根据估计值)

$$\hat{P}(G = 1 | \mathbf{X} = \mathbf{x}) > \hat{P}(G = 2 | \mathbf{X} = \mathbf{x}),$$

(连续变量的值两类概率
并无法分类)

where, for $k = 1, 2$, $\hat{P}(G = k | \mathbf{X} = \mathbf{x})$ is an estimator of $P(G = k | \mathbf{X} = \mathbf{x})$ computed from the training sample $(\mathbf{X}_1, G_1), \dots, (\mathbf{X}_n, G_n)$.

- One way to obtain $\hat{P}(G = k | \mathbf{X} = \mathbf{x})$ is through regression.

express these two groups into 2 regression curves

- To see how, for $k = 1, 2$ define a variable from a generic G by taking

对一样本: $Y_k = I\{G = k\} = \begin{cases} 1 & \text{if } G = k \\ 0 & \text{otherwise.} \end{cases}$

G 是类别标签
Indicator variable

indicator variable.
只能有两个值

$k=1 \Rightarrow Y_1 = I\{G=1\} = \begin{cases} 1 & \text{if } G=1 \\ 0 & \text{otherwise} \end{cases}$

$k=2 \Rightarrow Y_2 = I\{G=2\} = \begin{cases} 1 & \text{if } G=2 \\ 0 & \text{otherwise} \end{cases}$

- At the sample level, for $i = 1, \dots, n$, define

$Y_{ik} = I\{G_i = k\}$ (indicator)

Label of i th individual is the label of k th class.

\Rightarrow for class 1: ($Y_1=1, Y_2=0$)
for class 2: ($Y_2=1, Y_1=0$)

Ex: suppose we observe a training sample of size $n = 7$ where the G_i 's are

$$(G_1, \dots, G_7) = (1, 1, 1, 2, 2, 2, 1).$$

Then, for $i = 1, \dots, n$ and $k = 1, 2$, the Y_{ik} 's are given by:

$$(Y_{11}, \dots, Y_{71}) = (1, 1, 1, 0, 0, 0, 1)$$

$$(Y_{12}, \dots, Y_{72}) = (0, 0, 0, 1, 1, 1, 0).$$

Note that

$$Y_1 = I\{G=1\} = \begin{cases} 1 & \text{if } G=1 \\ 0 & \text{otherwise} \end{cases} \quad Y_2 = I\{G=2\} = \begin{cases} 1 & \text{if } G=2 \\ 0 & \text{otherwise} \end{cases}$$

$$P(G = k | \mathbf{X} = \mathbf{x}) = E(I\{G = k\} | \mathbf{X} = \mathbf{x}) = m_k(\mathbf{x}),$$

if we define the regression curve

$$m_k(\mathbf{x}) = E(Y_k | \mathbf{X} = \mathbf{x}).$$

Can use regression estimation techniques to estimate m_k from training data $(\mathbf{X}_1, Y_{1k}), \dots, (\mathbf{X}_n, Y_{nk})$ (Y_{ik} 's observed as computed from G_i 's).

Now every individual must belong to one of the two groups, so that for all $\mathbf{x} \in \mathbb{R}^p$

$$P(G = 1 | \mathbf{X} = \mathbf{x}) + P(G = 2 | \mathbf{X} = \mathbf{x}) = 1,$$

that is

$$m_2(\mathbf{x}) = 1 - m_1(\mathbf{x}).$$

Thus we only need to estimate m_1 (we can deduce m_2 from it).

Discrimination/classification boundary

Once we have estimated m_1 and $m_2 = 1 - m_1$ by estimators \hat{m}_1 and $\hat{m}_2 = 1 - \hat{m}_1$, the classification boundary is the boundary between the region where we classify the data in group 1 and the region where we classify the data in group 2. It is obtained by solving the equation

决策边界

$$\hat{m}_1(\mathbf{x}) - \hat{m}_2(\mathbf{x}) = 0.$$

which, is equivalent to

$$\hat{m}_1(\mathbf{x}) = 1/2.$$

Classify in group 1 if $\hat{m}_1(\mathbf{x}) > \hat{m}_2(\mathbf{x}) \iff \hat{m}_1(\mathbf{x}) > 1/2$, otherwise classify in group 2.

