

# 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:  $28 \times 28$  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  $\mathbf{x}$

- 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  $\pi_k$ 
    - \*  $\hat{\pi}_k = n_k/n$
    - \*  $\hat{\boldsymbol{\mu}}_k = n_k^{-1} \sum_{i=1}^n \mathbf{x}_i \cdot \mathbf{1}(y_i = k)$
    - \*  $\hat{\boldsymbol{\Sigma}} = (n - K)^{-1} \sum_{k=1}^K \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)^{\top} \cdot \mathbf{1}(y_i = k)$
  - Empirical LDA classifier

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

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

## Quadratic discriminant analysis (QDA)

- Assuming  $f_k(\mathbf{x}) = \text{density of } MVN_p(\boldsymbol{\mu}_k, \Sigma_k)$
  - QDA classifier
$$h(\mathbf{x}) = \arg \max_{k=1, \dots, K} \delta_k(\mathbf{x})$$
    - Discriminant functions  $\delta_k(\mathbf{x}) = -\mathbf{x}^\top \Sigma_k^{-1} \mathbf{x} + 2\mathbf{x}^\top \Sigma_k^{-1} \boldsymbol{\mu}_k - \boldsymbol{\mu}_k^\top \Sigma_k^{-1} \boldsymbol{\mu}_k + 2 \ln \pi_k - \ln \det \Sigma_k$ 
      - \* Quadratic functions with respect to  $\mathbf{x}$
  - 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  $\boldsymbol{\mu}_k$ ,  $\Sigma$  and  $\pi_k$ 
      - \*  $\hat{\pi}_k = n_k/n$
      - \*  $\hat{\boldsymbol{\mu}}_k = n_k^{-1} \sum_{i=1}^n \mathbf{x}_i \cdot \mathbf{1}(y_i = k)$
      - \*  $\hat{\Sigma}_k = (n_k - 1)^{-1} \sum_{i=1}^n (\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)(\mathbf{x}_i - \hat{\boldsymbol{\mu}}_k)^\top \cdot \mathbf{1}(y_i = k)$
    - Empirical classifier
$$\hat{h}(\mathbf{x}) = \arg \max_{k=1, \dots, K} \hat{\delta}_k(\mathbf{x})$$
      - \*  $\hat{\delta}_k(\mathbf{x}) = -\mathbf{x}^\top \hat{\Sigma}_k^{-1} \mathbf{x} + 2\mathbf{x}^\top \hat{\Sigma}_k^{-1} \hat{\boldsymbol{\mu}}_k - \hat{\boldsymbol{\mu}}_k^\top \hat{\Sigma}_k^{-1} \hat{\boldsymbol{\mu}}_k + 2 \ln \hat{\pi}_k - \ln \det \hat{\Sigma}_k$
  - Example (`iris` data, con'd)
- 

## Assessment

- Misclassification rate
  - Population:  $\text{err} = \Pr(Y \neq h(\mathbf{X}))$
  - Empirical:  $\widehat{\text{err}} = (n^*)^{-1} \sum_{i=1}^{n^*} \mathbf{1}\{y_i \neq \hat{h}(\mathbf{x}_i)\}$ 
    - \* Testing data:  $\mathbf{x}_i \in \mathbb{R}^p$  and  $y_i \in \{1, \dots, K\}$ ,  $i = 1, \dots, n^*$
- Cross validation (CV)
  - Leave-one-out CV
  - $M$ -fold CV
    - \* Leave-one-out  $\Leftrightarrow n$ -fold