

STAT 3690 Lecture 31

zhiyanggeezhou.github.io

Zhiyang Zhou (zhiyang.zhou@umanitoba.ca)

Apr 13, 2022

Classification

- Predictive task in which the response takes values across K discrete categories (i.e., not continuous)
 - For one subject, to predict its class label Y when its features \mathbf{X} is observed
 - Binary classification: $K = 2$
 - Having training data with known class labels
 - E.g.
 - * Given scanned handwritten digits: 28×28 grid of pixels each reflecting the value of grey scale; see Lec 23. From vectorized pictures determine what digit was written.
 - * Predicting the region of Italy in which a brand of olive oil was made, based on its chemical composition; see Lec 29.

-
- Bayes classifier
 - Classify according to posterior $\Pr(Y = k \mid \mathbf{X} = \mathbf{x}) = f_k(\mathbf{x})\pi_k / \sum_{\ell=1}^K f_{\ell}(\mathbf{x})\pi_{\ell}$, $k = 1, \dots, K$
 - * $f_k(\mathbf{x})$: the probability density/mass function of \mathbf{X} conditioning on Class k
 - * $\pi_k = \Pr(Y = k)$: prior probability of Class k
 - Bayes classifier

$$h(\mathbf{x}) = \arg \max_{k=1, \dots, K} \Pr(Y = k \mid \mathbf{X} = \mathbf{x}) = \arg \max_{k=1, \dots, K} f_k(\mathbf{x})\pi_k$$

Linear discriminant analysis (LDA)

- Assuming $f_k(\mathbf{x}) = \text{density of } MVN_p(\boldsymbol{\mu}_k, \boldsymbol{\Sigma})$
- LDA classifier

$$h(\mathbf{x}) = \arg \max_{k=1, \dots, K} \delta_k(\mathbf{x})$$

- Discriminant functions $\delta_k(\mathbf{x}) = \mathbf{x}^{\top} \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k - \frac{1}{2} \boldsymbol{\mu}_k^{\top} \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_k + \ln \pi_k$
 - * Linear functions with respect to \mathbf{x}

$$\begin{aligned}
& \max_{k=1, \dots, K} f_k(\mathbf{x}) \pi_k \\
&= \max_{k=1, \dots, K} (2\pi)^{-\frac{p}{2}} \left\{ \det(\Sigma) \right\}^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mu_k)^\top \Sigma^{-1} (\mathbf{x} - \mu_k) \right\} \pi_k \\
&= \max_{k=1, \dots, K} \exp \left\{ -\frac{1}{2} (\mathbf{x}^\top \Sigma^{-1} \mathbf{x} - 2 \mathbf{x}^\top \Sigma^{-1} \mu_k + \mu_k^\top \Sigma^{-1} \mu_k) \right\} \pi_k \\
&= \max_{k=1, \dots, K} \exp \left(\mathbf{x}^\top \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^\top \Sigma^{-1} \mu_k \right) \pi_k \\
&= \max_{k=1, \dots, K} \left(\mathbf{x}^\top \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^\top \Sigma^{-1} \mu_k + \ln \pi_k \right)
\end{aligned}$$

- Empirical version
 - Training data: $\mathbf{x}_i \in \mathbb{R}^p$ and $y_i \in \{1, \dots, K\}$, $i = 1, \dots, n$
 - * n_k : the number of training observations in class k , $k = 1, \dots, K$
 - Estimation for μ_k , Σ and π_k
 - * $\hat{\pi}_k = n_k/n$
 - * $\hat{\mu}_k = n_k^{-1} \sum_{i=1}^n \mathbf{x}_i \cdot \mathbf{1}(y_i = k)$
 - * $\hat{\Sigma} = (n-1)^{-1} \sum_{k=1}^K \sum_{i=1}^n (\mathbf{x}_i - \hat{\mu}_k)(\mathbf{x}_i - \hat{\mu}_k)^\top \cdot \mathbf{1}(y_i = k)$
 - Empirical LDA classifier

$$\hat{h}(\mathbf{x}) = \arg \max_{k=1, \dots, K} \hat{\delta}_k(\mathbf{x})$$

$$* \hat{\delta}_k(\mathbf{x}) = \mathbf{x}^\top \hat{\Sigma}^{-1} \hat{\mu}_k - \frac{1}{2} \hat{\mu}_k^\top \hat{\Sigma}^{-1} \hat{\mu}_k + \ln \hat{\pi}_k$$

-
- Example (Fisher's or Anderson's iris data)
 - 50 flowers from each of 3 species of iris: setosa, versicolor, and virginica.
 - Measurements in centimeters of the variables sepal length and width and petal length and width.

```

options(digits = 4)
set.seed(3690)
picked = sample.int(nrow(iris), size = floor(nrow(iris)/3))
train = iris[-picked,]
Xtrain = train[, !(names(train) %in% c("Species"))]
Ytrain = train$Species
test = iris[picked,]
Xtest = test[, !(names(test) %in% c("Species"))]
Ytest = test[, names(test) %in% c("Species")]

# follow formulas
labels = unique(Ytrain)
K = length(labels)
p = ncol(Xtrain)
n = nrow(Xtrain)
nks = numeric(K)
piks = numeric(K)
Muks = matrix(0, nrow = K, ncol = p)
Sigmaks = list()
for (k in 1:K){
  Xtrain_k = Xtrain[Ytrain == labels[k],]
  nks[k] = nrow(Xtrain_k)
  piks[k] = nks[k]/n
  Muks[k,] = colMeans(Xtrain_k)
  Sigmaks[[k]] = cov(Xtrain_k)
  if (k==1){