7.3.1 LINEAR REGRESSION CLASSIFIER

☛ When we use a linear regression classifier, we assume

Regression curve:

$$m_1(\mathbf{x}) = E(Y_1 | \mathbf{X} = \mathbf{x}) = P(G = 1 | \mathbf{X} = \mathbf{x}) = \beta_0 + \beta^T \mathbf{x}.$$

The probability to be in a class can be modelled by the linear regression

☛ Estimate m_1 using least squares (LS) estimator of β_0 and β (or PCA or PLS estimators if we need to reduce dimension, see page 260 for how to compute CV in classification) to find

$$\hat{m}_1(\mathbf{x}) = \hat{\beta}_0 + \hat{\beta}^T \mathbf{x}.$$

☛ Deduce $\hat{m}_2 = 1 - \hat{m}_1 = 1 - \hat{\beta}_0 - \hat{\beta}^T \mathbf{x}$.

☛ Classify new \mathbf{x} in group 1 if *based on the assumption that a reasonable model for modelling this probability.*

$$\hat{m}_1(\mathbf{x}) > \hat{m}_2(\mathbf{x}) \iff \hat{\beta}_0 + \hat{\beta}^T \mathbf{x} > 1 - \hat{\beta}_0 - \hat{\beta}^T \mathbf{x} \iff \hat{\beta}_0 + \hat{\beta}^T \mathbf{x} > 1/2;$$

otherwise classify \mathbf{x} in group 2.

☛ NOTE: this method relies on the validity of the linear regression model and is usually only an approximation in real life.

Example: Golub data (ESL, Section 18.4).

Golub TR et al. (1999), Molecular Classification of cancer: class Discovery and Class Prediction by gene expression monitoring, Science 286:531-7.

☞ Two categories of patients, corresponding to two types of leukemia ("ALL" and "AML").

☞ $\mathbf{X} = (X_1, X_2)^T$: gene expressions for two genes known to be connected to leukemia type.

☞ If we fit two linear models $m_1(\mathbf{x}) = \beta_0 + \beta^T \mathbf{x}$ and $m_2(\mathbf{x}) = \beta_{0,2} + \beta_2^T \mathbf{x}$ by LS, we find

$$\hat{m}_1(\mathbf{x}) = 0.9231 + 0.2454x_1 - 0.4800x_2$$

$$\hat{m}_2(\mathbf{x}) = 0.07691 - 0.24542x_1 + 0.47999x_2 \approx 1 - \hat{m}_1(\mathbf{x})$$

as expected (the \approx is due to numerical errors). Illustrates we don't need to fit m_2 , just take $\hat{m}_2 = 1 - \hat{m}_1$

➡ For the Golub data, classification boundary obtained by solving $\hat{m}_1(\mathbf{x}) = 1/2$ is

