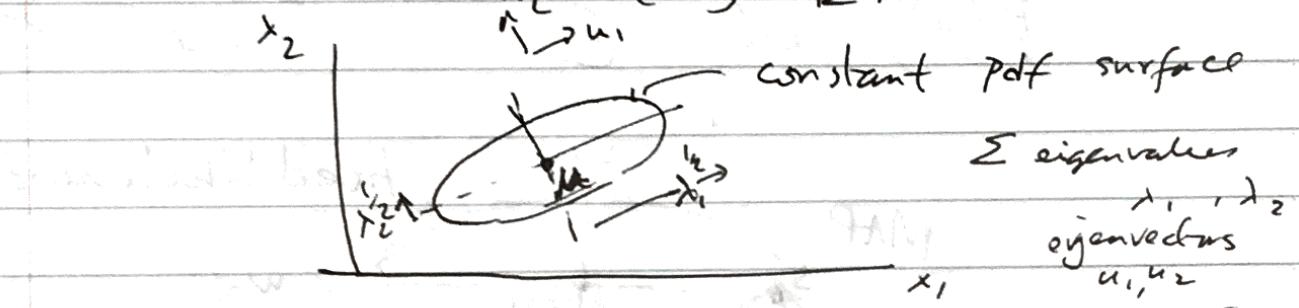


Gaussian Models & Classification

Multivariate Gaussian or multivariate normal (MVN) models.

$$N(x|\mu, \Sigma) \triangleq \frac{1}{(2\pi)^{\frac{D}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)}$$



$$\text{If } \Sigma u_i = \lambda_i u_i \quad \Sigma^{-1} = \sum_i \frac{1}{\lambda_i} u_i u_i^T$$

$$\text{Define } y_i = u_i^T (x - \mu)$$

$$(x - \mu)^T \Sigma^{-1} (x - \mu) = \sum_i \frac{y_i^2}{\lambda_i}$$

Mahalanobis distance.

$$\text{In two-dim } \frac{y_1^2}{\lambda_1} + \frac{y_2^2}{\lambda_2} = \text{constant}$$

→ Ellipses.

MLE for MVN

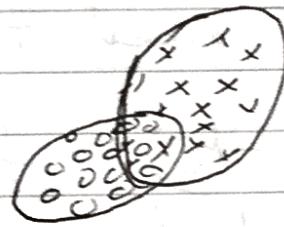
(look up derive in the book!)

$$\begin{aligned} \hat{\mu} &= \frac{1}{N} \sum_{i=1}^N x_i = \bar{x} \\ \hat{\Sigma} &= \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(x_i - \bar{x})^T \\ &= \frac{1}{N} \sum_{i=1}^N x_i x_i^T - \bar{x} \bar{x}^T \end{aligned}$$

Gaussian discriminant analysis

$$P(x|y=c, \theta) = \mathcal{N}(x|\mu_c, \Sigma_c)$$

We will use this model for several supervised (in this lecture) and unsupervised tasks (in later ones).



Imagine we can estimate the parameters (say with labeled data)

Classification

$$\hat{y}(x) = \arg \max_c \left[\log P(y=c|\pi) + \log P(x|\theta_c) \right]$$

$$\text{This is just } \hat{y}(x) = \arg \max_c \left\{ (x - \mu_c)^T \Sigma_c^{-1} (x - \mu_c) + \log |\Sigma_c| \right\}$$

for uniform π

The first term is the Mahalanobis distance from each centroid.

If Σ_c 's are different for each class then classes are decided by thresholding some

Quadratic function \Rightarrow Quadratic

$$P(y=c) = \frac{\pi_c \mu_c^T \Sigma^{-1} (x - \mu_c)}{\sum_{c' \neq c} \pi_{c'} \mu_{c'}^T \Sigma^{-1} (x - \mu_{c'})} \quad \text{Discriminant analysis (QDA)}$$

Life Simplifies if we assume all

$\Sigma_c = \Sigma$, namely, the covariance matrices are tied.

