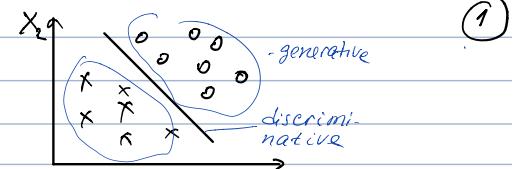


## Generative models



Bayes rule:

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

distribution of  $X$  in a class  $\rightarrow$  Likelihood

↓ class prior  
 ↑ normalizing constant  
 most probable class, given the data does not depend on  $k$

use more prior info  
⇒ often cheaper to compute

The MAP estimator minimizing Bayes risk under 0-1 Loss is

$$\hat{y} = \arg \max_k P\{Y=k|X\} = \arg \max_k f_k(x) \pi_k - \text{ignore the denominator}$$

Specification of  $f_k(x)$  define the classifier: variance-covariance

- $f_k(x)$  multivariate Gaussian, same  $\Sigma$  per class → LDA
- $f_k(x)$  multivariate Gaussian, different  $\Sigma_k$  → QDA
- $f(x) = \text{mixture of Gaussians}$  linear discriminant analysis
- $f_k(x)$  component of the mixture quadratic discriminant analysis
- $f_k(x) = \prod_{p=1}^P f_k(x_p)$  - naive Bayes
  - in each class, orthogonal dimensions
  - typically (but not always)  $f_k(x_p)$  is non-parametric - see hw3

$$\begin{aligned} \hat{y} &= \arg \max_k P\{Y=k|X\} = \arg \max_k [\log f_k(x) + \log \pi_k] \stackrel{\text{MVN}}{=} \arg \max_k \left\{ -\frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) + \log \pi_k \right\} \\ &= \arg \max_k \left\{ -\frac{1}{2} x^T \Sigma^{-1} x + \underbrace{x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k}_{\delta_k(x)} \right\} \end{aligned}$$

$\delta_k(x)$  - linear discriminant function

## LDA in high dimensions

### Nearest centroids (diagonal covariance LDA)

When  $P \gg N$ , estimation of  $\Sigma$  is unstable ⇒ assume indep. predictors

$$X_i \sim \text{MVN}\left(\mu_i, \begin{pmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_p^2 \end{pmatrix}\right), \quad \hat{\mu}_k = \begin{pmatrix} \bar{x}_{1k} \\ \vdots \\ \bar{x}_{pk} \end{pmatrix}, \quad \hat{\Sigma} = \begin{pmatrix} \hat{\sigma}_1^2 & 0 & \dots & 0 \\ 0 & \hat{\sigma}_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & 0 \\ 0 & 0 & \dots & \hat{\sigma}_N^2 \end{pmatrix}$$

over all  $i: Y_i = k$

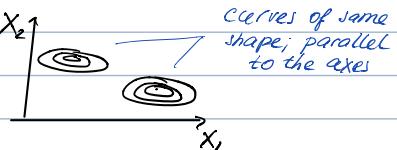
pooled variance across classes  
 $\hat{\sigma}_j^2 = \frac{1}{N-k} \sum_{k=1}^K \sum_{i: Y_i = k} (x_{ij} - \bar{x}_{jk})^2$

$$\hat{y}(x_i) = \arg \max_k \left\{ -\frac{1}{2} \sum_{j=1}^p \frac{(x_{ij} - \bar{x}_{jk})^2}{\hat{\sigma}_j^2} + \log \pi_k \right\}$$

$$\text{Estimation: } \hat{y}(x_i) = \arg \min_k \sum_{j=1}^p \frac{(x_{ij} - \bar{x}_{jk})^2}{\hat{\sigma}_j^2} - z \log \pi_k$$

new obs over  $\{i: Y_i = k\}$

⇒ nearest centroid classifier (after standardization and correction for class prior)



Note:  
connection to SVD

## Connection to Singular Value Decomposition (SVD)

Consider class  $k$ . To estimate  $\hat{\Sigma}_k$  with class-specific variance-covariance,