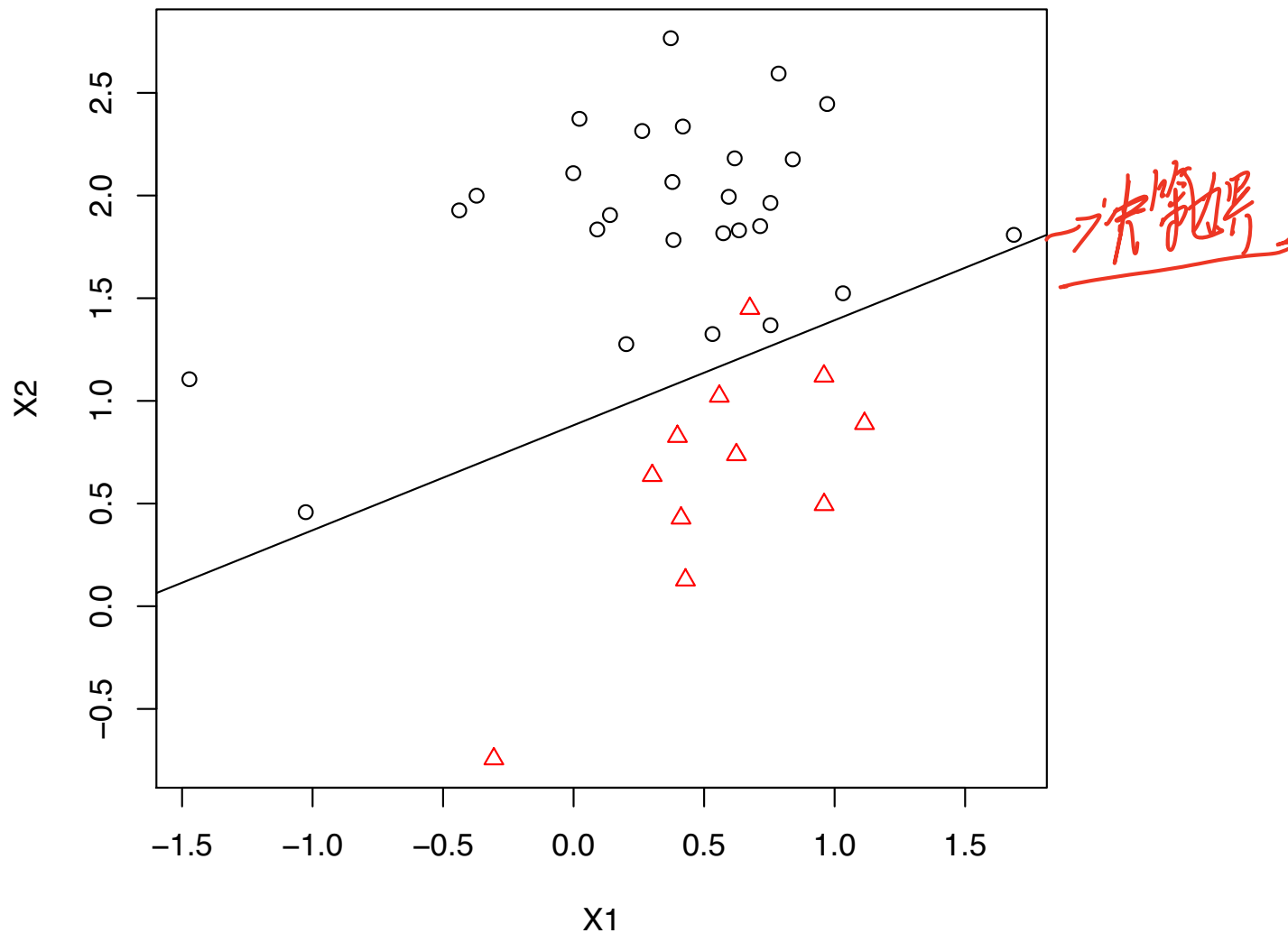
decision boundary $(0.9231 + 0.2454x_1 - 0.4800x_2 = 1/2.)$

➡ This can be expressed by the line

equational line: $x_2 =$
$$x_2 = \frac{0.9231 - 0.5}{0.4800} + \frac{0.2454}{0.4800}x_1 = 0.8815 + 0.5113x_1.$$