Then we just need to compare

$$\log \frac{\pi_c}{\sum_{c' \neq c} \pi_{c'}} (x - \mu_c)^T \Sigma^{-1} (x - \mu_c) \quad \text{for different classes [} |\Sigma| \text{ part is the same]}$$

$$(x - \mu_c)^T \Sigma^{-1} (x - \mu_c)$$

$$= x^T \Sigma^{-1} x - 2 \mu_c^T \Sigma^{-1} x + \mu_c^T \Sigma^{-1} \mu_c + b$$

\uparrow
independent
of class

$$\text{Call } \gamma_c = \log \frac{\pi_c}{\sum_{c' \neq c} \pi_{c'}} - \frac{1}{2} \mu_c^T \Sigma^{-1} \mu_c$$

$$\beta_c = \Sigma^{-1} \mu_c$$

$$P(y=c|x, \theta) = \frac{e^{\beta_c^T x + \gamma_c}}{\sum_{c'} e^{\beta_{c'}^T x + \gamma_{c'}}} = \sigma(\eta_c)$$

$$\eta = [\beta_1^T x + \gamma_1, \text{[redacted]}; \beta_C^T x + \gamma_C]$$

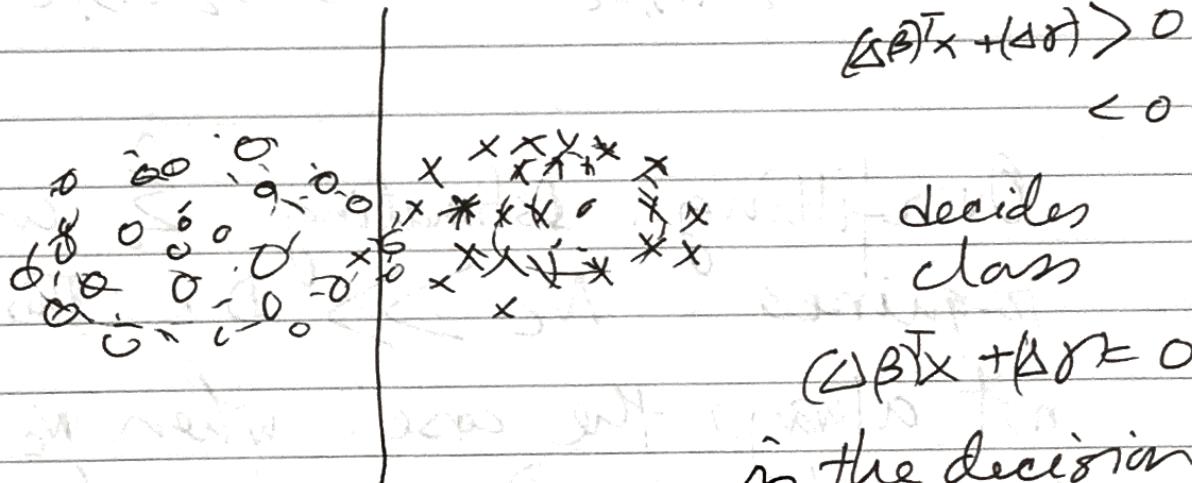
$$S(\eta)_c = \frac{e^{\eta_c}}{\sum_c e^{\eta_c}}$$

