UVA CS 4501: Machine Learning

Lecture 20: Generative Gaussian Bayes Classifiers

Dr. Yanjun Qi

University of Virginia

Department of Computer Science

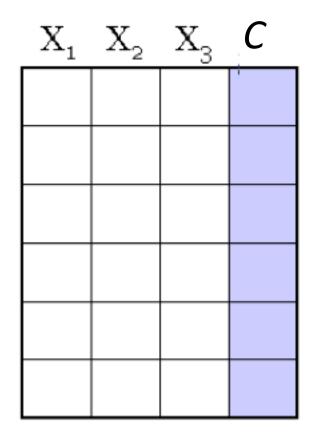
Where are we? Three major sections for classification

- We can divide the large variety of classification approaches into roughly three major types
 - 1. Discriminative
 - directly estimate a decision rule/boundary
 - e.g., support vector machine, decision tree



2. Generative:

- build a generative statistical model
- e.g., naïve bayes classifier, Bayesian networks
- 3. Instance based classifiers
 - Use observation directly (no models)
 - e.g. K nearest neighbors



A Dataset for classification-

$$f: [X] \longrightarrow [C]$$
Output as Discrete
Class Label
 $C_1, C_2, ..., C_L$

$$P(C \mid X)$$

- Data/points/instances/examples/samples/records: [rows]
- **Features**/attributes/dimensions/independent variables/covariates/predictors/regressors: [columns, except the last]
- Target/outcome/response/label/dependent variable: special column to be predicted [last column]

Today: Generative Bayes Classifiers

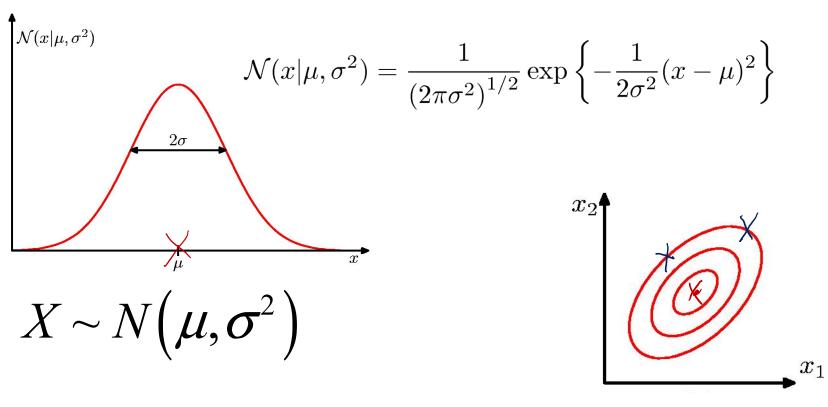


- ✓ Gaussian Bayes Classifiers
 - Gaussian distribution



- Naïve Gaussian BC
- Not-naïve Gaussian BC -> LDA, QDA

Gaussian Distribution



$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu},\boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{1\!\!\!P/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^{\mathrm{T}}\boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})\right\}$$
 Covariance Matrix

$$\underset{C}{\operatorname{argmax}} P(C \mid X) = \underset{C}{\operatorname{argmax}} P(X, C) = \underset{C}{\operatorname{argmax}} P(X \mid C) P(C)$$

Naïve Bayes Classifier

$$P(X | C) = P(X_{1}, X_{2}, \dots, X_{p} | C)$$

$$= P(X_{1} | X_{2}, \dots, X_{p}, C) P(X_{2}, \dots, X_{p} | C)$$

$$= P(X_{1} | C) P(X_{2}, \dots, X_{p} | C)$$

$$= P(X_{1} | C) P(X_{2} | C) \dots P(X_{p} | C)$$

$$\underset{C}{\operatorname{argmax}} P(C \mid X) = \underset{C}{\operatorname{argmax}} P(X, C) = \underset{C}{\operatorname{argmax}} P(X \mid C) P(C)$$

Naïve Bayes Classifier

$$P(X | C) = P(X_1, X_2, \dots, X_p | C)$$

$$= P(X_1 | X_2, \dots, X_p, C) P(X_2, \dots, X_p | C)$$

$$= P(X_1 | C) P(X_2, \dots, X_p | C)$$

$$= P(X_1 | C) P(X_2 | C) \dots P(X_p | C)$$

$$\hat{P}(X_j \mid C = c_i) = \frac{1}{\sqrt{2\pi}\sigma_{ji}} \exp\left(-\frac{(X_j - \mu_{ji})^2}{2\sigma_{ji}^2}\right)$$

 μ_{ji} : mean (avearage) of attribute values X_j of examples for which $C = c_i$

 σ_{ii} : standard deviation of attribute values X_i of examples for which $C = c_i$