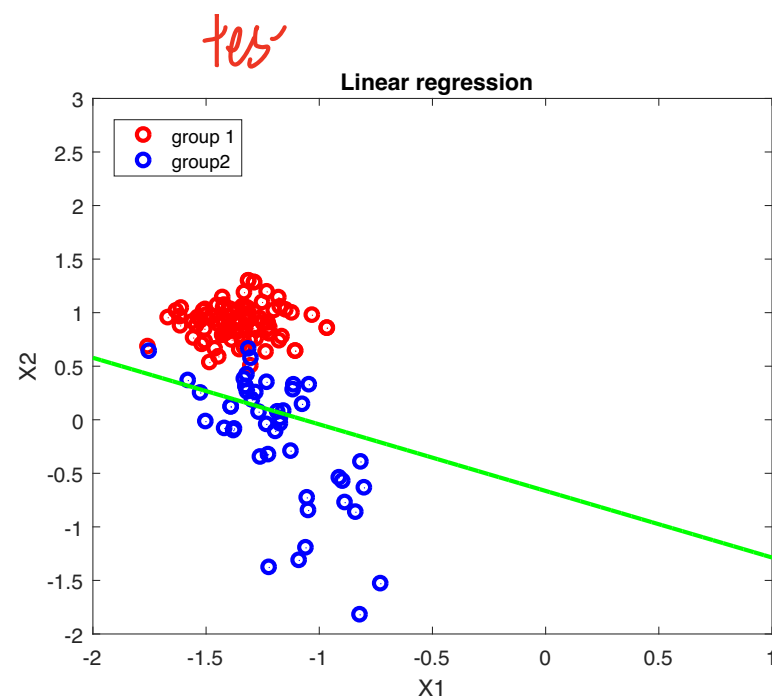
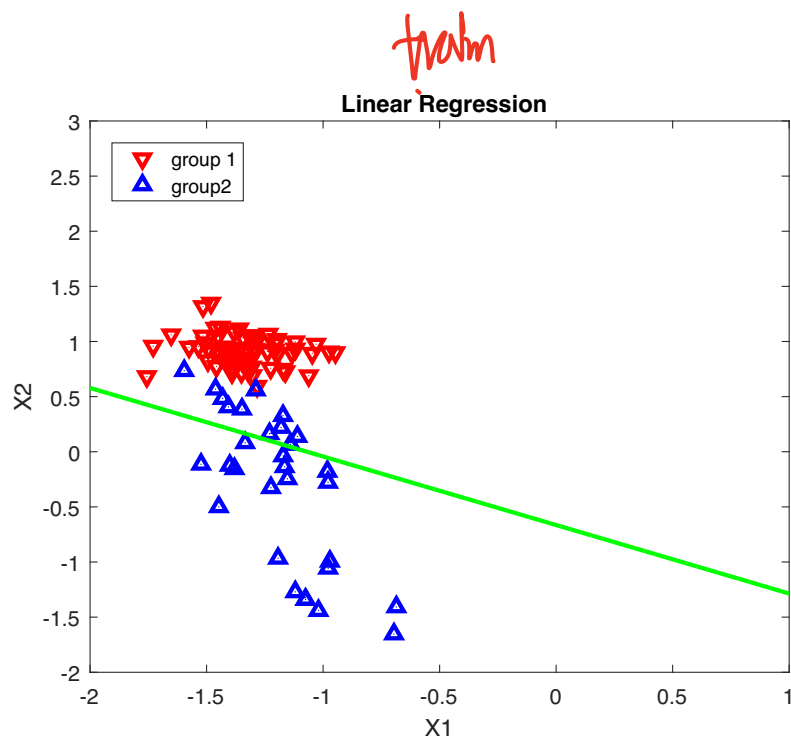
On one side of the line we classify in one class, on the other side we classify in the other class.

Discrimination boundary shown by the line. Training X_i 's from two groups displayed in different colours/symbols. New observations falling below line will be classified in red group; others will be classified in black group.



Other example

- Left: training X_i 's, different groups shown by blue/red.
- Right: set of new data to be classified
- Here it's just an illustration. Unlike real life we know the truth for the new data (blue/red on right fig): can check how classifier performs.
- Green line shows decision boundary $\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 = 1/2$ constructed from training X_i 's.
- Data classified in group 2 if below line, in group 1 if above line (classifier makes mistakes: some blues are above line).



