

Linear Discriminant Analysis

- LDA can be used for multi-class classification unlike logistic regression.
- LDA is more stable than logistic regression when N is small, and when classes are well separated.

Model

LDA tries to approximate the Bayes's classifier (the classifier that has the lowest error rate) presented below:

$$P(Y=k | X=\mathbf{x}) = \frac{\overset{\text{(prior)}}{P(Y=k)} \cdot \overset{\text{(likelihood)}}{P(X=\mathbf{x} | Y=k)}}{\sum_{l=1}^K P(Y=l) \cdot P(X=\mathbf{x} | Y=l)}$$

where, k represents a class; K is the total number of classes
 \mathbf{x} represents the feature vector

$P(Y=k | X=\mathbf{x})$ is the conditional probability that \mathbf{x} belongs to class k .

$P(Y=k)$ is the marginal probability of the response.

This is also known as the prior.

$P(X=\mathbf{x} | Y=k)$ is the probability density function of class K . It gives the likelihood of \mathbf{x} given the distribution of K .

Fitting the Model

We need two quantities to come up with the Baye's classifier:

1. $P(Y=k)$: To compute the prior, we simply compute the fraction (prior) of data points present in each class.

Sometimes, priors could also be computed based on some information given to us.

2. $P(X=x | Y=k)$:

We assume this density function to be a Gaussian distributed (univariate in case of single predictor; multivariate in case of multiple predictors). Number of distributions = K .

$$P(X=x | Y=k) = \frac{1}{(2\pi)^{p/2} |\Sigma_k|^{1/2}} \exp\left(-\frac{1}{2} (x_k - \mu_k)^T \Sigma_k^{-1} (x_k - \mu_k)\right)$$

i.e. observations from K^{th} class are drawn from $N(\mu_k, \Sigma_k)$
where $x, \mu_k, \Sigma_k \in \mathbb{R}^p$.

Estimating $P(X|Y)$ means estimating Gaussian means estimating μ_k and Σ_k :

μ_k is mean of feature vectors for class k .

Σ_k is covariance matrix of features for class k . It is assumed to be equal for all classes, i.e. $\Sigma_1 = \Sigma_2 = \dots = \Sigma_K = \Sigma$

Now, given the prior $P(Y)$ and density functions (Gaussian with μ_k and Σ_k), we can compute class K for which $P(Y|X)$ is largest.

Computing $\arg\max_K (\text{Baye's classifier}) = \text{computing } \arg\max_K (\text{discriminant function})$

$$\Rightarrow K = \arg\max_K \delta(x)$$

$$\text{where } \underset{\substack{\text{(discriminant} \\ \text{function)}}}{\delta(x)} = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log p(Y=k)$$

The discriminant function is linear in x , hence the name "linear" discriminant analysis.

Assumptions

All assumptions lie in estimating the density function $P(X|Y)$.

1. Each class has Gaussian distribution — huge assumption.
2. All classes share covariance matrix Σ . This is addressed by quadratic discriminant analysis.

Quadratic Discriminant Analysis

- Model is same as LDA.
- Fitting is same as LDA except that QDA does not assume covariance matrices for each class to be the same i.e. $\Sigma_1 \neq \Sigma_2 \neq \dots \neq \Sigma_K$.

The above assumption leads to change in the discriminant function which now becomes quadratic in x :

$$K = \underset{K}{\operatorname{argmax}} \delta(K)$$

where,
$$\delta(K) = -\frac{1}{2} (x - \mu_K)^T \underset{\substack{\uparrow \\ \text{quadratic} \\ \text{function}}}{\Sigma_K^{-1}} (x - \mu_K) + \log p(Y=K)$$

 \uparrow
each class has
a different Σ

- Assumption: The only assumption is that the density functions are Gaussian.

Note: LDA is better than QDA where the data has linear decision boundary, and when there are more predictors (per observation).

QDA is better when decision boundary is non-linear but needs large N to estimate each Σ properly.

LDA - high bias, less variance

QDA - less bias, high variance

- When true decision boundary is linear: LDA, log reg > QDA, KNN
- When X are normally distributed: LDA, QDA > log reg