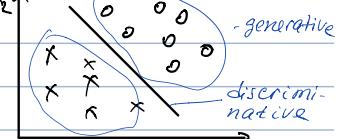


Day 10, 10/9/2018

Generative models

Generative models describe distributions of X in each class



$$\text{Bayes rule: } p(Y|X) = \frac{P(X|Y) P(Y)}{P(X)} \propto \underbrace{P(X|Y) P(Y)}_{\substack{\text{Linear regression} \\ \text{Logistic regression}}} \quad \text{generative models}$$

→ discriminative models

In classification, Y is discrete

$$P\{Y=k|X\} = \frac{f_k(x) \pi_k}{\sum_{l=1}^K f_l(x) \pi_l} \quad \begin{array}{l} \text{distribution of } X \rightarrow \text{likelihood} \\ \text{most probable class, given the data} \end{array}$$

Goal:

find class K maximizing the posterior probability

MAP (maximum a posteriori learning)

$$\text{Data: } \{(x_i, y_i)\}_{y_i \in \{1, \dots, K\}}$$

More formally: minimize the 0-1 loss function

$$J(Y, \hat{Y}(x)) = \begin{cases} 0, & \hat{Y}(x) = Y \text{ - correct prediction} \\ 1, & \hat{Y}(x) \neq Y \text{ - wrong prediction} \end{cases}$$

In Linear reg.:
 $J(Y, \hat{Y}) = \sum_{i=1}^N (Y_i - \hat{Y}_i)^2$

Bayes risk: loss function, averaged over $f_k(x)$ and K :

$$R(\hat{Y}(x)) = E\left\{E\left\{J(Y, \hat{Y}(x))\right\}\right\} \quad \begin{array}{l} \text{over } f(x|Y=k) - \text{all data of a class} \\ \text{over } \pi_k - \text{all classes} \end{array}$$

In Linear regression:
 $E\left\{E\left\{J(Y, \hat{Y}(x))\right\}|\{x_i, y_i\}\right\} \quad \begin{array}{l} \text{over all future data, given training set} \\ \text{over all training sets} \end{array}$

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

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

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

- $f_k(x)$ multivariate Gaussian, same Σ per class → LDA

- $f_k(x)$ multivariate Gaussian, different Σ_k → QDA

$f(x) = \text{mixture of Gaussians}$
 $f_k(x)$ - component of the mixture

linear discriminant analysis

- $f_k(x) = \prod_{p=1}^P f_k(x_p)$ - naïve Bayes

in each class, orthogonal dimensions

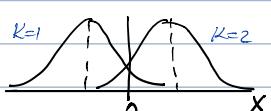
quadratic discriminant analysis

typically (but not always) $f_k(x_p)$ is non-parametric - see hw3

Example: Univariate Gaussian

Population:

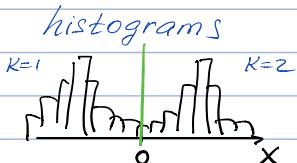
$$\mathcal{N}(\mu_k, \sigma^2), \text{ same}$$



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

$$P\{Y=k|X\} = \frac{\frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{1}{2\sigma^2} (x-\mu_k)^2\right) \cdot \pi_k}{\sum_{l=1}^2 \frac{1}{\sqrt{2\pi} \sigma} \exp\left(-\frac{1}{2\sigma^2} (x-\mu_l)^2\right) \cdot \pi_l} \quad \begin{array}{l} \text{same Variance} \\ \text{does not depend on } K \end{array}$$

Data:



Q: what is the optimal decision boundary?

