

Logistic Regression and Discriminant Analysis

Tathagata Basu

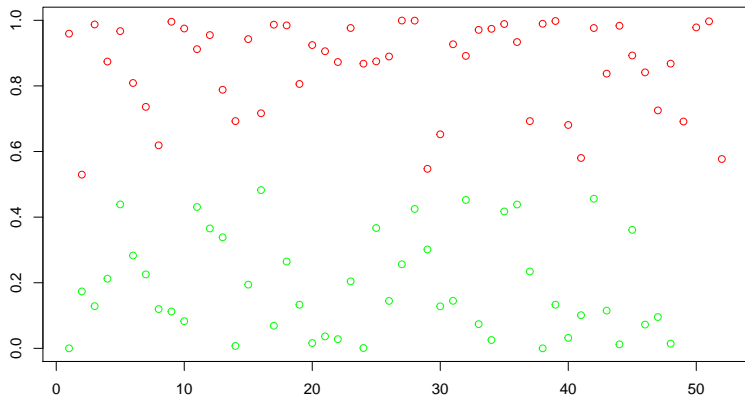
UE de Master 2, AOS1

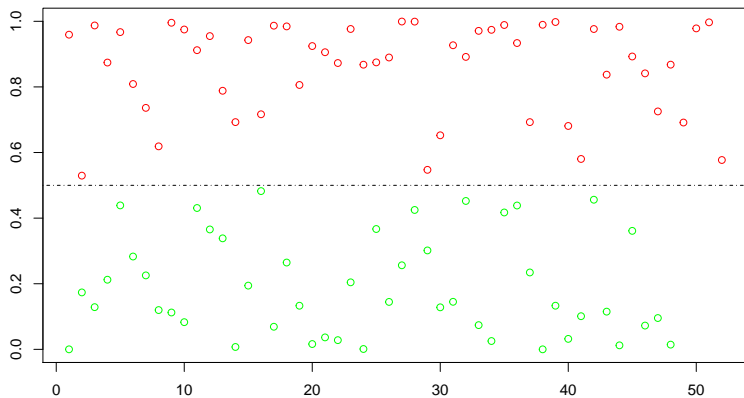
Autumn 2022

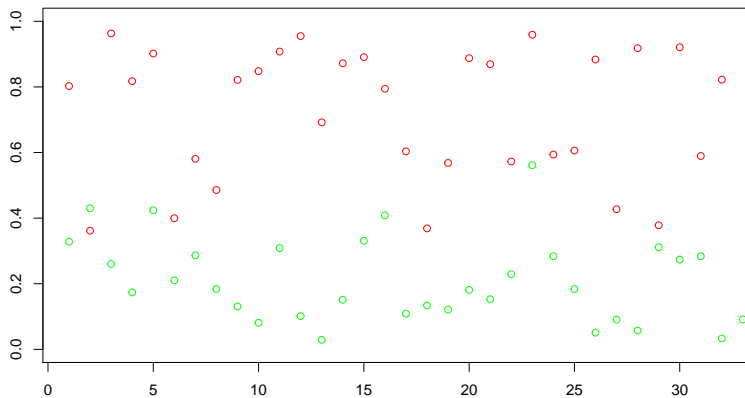
Outline I

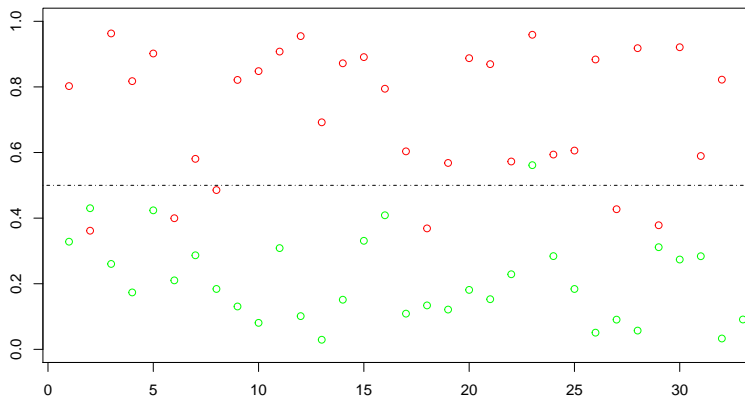
1 Classification

2 Discriminant Analysis









Classification

Let $C = (c_1, \dots, c_n)$ such that $c_i \in \{0, 1\}$. Let $a := (a_1, a_2, \dots, a_p)$ denote attributes.