7.3.2 LOGISTIC REGRESSION CLASSIFIER

in linear case: you might estimate probabilities that's weird 已经问过的范围在[-10, 10]而非[0, 1]
线性分类不能保证这个值
A problem with the linear regression classifier is that $\hat{\beta}_0 + \hat{\beta}^T \mathbf{x}$, our estimator of $m_1(\mathbf{x}) = P(G = 1 | \mathbf{X} = \mathbf{x})$, is not guaranteed to be between 0 and 1, whereas it estimates a probability. 始终在[0, 1]之间

One way to ensure to be in $[0, 1]$ is use logistic regression, i.e. assume 线性分类不能保证始终在[0, 1]之间. 用-10, 10

$$m_1(\mathbf{x}) = P(G = 1 | \mathbf{X} = \mathbf{x}) = E(Y_1 | \mathbf{X} = \mathbf{x}) = \frac{\exp(\beta_0 + \beta^T \mathbf{x})}{1 + \exp(\beta_0 + \beta^T \mathbf{x})}.$$

Estimate β_0, β using maximum likelihood (ML) estimators $\hat{\beta}_0, \hat{\beta}$ computed from the training data

$$(\mathbf{X}_1, Y_{11}), \dots, (\mathbf{X}_n, Y_{n1}).$$

(or replace the data by their PCA or PLS components if needed, see page 260 for how to compute CV in classification).

NOTE: here too, the method relies on the logistic model assumption which is usually only an approximation in real life.

As noted at slide 234, $m_2 = 1 - m_1$, i.e. in this model:

$$m_1(\mathbf{x}) = P(G = 1 | \mathbf{X} = \mathbf{x}) = \frac{\exp(\beta_0 + \beta^T \mathbf{x})}{1 + \exp(\beta_0 + \beta^T \mathbf{x})} \quad \text{range } (0,1) \text{ 之间}$$
$$m_2(\mathbf{x}) = P(G = 2 | \mathbf{X} = \mathbf{x}) = 1 - m_1(\mathbf{x}) = \frac{1}{1 + \exp(\beta_0 + \beta^T \mathbf{x})}.$$

Classify new \mathbf{x} in group 1 if

$$\hat{P}(G = 1 | \mathbf{X} = \mathbf{x}) > \hat{P}(G = 2 | \mathbf{X} = \mathbf{x})$$
$$\iff \exp(\hat{\beta}_0 + \hat{\beta}^T \mathbf{x}) > 1 \iff \hat{\beta}_0 + \hat{\beta}^T \mathbf{x} > 0$$

in linear classification: $\hat{\beta}_0 + \hat{\beta}^T \mathbf{x} > \frac{1}{2}$

and otherwise classify \mathbf{x} in group 2.

As before, the classification boundary is obtained by solving

$$\hat{m}_1(\mathbf{x}) = 1/2 \iff \exp(\hat{\beta}_0 + \hat{\beta}^T \mathbf{x}) = 1 \iff \hat{\beta}_0 + \hat{\beta}^T \mathbf{x} = 0.$$

Classify in group 1 if $\hat{m}_1(\mathbf{x}) > \hat{m}_2(\mathbf{x}) \iff \hat{m}_1(\mathbf{x}) > 1/2 \iff \hat{\beta}_0 + \hat{\beta}^T \mathbf{x} > 0$, otherwise classify in group 2.

Example: Golub data

Two categories of patients, corresponding to two types of leukemia ("ALL" and "AML"); $\mathbf{X} = (X_1, X_2)^T$: expressions of two genes.

From the estimated logistic model we have

$$\hat{m}_1(\mathbf{x}) = \frac{\exp(0.9432 + 0.5487x_1 - 10.714x_2)}{1 + \exp(0.9432 + 0.5487x_1 - 10.714x_2)}.$$