is the softmax function

If just have two classes

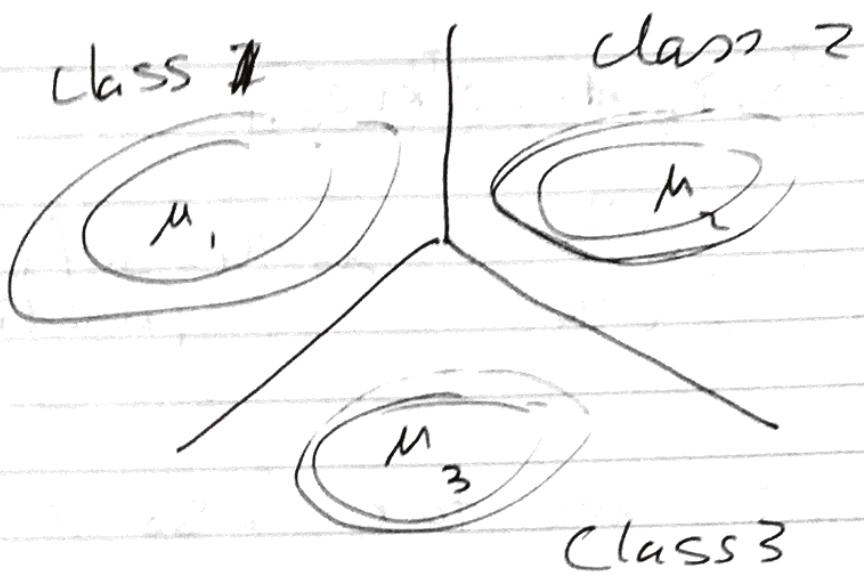
$$c=1, 2$$

$$P(c=1|x, \theta) = \frac{e^{\beta_1^T x + \gamma_1}}{e^{\beta_1^T x + \gamma_1} + e^{\beta_2^T x + \gamma_2}} = \frac{e^{(\beta_1 - \beta_2)^T x + \gamma_1 - \gamma_2}}{1 + e^{(\beta_1 - \beta_2)^T x + \gamma_1 - \gamma_2}}$$



Linear ~~discriminant analysis~~ (LDA)

Even in the multiclass case, we have



For MLE parameter fit with supervised learning for full-fledged gaussian model:

$$\hat{\mu}_c = \frac{1}{N_c} \sum_{i:y_i=c} x_i$$

$$\hat{\Sigma} = \frac{1}{N_c} \sum_{i:y_i=c} (x_i - \hat{\mu}_c)(x_i - \hat{\mu}_c)^T$$

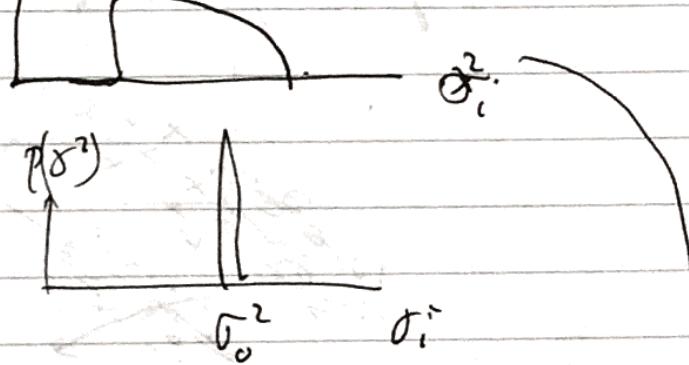
Overfitting! Estimating $\hat{\Sigma}$ well requires $N_c \gg D$. This is not always the case. When N_c is comparable to D we do not estimate $\hat{\Sigma}$ well.

For example, if I give you $x_i \sim N(0, \sigma^2 I_D)$, and you estimate $i=1, \dots, N$, x_i independent

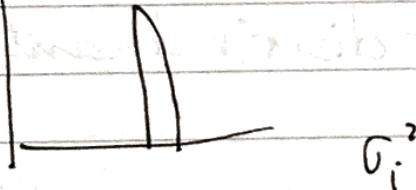
$$\hat{\Sigma} = \frac{1}{N} \sum_i (x_i - \bar{x})(x_i - \bar{x})^T$$

$\hat{\Sigma}$ would have eigenvalues distributed according to the ~~Marčenko-Pastur~~ $P(\sigma_i^2)$ distribution for large N, D with N/D fixed

$$\text{True } \Sigma = \sigma_0^2 I_D$$



Only when $N/D \rightarrow \infty$ do you get



One more we could be fitting noise if $N \sim D$.