We define

$$\pi(a) := E(C \mid a) = P(C = 1 \mid a). \quad (1)$$

Therefore, we can say, C is Bernoulli distributed with parameter $\pi(a)$ such that, the likelihood is given by:

$$\mathcal{L}(\pi \mid C, a) = \prod_{i=1}^n \pi(a)^{c_i} (1 - \pi(a))^{1-c_i} \quad (2)$$

Regression model

Let $b := (b_1, b_2, \dots, b_p)^T$ denote regression coefficients. Then $\pi(a)$ can be defined by:

$$\pi(a) = h(a^T b) \quad (3)$$

where h acts as a 'link' function.

For **logistic regression** the link function is given by:

$$h(x) := \frac{\exp(x)}{1 + \exp(x)}. \quad (4)$$

Logistic Regression

Now, replacing our regression model, in the likelihood we get,

$$\mathcal{L}(\pi \mid C, a) = \mathcal{L}(b \mid C, a) = \prod_{i=1}^n \left[h(a^T b) \right]^{c_i} \left[1 - h(a^T b) \right]^{1-c_i} \quad (5)$$

Then the log likelihood is given by:

$$\log(\mathcal{L}(b \mid C, a)) = \sum_{i=1}^n \left(-C_i (a_i^T b) - \log(1 + \exp(a_i^T b)) \right). \quad (6)$$

For, maximum likelihood estimates we find

$$\hat{b} := \arg \max_b \{ \log(\mathcal{L}(b \mid C, a)) \}. \quad (7)$$

Logistic Regression

Now, replacing our regression model, in the likelihood we get,

$$\mathcal{L}(\pi \mid C, a) = \mathcal{L}(b \mid C, a) = \prod_{i=1}^n \left[h(a^T b) \right]^{c_i} \left[1 - h(a^T b) \right]^{1-c_i} \quad (5)$$

Then the log likelihood is given by:

$$\log(\mathcal{L}(b \mid C, a)) = \sum_{i=1}^n \left(-C_i (a_i^T b) - \log(1 + \exp(a_i^T b)) \right). \quad (6)$$

For, maximum likelihood estimates we find

$$\hat{b} := \arg \max_b \{ \log(\mathcal{L}(b \mid C, a)) \}. \quad (7)$$

Regularisation

In high dimensional problems ($p > n$), may contribute to overfitting.

So we can use **Bayesian regularisation**. A natural choice for prior for b is Gaussian distribution, then our log-posterior becomes

$$\log(P(b \mid C, a) \equiv \log(\mathcal{L}(b \mid C, a)) - \lambda \|b\|^2 \quad (8)$$

By maximising our log-posterior, we can obtain **MAP** estimates of b .

Motivation I

- Classical logistic regression models are easy but has limited expressiveness.
- We can construct a probabilistic model with Gaussian assumption instead.

Mixture Model

Let (x_1, \dots, x_n) be a p -dimensional in \mathbb{R}^p and z be corresponding classes such that for $1 \leq i \leq n$, $z_i \in \{1, 2, \dots, M\}$.

We assume that,

$$x_i \mid z_i = k, \mu_k, \Sigma_k \sim \mathcal{N}(\mu_k, \Sigma_k). \quad (9)$$

Then, we have

$$x_i, z_i = k \mid \theta \sim P(z_i = k)P(x_i \mid z_i = k, \mu_k, \Sigma_k) \quad (10)$$

where $\theta \equiv (\mu_1, \dots, \mu_M; \Sigma_1, \dots, \Sigma_M)$

Joint Likelihood

The joint likelihood is given by:

$$L(\theta \mid X, Z) = \prod_i P(X = x_i, z = z_i \mid \theta) \quad (11)$$

$$= \prod_i \prod_k [P(X = x_i, z = k \mid \theta) P(z_i = k)]^{z_{ik}} \quad (12)$$

where, $z_{ik} = \mathbb{I}(z_i = k)$.

So, we have

$$L(\theta \mid X, Z) = \prod_i \prod_k [\pi_k f_k(x_i)]^{z_{ik}} \quad (13)$$

where $P(z = k) = \pi_k$ and $f_k(x_i) = P(X = x_i, z = k \mid \theta)$.