```

```

    SigmaPool = Sigmaks[[k]] * (nks[k]-1)
  }else{
    SigmaPool = SigmaPool + Sigmaks[[k]] * (nks[k]-1)
  }
}
SigmaPool = SigmaPool/(n-1)
SigmaPoolInv = solve(SigmaPool)

deltaksLda = matrix(0, nrow = nrow(Xtest), ncol = K)
for (k in 1:K){
  deltaksLda[,k] = as.matrix(Xtest) %*% SigmaPoolInv %*% Muks[k,] -
    .5* as.vector(t(Muks[k,]) %*% SigmaPoolInv %*% Muks[k,]) +
    log(piks[k])
}
(resLda1 = apply(deltaksLda, 1, FUN = function(x){labels[which.max(x)]}))
# use MASS::lda
objLda = MASS::lda(Xtrain, Ytrain, method = "moment")
(resLda2 = predict(objLda, Xtest)$class)
# comparison
mean(resLda1 == resLda2)
mean(Ytest != resLda1)

```

Quadratic discriminant analysis (QDA)

- Assuming $f_k(\mathbf{x})$ = density of $MVN_p(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$
- QDA classifier

$$h(\mathbf{x}) = \arg \max_{k=1, \dots, K} \delta_k(\mathbf{x})$$

- Discriminant functions $\delta_k(\mathbf{x}) = -\mathbf{x}^\top \boldsymbol{\Sigma}_k^{-1} \mathbf{x} + 2\mathbf{x}^\top \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k - \boldsymbol{\mu}_k^\top \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k + 2 \ln \pi_k - \ln \det \boldsymbol{\Sigma}_k$
 * Quadratic functions with respect to \mathbf{x}

$$\begin{aligned}
 & \max_{k=1, \dots, K} f_k(\mathbf{x}) \pi_k \\
 &= \max_{k=1, \dots, K} \frac{1}{(2\pi)^{-\frac{p}{2}} |\det(\boldsymbol{\Sigma}_k)|^{-\frac{1}{2}}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^\top \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \right\} \pi_k \\
 &= \max_{k=1, \dots, K} \left\{ \det(\boldsymbol{\Sigma}_k) \right\}^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} \left(\mathbf{x}^\top \boldsymbol{\Sigma}_k^{-1} \mathbf{x} - 2\mathbf{x}^\top \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k + \boldsymbol{\mu}_k^\top \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k \right) \right\} \pi_k \\
 &= \max_{k=1, \dots, K} \left\{ -\mathbf{x}^\top \boldsymbol{\Sigma}_k^{-1} \mathbf{x} + 2\mathbf{x}^\top \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k - \boldsymbol{\mu}_k^\top \boldsymbol{\Sigma}_k^{-1} \boldsymbol{\mu}_k + 2 \ln \pi_k - \ln \det(\boldsymbol{\Sigma}_k) \right\}
 \end{aligned}$$

- Empirical version
 - Training data: $\mathbf{x}_i \in \mathbb{R}^p$ and $y_i \in \{1, \dots, K\}$, $i = 1, \dots, n$
 - * n_k : the number of training observations in class k , $k = 1, \dots, K$
 - Estimation for $\boldsymbol{\mu}_k$, $\boldsymbol{\Sigma}$ and π_k
 - * $\hat{\pi}_k = n_k/n$
 - * $\hat{\boldsymbol{\mu}}_k = n_k^{-1} \sum_{i=1}^n \mathbf{x}_i \cdot \mathbf{1}(y_i = k)$
 - * $\hat{\boldsymbol{\Sigma}}_k = (n_k - 1)^{-1} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)^\top \cdot \mathbf{1}(y_i = k)$
 - Empirical classifier

$$\hat{h}(\mathbf{x}) = \arg \max_{k=1, \dots, K} \hat{\delta}_k(\mathbf{x})$$

$$* \hat{\delta}_k(\mathbf{x}) = -\mathbf{x}^\top \hat{\boldsymbol{\Sigma}}_k^{-1} \mathbf{x} + 2\mathbf{x}^\top \hat{\boldsymbol{\Sigma}}_k^{-1} \hat{\boldsymbol{\mu}}_k - \hat{\boldsymbol{\mu}}_k^\top \hat{\boldsymbol{\Sigma}}_k^{-1} \hat{\boldsymbol{\mu}}_k + 2 \ln \hat{\pi}_k - \ln \det \hat{\boldsymbol{\Sigma}}_k$$

- Example (iris data, con'd)

```
# follow formulas
deltaksQda = matrix(0, nrow = nrow(Xtest), ncol = K)
for (k in 1:K){
  SigmakInv = solve(Sigmaks[[k]])
  deltaksQda[,k] = -diag(as.matrix(Xtest) %*% SigmakInv %*% t(as.matrix(Xtest))) +
    2* as.matrix(Xtest) %*% SigmakInv %*% Muks[k,] -
    as.vector(t(Muks[k,]) %*% SigmakInv %*% Muks[k,]) +
    2* log(piks[k]) -
    log(det(Sigmaks[[k]]))
}
(resQda1 = apply(deltaksQda, 1, FUN = function(x){labels[which.max(x)]}))
# use MASS::qda
objQda = MASS::qda(Xtrain, Ytrain, method = "moment")
(resQda2 = predict(objQda, Xtest)$class)
# comparison
mean(resQda1 == resQda2)
mean(Ytest != resQda1)
```