

Linear Discriminant Analysis

1.Theory

From Bayes Rules, the condition probability can be formulated as:

$$p(Y = k|X = x) = \frac{\pi_k f_k(x)}{\sum_{l=1}^K \pi_l f_l(x)}$$

Where $f_k(x)$ is the conditional density function of $X|Y=k$, and $\pi_k = P(Y = k)$ is the prior probability.

The best prediction is picking the one that maximizing the posterior:

$$\arg \max_k \pi_k f_k(x)$$

LDA model $f_k(x)$ as a normal distribution. Suppose we model each class density as multivariate Gaussian $N(u_k, \Sigma_k)$, and assume that the covariance matrices are the same across all k , i.e., $\Sigma_k = \Sigma$

Then, the probability function for class k is:

$$f_k(x) = \frac{1}{(2\pi)^{\frac{p}{2}} \Sigma^{\frac{1}{2}}} \exp\left[-\frac{1}{2}(x - u_k)^T \Sigma^{-1}(x - u_k)\right]$$

The log-likelihood function for the conditional distribution is:

$$\begin{aligned} \log f_k(x) &= -\log((2\pi)^{\frac{p}{2}} \Sigma^{\frac{1}{2}}) - \frac{1}{2}(x - u_k)^T \Sigma^{-1}(x - u_k) \\ &= -\frac{1}{2}(x - u_k)^T \Sigma^{-1}(x - u_k) + \text{constant} \end{aligned}$$

Hence, we just need to select the category that attains the highest posterior density:

$$\begin{aligned} Y_{pred} &= \arg \max_k \log(\pi_k f_k(x)) \\ &= \arg \max_k -\frac{1}{2}(x - u_k)^T \Sigma^{-1}(x - u_k) + \log(\pi_k) \end{aligned}$$

Noticing that quadratic term can be simplified to:

$$\begin{aligned} &-\frac{1}{2}(x - u_k)^T \Sigma^{-1}(x - u_k) \\ &= x^T \Sigma^{-1} u_k - \frac{1}{2} u_k^T \Sigma^{-1} u_k + \text{irrelevant things} \end{aligned}$$

Then the discriminant function is defined as:

$$\delta_k(x) = x^T \Sigma^{-1} u_k - \frac{1}{2} u_k^T \Sigma^{-1} u_k + \log(\pi_k)$$

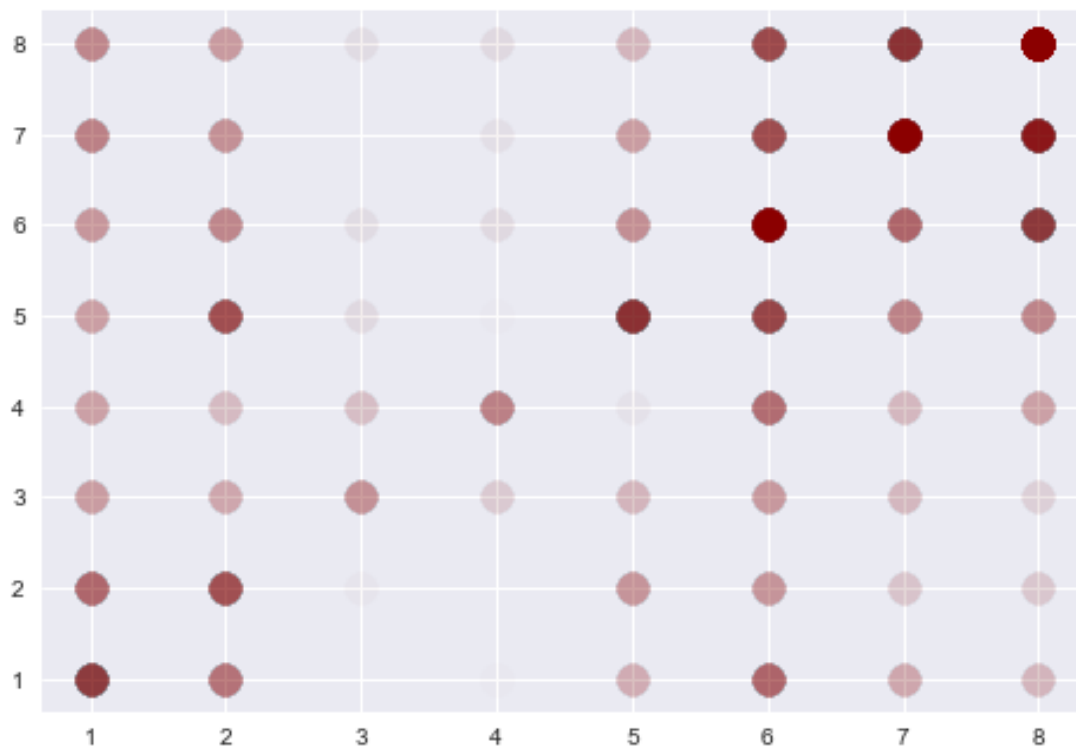
$$= w_k^T x + b_k$$

We can just calculate w_k and b_k for each class k .

2.Result

In LDA.py, I did the preprocessing and wrote the function for LDA. Then I implemented LDA (select all the attributes), it only cost 0.16 seconds. I got the prediction in variable “pred_kaggle” and I got score 0.51 on Kaggle.

Here is accuracy plot (by splitting train data into train and test):



Plot: Actual(X) vs Predict(Y)