(2)

Predict class K :

$$\hat{Y}(x) = \arg \max_k P\{Y=k|X\} = \arg \max_k \log P\{Y=k|X\}$$

$$= \arg \max_k \left\{ -\frac{1}{2\sigma^2} (x^2 - 2x\mu_k + \mu_k^2) + \log \pi_k \right\}$$

does not depend on k

$$= \arg \max_k \left\{ \underbrace{\frac{x\mu_k}{\sigma^2} - \frac{\mu_k^2}{2\sigma^2}}_{\delta_k(x), \text{ Linear in } x} + \log \pi_k \right\} = \arg \max_k \delta_k(x)$$

$\delta_k(x)$, Linear in $x \rightarrow$ Linear discriminant function

Assume $\pi_1 = \pi_2 = \frac{1}{2}$

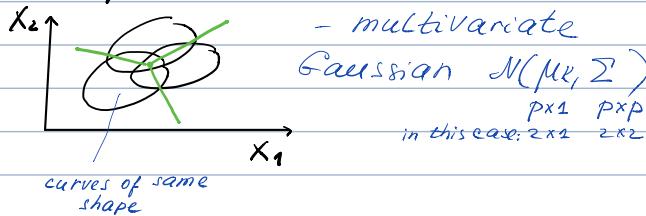
$$\hat{Y}(x) = \arg \max_k \left\{ x\mu_k - \frac{1}{2}\mu_k^2 \right\}. \quad \text{The boundary is } \{x\} \text{ where both classes have same posterior prob.}$$

$$\delta_1(x) = x\mu_1 - \frac{1}{2}\mu_1^2 = x\mu_2 - \frac{1}{2}\mu_2^2 = \delta_2(x)$$

$$x = \frac{\mu_1^2 - \mu_2^2}{2(\mu_1 - \mu_2)} = \frac{\mu_1 + \mu_2}{2}$$

In practice, estimate μ_k , σ and π_k from the dataExample: multivariate Gaussian, $K > 2$ classes

Population:



$$\hat{Y}(x) = \arg \max_k P\{Y=k|X\}$$

$$\stackrel{\text{Bayes rule}}{=} \arg \max_k \frac{f_k(x) \pi_k}{\sum_{l=1}^K f_l(x) \pi_l}$$

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

$$= \arg \max_k [\log f_k(x) + \log \pi_k]$$

Next, substitute multivariate Normal distribution:

$$\hat{Y}(x) = \arg \max_k \left\{ -\log(2\pi)^{\frac{p}{2}} |\Sigma|^{\frac{1}{2}} - \frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) + \log \pi_k \right\}$$

does not depend on K

$$\stackrel{\text{open}}{=} \arg \max_k \left\{ -\frac{1}{2} x^T \Sigma^{-1} x + x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k \right\}$$

Linear expression in $X \rightarrow \delta_k(x)$ - Linear discriminant function

Decision boundary between classes k and l :

$$\{X: P\{Y=k|X\} = P\{Y=l|X\}\} = \{X: \delta_k(x) = \delta_l(x)\} = \{X: \delta_k(x) - \delta_l(x) = 0\}$$

$$\stackrel{\text{MVN}}{=} \{X: x^T \Sigma^{-1} (\mu_k - \mu_l) - \frac{1}{2} (\mu_k - \mu_l)^T \Sigma^{-1} (\mu_k - \mu_l) + \log \frac{\pi_k}{\pi_l} = 0\}$$

For example, $K=2$:

$$\log \frac{\pi_1}{\pi_2} - \frac{1}{2} (\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 - \mu_2) = a_0$$

new notation

$$\underbrace{\Sigma^{-1} (\mu_1 - \mu_2)}_{\substack{\text{new notation} \\ \text{PXP}}} = (a_1 \ a_2 \dots \ a_p)^T \Rightarrow \text{classify } \hat{Y}=1 \text{ if } a_0 + \sum_{j=1}^p a_j x_j > 0$$

Linear combination of all predictors (i.e., no feature selection!)

$$\text{Report posterior probability } P\{y=k | X\} = \frac{f_k(x)\pi_k}{\sum_{i=1}^k f_i(x)\pi_i} = \frac{\exp\left\{-\frac{1}{2}\delta_k(x)\right\}}{\sum_{i=1}^k \exp\left\{-\frac{1}{2}\delta_i(x)\right\}}$$