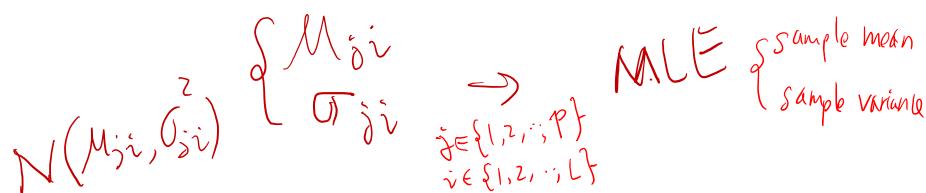
- Continuous-valued Input Attributes
 - Conditional probability modeled with the normal distribution

$$\hat{P}(X_j \mid C = c_i) = \frac{1}{\sqrt{2\pi}\sigma_{ji}} \exp\left(-\frac{(X_j - \mu_{ji})^2}{2\sigma_{ji}^2}\right)$$

 μ_{ii} : mean (avearage) of attribute values X_i of examples for which $C = c_i$

 σ_{ii} : standard deviation of attribute values X_i of examples for which $C = c_i$

- Learning Phase: for $\mathbf{X}=(X_1,\cdots,X_p), \quad C=c_1,\cdots,c_L$ Output: $p\times L$ normal distributions and $\quad P(C=c_i) \quad i=1,\cdots,L$



$$= \sum_{j \in \{1,2,3,7\}}$$

$$i, \in \{1,2,3,4\}$$

- Continuous-valued Input Attributes
 - Conditional probability modeled with the normal distribution

$$\hat{P}(X_{j} \mid C = c_{i}) = \frac{1}{\sqrt{2\pi}\sigma_{ii}} \exp\left(-\frac{(X_{j} - \mu_{ji})^{2}}{2\sigma_{ji}^{2}}\right)$$

 μ_{ji} : mean (avearage) of attribute values X_j of examples for which $C = c_i$

 σ_{ji} : standard deviation of attribute values X_j of examples for which $C = c_i$

- Learning Phase: for $\mathbf{X}=(X_1,\cdots,X_p), \quad C=c_1,\cdots,c_L$ Output: $p\times L$ normal distributions and $\quad P(C=c_i) \quad i=1,\cdots,L$
- Test Phase: for $\mathbf{X}' = (X_1', \dots, X_p')$

- Calculate conditional probabilities with all the normal distributions
- Apply the MAP rule to make a decision $\text{Argmax} \gamma((=(i))\gamma(x_i|i)) \gamma(x_i|i)$

Today: Generative Bayes Classifiers



- ✓ Gaussian Bayes Classifiers
 - Gaussian distribution
 - Naïve Gaussian BC



Not-naïve Gaussian BC -> LDA, QDA

Naïve Gaussian means?

Naïve
$$P(X_1, X_2, \dots, X_p \mid C = c_j) = P(X_1 \mid C)P(X_2 \mid C) \dots P(X_p \mid C)$$

$$= \prod_i \frac{1}{\sqrt{2\pi}\sigma_{ji}} \exp\left(-\frac{(X_j - \mu_{ji})^2}{2\sigma_{ji}^2}\right) \sum_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma_{ji}} \exp\left(-\frac{(X_i - \mu_{ji})^2}{2\sigma_{ji}^2}\right) \sum_{i=1}^{n} \frac{1}{\sqrt{2\pi}\sigma_{ji}} \exp\left($$

Diagonal Matrix
$$\sum c_k = \Lambda c_k$$

covariance matrix is diagonal

Not Naïve Gaussian means?



Not Naïve

$$\begin{split} P(X_1, X_2, \cdots, X_p \mid C) = \\ \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= \frac{1}{(2\pi)^{D/2}} \frac{1}{|\boldsymbol{\Sigma}|^{1/2}} \exp\left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^{\mathrm{T}} \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\} \end{split}$$

 $P(X_1, X_2, \dots, X_p \mid C = c_j) = P(X_1 \mid C)P(X_2 \mid C) \dots P(X_p \mid C)$

Naïve

44122188

$$= \prod_{i} \frac{1}{\sqrt{2\pi}\sigma_{ji}} \exp\left(-\frac{(X_{j} - \mu_{ji})^{2}}{2\sigma_{ji}^{2}}\right) \ge \left(-\frac{(X_{j} - \mu_{ji})^{2}}{2\sigma_{ji}^{2}}\right)$$