(1) center each feature

$$\begin{matrix} \text{centered} \\ X' : \end{matrix} \quad \begin{array}{c|cc} & X_1 & \dots & X_p \\ \hline 1 & X_{11} - \bar{x}_1 & \dots & X_{1p} - \bar{x}_p \\ 2 & X_{21} - \bar{x}_1 & \dots & X_{2p} - \bar{x}_p \\ \vdots & \vdots & \ddots & \vdots \\ N_k & X_{N_k 1} - \bar{x}_1 & \dots & X_{N_k p} - \bar{x}_p \\ \hline & \bar{x}_1 & \dots & \bar{x}_p \end{array}$$

$$(2) (N_k - 1)\hat{\Sigma}_k = X'^T X'$$

$\begin{pmatrix} \lambda_1 & & 0 \\ 0 & \lambda_2 & & 0 \\ & & \ddots & 0 \\ & & & \lambda_p \end{pmatrix}$  number of  $\lambda_j \neq 0 = \text{rank of } X_k$

diagonal

Singular Value Decomposition:  $X' = U \Lambda V^T$

$$\begin{matrix} N_k \times p & N_k \times p & p \times p \\ \underbrace{U U^T = I}_{p \times p} & \underbrace{V V^T = I}_{p \times p} & - V \text{ is orthonormal} \end{matrix}$$

$U$  is an orthonormal matrix

Can write  $X'V = U\Lambda$

projection of columns of  $X'$  (i.e., features) onto directions defined by columns of  $V$

Using the SVD decomposition in QDA:

$$\delta_k(X^{\text{new}}) = -\frac{1}{2} \underbrace{(X^{\text{new}} - \mu_k)^T}_{1 \times p} \underbrace{\Lambda^{-1}}_{p \times p} \underbrace{\Sigma_k^{-1}(X^{\text{new}} - \mu_k)}_{p \times 1} + \log \pi_k$$

substitute SVD

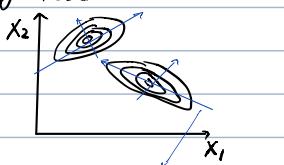
$$\hat{\Sigma}_k = X_k'^T X_k' = V \Lambda U^T U \Lambda V^T = V \Lambda^2 V^T$$

$$\delta_k(X^{\text{new}}) = -\frac{1}{2} \underbrace{[\Lambda^{-1} V^T (X^{\text{new}} - \mu_k)]^T}_{\text{Linear transform}} \underbrace{[\Lambda^{-1} V^T (X^{\text{new}} - \mu_k)]}_{\text{Linear transform}}$$

level curves are spheres on the transformed scale  $\Rightarrow$  classify the new point as the class that has the closest centroid

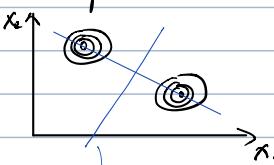
Graphical interpretation:

Original scale:



directions defined by columns of  $V$ , scaled by the corresponding  $\lambda_i$

Transformed scale:



decision boundary on the transformed scale (assuming all  $\pi_k$  equal)

In LDA, same approach; subtract class-specific feature means in (1)

## Regularized LDA Middle ground between LDA and QDA

$$\hat{\Sigma}_k(\alpha) = \alpha \hat{\Sigma}_k + (1-\alpha) \hat{\Sigma}_{\text{pooled}}$$

Choose  $\alpha$  by cross-validation

Middle ground between LDA and diagonal-covariance LDA

$$\hat{\Sigma}(\alpha) = \alpha \hat{\Sigma} + (1-\alpha) I^p$$

$p \times p$  identity matrix

Nearest shrunken centroids:

$$\hat{\mu}_k = \begin{pmatrix} \bar{x}_{11} \\ \vdots \\ \bar{x}_{1P} \end{pmatrix}$$

over all  $i: y_i=k$

centroid of class  $k$

$$\hat{\mu} = \begin{pmatrix} \bar{x}_{11} \\ \vdots \\ \bar{x}_{1P} \end{pmatrix}$$

overall centroid all  $i$

$$\text{Define } d_{jk} = \frac{\hat{\mu}_{jk} - \hat{\mu}_j}{\sqrt{\hat{\sigma}_{kk} - \hat{\sigma}_{jj}} \cdot s_j}$$

