

LDA는 정규분포를 가정하며, 각 클래스는

평균이 다르다라고 분석은 동일함.

→ 각 클래스의 PDF는 아래와 같음.

$$f_k(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu_k)^2\right)$$

최종적인 LDA의 모델은 아래와 같음.

$$\rightarrow f_k(x) = \frac{\pi_k \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu_k)^2\right)}{\sum_{k=1}^K \pi_k \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu_k)^2\right)}$$

하나, 모든 클래스에 대해서 분포와 상수는 동일하기 때문에 클래스 수에 따라 변하는 식만 남기게 되면, 아래와 같다.

$$\rightarrow \pi_k \exp\left(-\frac{1}{2\sigma^2}(x-\mu_k)^2\right)$$

이 값을 로그변환하면 아래와 같다.

$$\rightarrow \log \pi_k - \frac{x^2}{2\sigma^2} + \frac{\mu_k x}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2}$$

이 중 $\frac{x^2}{2\sigma^2}$ 는 모든 클래스에 대해 동일하므로 최종적으로는 아래만

비교하여 가장 큰 값을 갖는 클래스로
수를 분류하게 된다.

$$\rightarrow \log \pi_k + \frac{\mu_k x}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2}$$

최종 판별함수는 아래와 같다.

$$\rightarrow \delta_k(x) = x \frac{\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2} + \log(\pi_k)$$

이 함수가 선형 함수이기 때문에 선형판별 LDA의 결정경계

부르는 것임.

LDA의 결정경계

k가 2인 경우 두확률이 동일하게 50%인 지점이
있고 이를 표현하면 아래와 같다.

$$\rightarrow \pi_1 = \pi_2 = 0.5$$

→ $\delta_1(x) = \delta_2(x)$ 가 되고, 이를 계산하면

$$x = (\mu_1 + \mu_2) / 2 \text{ 가 된다.}$$

즉, 결정경계 값은 $(\mu_1 + \mu_2) / 2$ 임.

LDA의 파라미터(π) 추정

최대우도 추정법보다 간단하게 접근이 가능함.

$$\rightarrow \hat{\mu}_k = \frac{1}{n_k} \sum_{i: y_i = k} x_i$$

$$\hat{\sigma}^2 = \frac{1}{n-k} \sum_{k=1}^K \sum_{i: y_i = k} (x_i - \hat{\mu}_k)^2$$

$$\hat{\pi}_k = \frac{n_k}{n}$$

다변수 LDA.

$X = (X_1, \dots, X_p)$ 이고, 조건부 가우시안 분포로부터
생성되었다고 가정할 가우시안 분포는 아래와 같다.

$X \sim \mathcal{N}(\mu, \Sigma)$, 평균은 μ , 분포는 Σ 임.

$$\rightarrow \mu = \begin{bmatrix} E[X_1] \\ \vdots \\ E[X_p] \end{bmatrix}$$

$$\Sigma = E[XX^T] = \begin{bmatrix} E[X_1^2] & \dots & E[X_1 X_p] \\ \vdots & & \vdots \\ E[X_p X_1] & \dots & E[X_p^2] \end{bmatrix}$$

다변량 가우시안 분포의 PDF는 아래와 같다.

$$\rightarrow f(x) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right)$$

다변수 LDA의 판별함수

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

마찬가지로 X 에 대한 선형함수의 형태임.

$\delta_k(x) = \delta_l(x)$ 를 통해 x 값을 찾게 되면

승수가 아닌 선형함수가 형태로 드러남.



→ 데이터 포인트는 조건부 가우시안 분포로부터
왔기 때문에 결론만 다르고 공변성이
동일함

→ 결정경계는 선의 형태로 드러남.