Estimation I

a) π_k

Lagrangian:

$$\mathcal{L}(\theta) = \ln L(\theta \mid \dots) - \lambda \left(\sum_k \pi_k - 1 \right)$$

$$\begin{cases} \frac{\partial \mathcal{L}(\theta)}{\partial \pi_k} = \sum_i \frac{z_{ik}}{\pi_k} - \lambda \\ \frac{\partial \mathcal{L}(\theta)}{\partial \lambda} = 1 - \sum_k \pi_k \end{cases}$$

$$\left. \begin{aligned} \frac{\partial \mathcal{L}(\theta)}{\partial \pi_k} = 0 &\Rightarrow \pi_k = \frac{1}{\lambda} \sum_i z_{ik} \\ \frac{\partial \mathcal{L}(\theta)}{\partial \lambda} = 0 &\Rightarrow \sum_k \pi_k = 1 \end{aligned} \right\} = \frac{1}{\lambda} \sum_i \sum_k z_{ik} = 1 \Leftrightarrow \lambda = \sum_{i,k} z_{ik} = n$$

Estimation II

$$\Rightarrow \mathcal{L}(\theta) = 0 \Leftrightarrow \hat{\pi}_k = \frac{1}{n} \sum_i z_{ik}$$

b) μ_k

$$\frac{\partial \ln L(\theta)}{\partial \mu_k} = \frac{\partial}{\partial \mu_k} \left[\sum_{i,k} z_{ik} \ln f_k(x_i; \theta) \right] \quad (14)$$

$$= C_1 \frac{\partial}{\partial \mu_k} \left[\sum_{i,k} z_{ik} (x_i - \mu_k)^T \Sigma_k^{-1} (x_i - \mu_k) \right] \quad (15)$$

$$= C_2 \sum_i z_{ik} (x_i - \mu_k) \quad (16)$$

Matrix Differentiation

Let $M = [m_{ij}]_{p \times p}$ be a $p \times p$ matrix such that

$$\frac{\partial f(M)}{\partial M} = \begin{bmatrix} \frac{\partial f(M)}{\partial m_{11}} & \cdot & \frac{\partial f(M)}{\partial m_{1p}} \\ \vdots & & \\ \frac{\partial f(M)}{\partial m_{p1}} & \cdots & \frac{\partial f(M)}{\partial m_{pp}} \end{bmatrix}. \quad (17)$$

Then the following holds

- $\frac{\partial |M|}{\partial M} = |M|(M^{-1})^T$
- $\frac{\ln f(M)}{\partial M} = \frac{1}{f(M)} \frac{\partial f(M)}{\partial M}$
- $\frac{\partial \text{tr}(MA)}{\partial M} = A^T$

Estimation III

c) Σ_k

Using the rule of matrix differentiation, we can show that

$$\hat{\Sigma}_k = \frac{\sum_i z_{ik} \hat{B}_{ik}}{\sum_i z_{ik}} \quad (18)$$

where

$$B_{ik} = (x_i - \mu_k)(x_i - \mu_k)^T. \quad (19)$$

$$\hat{B}_{ik} = (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T. \quad (20)$$

Prediction

Reminder

$$P(Z | X) = \frac{P(X | Z)P(Z)}{P(X)} = \frac{P(X | Z)P(Z)}{\sum_z P(X | Z = z)P(Z = z)}$$

$$\Rightarrow P(Z = k | X = x) = \frac{\pi_k f_k(x)}{\sum_k \pi_k f_k(x)} \text{ for any } x$$

\Rightarrow our posterior probability estimates are

$$\hat{P}(Z = k | x) = \frac{\hat{\pi}_k f_k(x | \hat{\mu}_k, \hat{\Sigma}_k)}{\sum_{\ell} \hat{\pi}_{\ell} f_{\ell}(x | \hat{\mu}_{\ell}, \hat{\Sigma}_{\ell})}$$

Bayesian Regularisation

In some cases, we might not have enough data to train our model. Instead, we can take help of Bayesian methods. We can use

- **Multivariate Gaussian** for μ_k so that $\mu_k \sim \mathcal{N}(\mu_{kp}, \Sigma_k/K_{kp})$
- **Inverse Wishart** for Σ_k so that $\Sigma_k \sim IW(\nu_{kp}, \Lambda_{kp})$

Both of these are conjugate priors and give us simple closed form expressions.

Posterior Estimates

Recall $\hat{B}_{ik} = (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T$. Then we have the following posterior estimates

$$\hat{\mu}_k = \frac{\sum_i z_{ik} x_i + K_{kp} \mu_{kp}}{\sum_i z_{ik} + K_{kp}} \quad (21)$$

$$\hat{\Sigma}_k = \frac{\sum_i z_{ik} \hat{B}_{ik} + K_{kp} (\hat{\mu}_{kp} - \mu_{kp})(\hat{\mu}_{kp} - \mu_{kp})^T + \Lambda_{kp}^{-1}}{\sum_i z_{ik} + \nu_{kp} + p + 2} \quad (22)$$