standardized distance of class centroid from overall centroid in dimension  $j$

Large  $d_{jk} \rightarrow k$  is a discriminative dimension for class  $k$

$\Rightarrow$

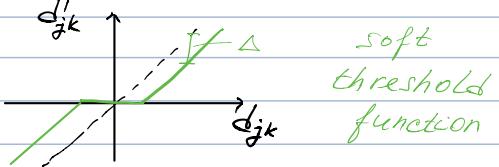
$$\hat{\mu}_{jk} = \hat{\mu}_j + m_k \cdot s_j \cdot d_{jk}$$

when small, dimension  $j$  is not informative for class  $k$

$$\text{Replace with } \hat{\mu}'_{jk} = \text{sign}(d_{jk}) \cdot \left\{ |d_{jk}| - \Delta \right\}_+$$

$$\hat{\mu}'_{jk} = \hat{\mu}_j + m_k \cdot s_j \cdot d_{jk}$$

choose by cross-validation



$\hookrightarrow$  remove dimensions that are non-distinguishable from the overall centroid

Then the discriminant function is

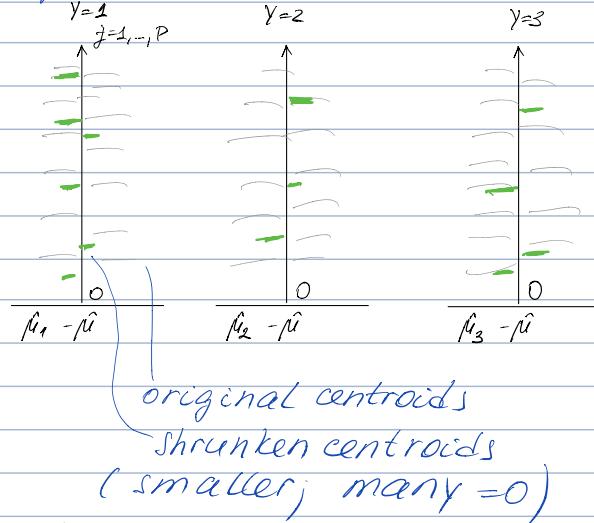
$$\delta_k(x) = \sum_{j=1}^P \frac{x_j - \hat{\mu}'_{jk}}{s_j^2} - 2 \log \pi_k$$

new data point

shrunken centroid

class probabilities:

$$\hat{\pi}_k(x) = \frac{\exp \left\{ -\frac{1}{2} \delta_k(x) \right\}}{\sum_{k=1}^K \exp \left\{ -\frac{1}{2} \delta_k(x) \right\}}$$



Naive Bayes with categorical predictors:

$$P\{Y=k | X\} = \frac{f_k(x) \pi_k}{f(x)} = \frac{\prod_{j=1}^J f_k(x_j) \pi_k}{f(x)} \rightarrow \text{estimate empirically (e.g., from histograms or frequencies)}$$

## Other uses of Bayesian inference

Discriminative models:

Linear regression:  $Y|X \sim N(\beta_0 + \beta_1 X, \sigma^2)$

Logistic regression:  $Y|X \sim \text{Bernoulli} \left( \frac{1}{1 + e^{-(\beta_0 + \beta_1 X)}} \right)$



Generative models

Any  $Y$ :  $P\{Y|X\} = \frac{P(X|Y) P(Y)}{P(X)}$

Categorical  $Y$ :  $P\{Y=k|X\} = \frac{f_k(x) \pi_k}{\sum_{l=1}^L f_l(x) \pi_l}$

apply Bayes rule

Can leverage Bayesian methods in discriminative models. leads to regularization

## Example: Linear regression

$$Y_i = X_i^\top \delta + \varepsilon_i, \quad \varepsilon_i \sim N(0, \sigma^2) \Rightarrow Y_i \sim N(X_i^\top \delta, \sigma^2)$$