Non-linear decision boundaries: create higher order features $\tilde{x}_1^2, \tilde{x}_2^2, \tilde{x}_1\tilde{x}_2$ etc

Estimation:

pooled variance-covariance across all classes

$$\hat{\pi}_k = \frac{N_k}{N} \quad \begin{matrix} \# \text{obs} \\ \text{in class } k \\ \text{total} \end{matrix}$$

$$\hat{\mu}_k = \begin{pmatrix} \bar{x}_{1k} \\ \bar{x}_{2k} \\ \vdots \\ \bar{x}_{pk} \end{pmatrix} \quad \begin{matrix} \text{average over all obs.} \\ \text{of class } k \end{matrix}$$

$$\hat{\Sigma} = \frac{1}{N-K} \sum_{k=1}^K \sum_{\{i: y_i=k\}} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T$$

Summary of LDA:

- decision boundary is a linear combination of all features
- no feature selection
- multivariate Gaussian
- Large # of parameters (esp. when p is large and N is small)
- without assuming common Σ across classes \rightarrow quadratic decision boundary (QDA)
- non-linear decision boundaries arise from (1) x^2 etc, or (2) QDA
- subject to bias-variance trade-off

LDA in high dimensions

Nearest centroids (diagonal covariance LDA)

When $p \gg N$, estimation of variance-covariance is unstable

Simplification: assume that the predictors are independent

$$\hat{\mu}_k = \begin{pmatrix} \bar{x}_{1k} \\ \vdots \\ \bar{x}_{pk} \end{pmatrix}, \quad \hat{\Sigma} = \begin{pmatrix} s_1^2 & 0 & \dots & 0 \\ 0 & s_2^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & s_p^2 \end{pmatrix} \quad \begin{matrix} \text{pooled variance across classes} \\ s_j^2 = \frac{1}{N-K} \sum_{k=1}^K \sum_{\{i: y_i=k\}} (x_{ij} - \bar{x}_{jk})^2 \end{matrix}$$

The discriminant function becomes $\delta_k(x_i) = \sum_{j=1}^p \frac{(x_{ij} - \bar{x}_{jk})^2}{s_j^2} - 2 \log \pi_k$

\Rightarrow nearest centroid classifier

(after standardization and correcting for class prior)

$$\hat{Y}(x_i) = \arg \max_k \delta_k(x_i)$$

Regularized LDA

Middle ground between LDA and QDA

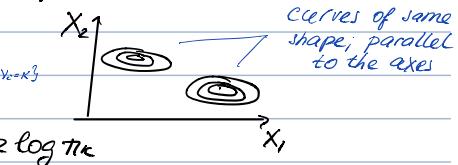
$$\hat{\Sigma}_k(\alpha) = \alpha \hat{\Sigma}_k + (1-\alpha) \hat{\Sigma}_{\text{pooled}}$$

Middle ground between LDA and diagonal-covariance LDA

$$\hat{\Sigma}(\alpha) = \alpha \hat{\Sigma} + (1-\alpha) \alpha^2 I$$

$p \times p$ identity matrix

Choose α by cross-validation



Nearest shrunken centroids:

$$\hat{\mu}_k = \begin{pmatrix} \bar{x}_{1,k} \\ \vdots \\ \bar{x}_{p,k} \end{pmatrix}$$

overall
Centroid
of class k

$$\hat{\mu} = \begin{pmatrix} \bar{x}_{1,\cdot} \\ \vdots \\ \bar{x}_{p,\cdot} \end{pmatrix}$$

overall
Centroid
all c

$$\text{Define } d_{jk} = \frac{\hat{\mu}_{jk} - \hat{\mu}_j}{\sqrt{m_k - m_j} \cdot s_j}$$

standardized distance
of class centroid from
overall centroid in dimension j

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

⇒

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

when small, dimension j
is not informative for
class k

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

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

choose by
cross-validation

Remove dimensions that are non-distinguishable from the overall centroid

Then the discriminant function is

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

new data
point

shrunken
centroids

class probabilities:

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