In this example the classification boundary is obtained by solving

决策边界 $\underbrace{0.9432 + 0.5487x_1 - 10.714x_2 = 0,}$

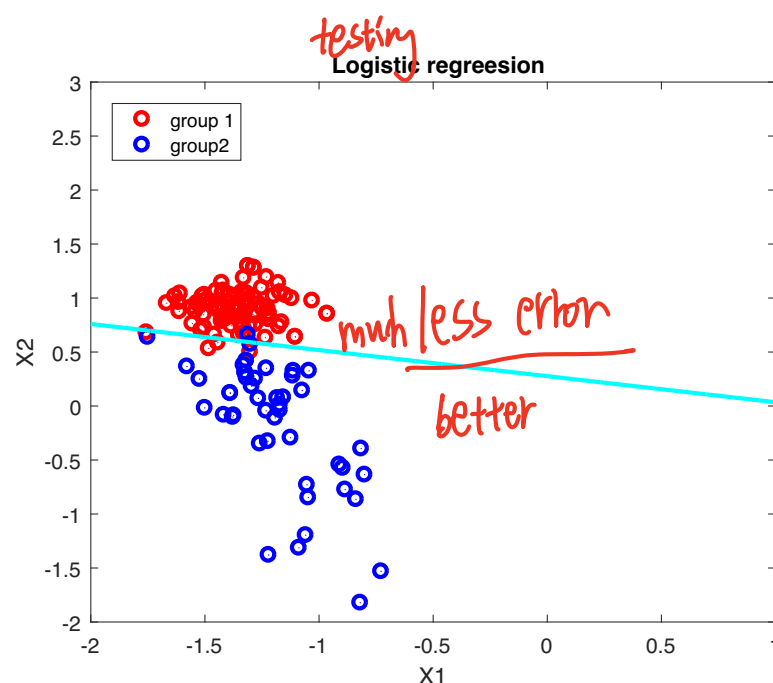
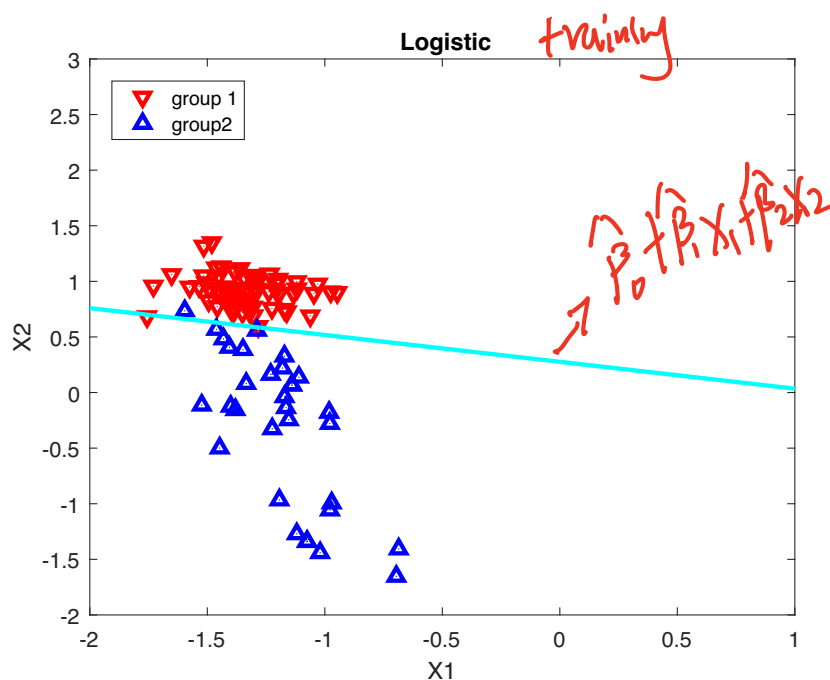
which we can express by the line

在二维图上划的线就要写成这样子

$$x_2 = \frac{0.9432}{10.714} + \frac{0.5487}{10.714}x_1 = 0.08803435 + 0.05121337x_1$$

Other example (same as page 239)

- Left: training X_i 's, different groups shown by blue/red.
- Right: set of new data to be classified
- Here unlike real life we know the truth for the new data (blue/red on right fig): can check how classifier performs.
- Cyan line shows decision boundary $\hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 = 0$ constructed from training X_i 's.
- Data classified in group 2 if below line, in group 1 if above line (classifier makes mistakes: some blues are above line).



Comparison linear and logistic regression