### Frequentist estimation

$$\log L(\delta | X, Y) = \text{Log-likelihood} \quad \text{Log } P\{Y | X, \delta\} = \sum_{i=1}^N \left[ \log \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(Y_i - X_i^\top \delta)^2} \right] = \sum_{i=1}^N \left[ -\frac{1}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2}(Y_i - X_i^\top \delta)^2 \right]$$

$$\underline{\text{MLE}}: \quad \hat{\delta}_{\text{MLE}} = \arg \max_{\delta} \left[ -\frac{1}{2\sigma^2} \sum_{i=1}^N (Y_i - X_i^\top \delta)^2 \right] = \arg \min_{\delta} \sum_{i=1}^N (Y_i - X_i^\top \delta)^2 \quad \text{- analytical solution exists}$$

### Bayesian estimation

View some quantities of interest as random variables. Here  $\delta \sim N(0, \frac{1}{2} I)$

$$P(\delta) = \frac{1}{(2\pi)^{p/2}} \cdot e^{-\frac{1}{2} \sum_{j=1}^p \delta_j^2} = \delta^T \delta \quad \begin{matrix} \text{assumed prior distr. of } \delta \\ \text{pre-defined } p \times p \text{ identity matrix} \\ \text{constant} \end{matrix}$$

$$\underline{\text{Log-posterior prob.}} = \frac{\log P\{Y | X, \delta\} \cdot P\{\delta\}}{P\{Y | X\}}$$

continue to view  $X$  as fixed

$$\underline{\text{MAP}}: \quad \hat{\delta}_{\text{MAP}} = \arg \max_{\delta} P\{\delta | X, Y\} = \arg \max_{\delta} \left[ \underbrace{\log P\{Y | X, \delta\}}_{\text{log-likelihood}} + \log P(\delta) \right]$$

$$= \arg \max_{\delta} \left[ \sum_{i=1}^N \left\{ -\frac{1}{2} \log 2\pi\sigma^2 - \frac{1}{2\sigma^2}(Y_i - X_i^\top \delta)^2 \right\} - \frac{p}{2} \log 2\pi - \frac{1}{2} \delta^T \delta \right]$$

$$= \arg \min_{\delta} \frac{1}{\sigma^2} \sum_{i=1}^N (Y_i - X_i^\top \delta)^2 + \frac{1}{2} \delta^T \delta = \hat{\delta}_{\text{ridge}}$$

⇒ Bayesian estimation with priors on parameters  $\delta$  is another motivation for  $L_2$  norm regularization. The prior acts as a regularizer.

## Example: Logistic regression

$$Y_i \sim \text{Bernoulli}(\pi_i), \quad \pi_i = \frac{1}{1 + e^{-X_i^\top \delta}}$$

### Frequentist estimation

$$\log L(\delta | X, Y) = \text{Log-likelihood} \quad \text{Log } P\{Y | X, \delta\} = \sum_{i=1}^N \left[ Y_i \log \pi_i + (1 - Y_i) \log (1 - \pi_i) \right] = \sum_{i=1}^N \left[ Y_i \log \frac{1}{1 + e^{-X_i^\top \delta}} + (1 - Y_i) \log \frac{1}{1 + e^{X_i^\top \delta}} \right]$$

$$\underline{\text{MLE}}: \quad \hat{\delta}_{\text{MLE}} = \arg \max_{\delta} \log L(\delta | X, Y) \rightarrow \text{no analytical solution; gradient descent}$$

### Bayesian estimation

$$P(\delta) = \frac{1}{(2\pi)^{p/2}} \cdot e^{-\frac{1}{2} \sum_{j=1}^p \delta_j^2} = \delta^T \delta \quad \begin{matrix} \text{log-likelihood} \\ \text{extra term involves } \delta \end{matrix}$$

$$\underline{\text{MAP}}: \quad \hat{\delta}_{\text{MAP}} = \arg \max_{\delta} P\{\delta | X, Y\} = \arg \max_{\delta} \left[ \underbrace{\log P\{Y | X, \delta\}}_{\text{use in gradient descent}} + \log P(\delta) \right]$$

$$= \arg \max_{\delta} \left[ \log P\{Y | X, \delta\} + \frac{1}{2} \delta^T \delta \right] = \hat{\delta}_{\text{ridge}}$$

(4)

Similarly, use Laplace(0, 6) prior for L<sub>1</sub> regularization

$$P(\delta) = \frac{1}{2\sigma} e^{-\frac{|\delta|}{\sigma}}$$