- Various solution:
 - a) Parameter tying : $\Sigma_c = \Sigma$ (accept bias to reduce variance)
 - b) Insist on diagonal Σ (like naïve Bayes)
 - c) Integrate over Σ .
 - d) project data to lower dim spaces.

The last idea: projecting to the maximally discriminative subspace, is an important topic in itself.

One approach could be to just find PCA components before classifying and keep a small number of principal components.

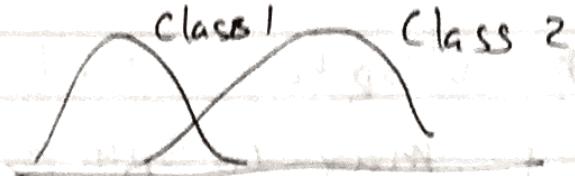


However, that approach does not pay attention to class labels.

Fisher linear discriminant analysis (FLDA)

For two classes, it is just a matter of finding a w ($= \beta_1 - \beta_0$, weight vector) that does the best discrimination.

For x_i , from a particular class we expect $\beta_0 + w^T x_i$ would be normally distributed



We should choose w s.t. Class 1 class 2

$$\boxed{m_c} = \frac{1}{N_c} \sum_{i:y_i=c} x_i$$

$$s_c^2 = \sum_{i:y_i=c} (x_i - m_c)^2$$

↑
Sum sq.
error

Construct a measure

$$J(w) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} = \frac{\text{centrodist}^2}{\text{intraclass variability}}$$

This is proportional to the pooled t-statistic

$$J(w) = \frac{w^T S_B w}{w^T S_W w}$$

$$S_B = (\mu_2 - \mu_1)(\mu_2 - \mu_1)^T$$

$$S_W = \sum_{i:y_i=1} (x_i - \mu_1)(x_i - \mu_1)^T + \sum_{i:y_i=2} (x_i - \mu_2)(x_i - \mu_2)^T$$

$$\boxed{\arg\max_w J(w)}$$

We need the direction,
So, we could choose different
normalizations of w .

$$\text{Max } w^T S_B w \quad \text{s.t. } w^T S_w w = 1$$

$$\max w^T S_B w - \lambda (w^T S_w w - 1)$$

~~$\max_{w \in \mathbb{R}^n} w^T S_B w - \lambda (w^T S_w w - 1)$~~

With the Lagrange multiplier adjusted to satisfy the constraint

$$S_B w = \lambda S_w w$$

Generalized eigenvalue problem

If S_w is invertible:

$S_w^{-1} S_B w = \lambda w$ as a conventional eigenvalue problem

Find the largest eigenvalue.

Because S_B is only rank one for two classes, making further simplification possible

$$S_B w = (\mu_2 - \mu_1)(\mu_2 - \mu_1)^T w = (\mu_2 - \mu_1)(m_2 - m_1)$$

$$\lambda w = S_w(\mu_2 - \mu_1)(m_2 - m_1) \leftarrow \text{scalar} \rightarrow$$

$$w \propto S_w^{-1}(\mu_2 - \mu_1)$$

This method could be generalized to more classes. Since S_B is a rank $C-1$ matrix. We could only find $L \leq C-1$ dimensions. This could be a limitation in some applications.

For completeness: $z_i = w x_i$
 [rows of W are like w^T] $z_i \in \mathbb{R}^L, x_i \in \mathbb{R}^D$

$$m_c = \frac{1}{N_c} \sum_{i:y_i=c} z_i \quad m \in \mathbb{R}^L$$

$$m = \frac{1}{N} \sum_c N_c m_c \quad \text{over all mean}$$

$$J(w) = \frac{\left| \sum_c N_c (m_c - m)(m_c - m)^T \right|}{\left| \sum_{c=1}^C \sum_i (z_i - m_c)(z_i - m_c)^T \right|}$$