Diagonal Matrix

$$\sum c_k = \Lambda c_k$$

Each class' covariance matrix is diagonal

Dr. Yanjun Qi / UVA CS

Today: Generative Bayes Classifiers



- Gaussian distribution
- Naïve Gaussian BC
- Not-naïve Gaussian BC



- LDA: Linear Discriminant Analysis
- QDA: Quadratic Discriminant Analysis

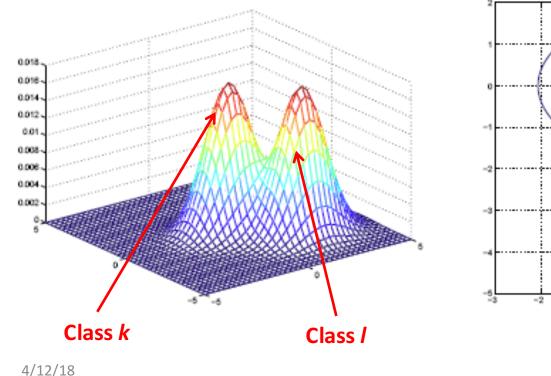
(1) covariance matrix are the same across classes

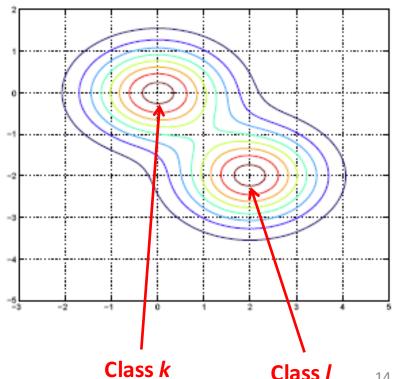
→ LDA (Linear Discriminant Analysis)

Linear Discriminant Analysis: $\sum_{k} = \sum_{k} \forall k$

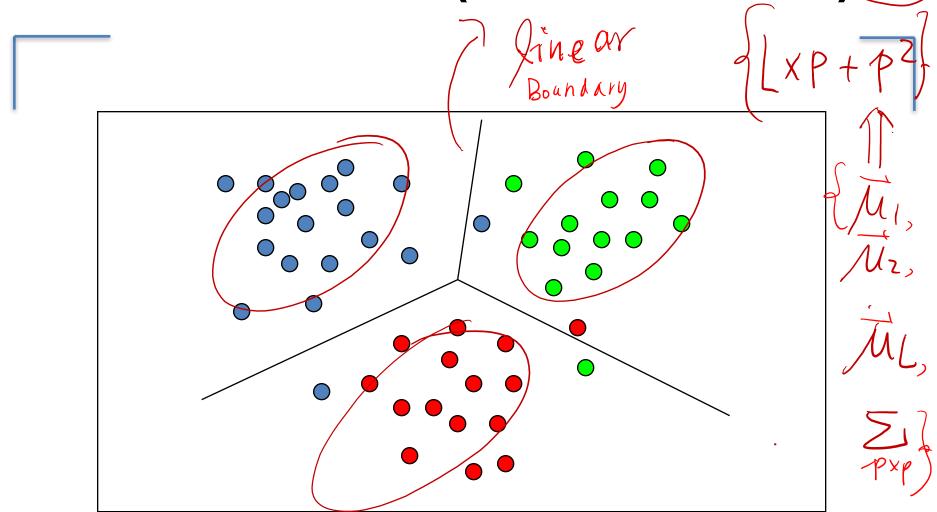
Each class' covariance

The Gaussian Distribution are shifted versions of each other





Visualization (three classes)



$$\underset{k}{\operatorname{argmax}} P(C_{k}|X) = \underset{k}{\operatorname{argmax}} P(X,C_{k}) = \underset{k}{\operatorname{argmax}} P(X|C_{k}) P(C_{k})$$

$$= \underset{k}{\operatorname{argmax}} \{P(X|C_{k}) P(C_{k})\}$$

Decision Boundary means those points Satisfying: p(Ci(X) = p(Cj(X)) $\frac{p(c_i|x)}{p(c_i|x)} = 1$ $p(c_i|x)$ = 0 $p(c_i|x)$ = 0

$$\operatorname{argmax} P(C_{k}|X) = \operatorname{argmax} P(X,C_{k}) = \operatorname{argmax} P(X|C_{k})P(C_{k})$$

$$= \operatorname{argmax} \log \{P(X|C_{k})P(C_{k})\}$$

$$= \operatorname{argmax} P(X|C_{k})P(C_{k})\}$$

$$= \operatorname{argmax} P(X|C_{k})P(C_{k})$$

$$\log \frac{P(C_k|X)}{P(C_l|X)} = \log \frac{P(X|C_k)}{P(X|C_l)} + \log \frac{P(C_k)}{P(C_l)}$$

$$= \log \frac{\pi_k}{\pi_\ell} - \frac{1}{2} (\mu_k + \mu_\ell)^T \mathbf{\Sigma}^{-1} (\mu_k - \mu_\ell) + x^T \mathbf{\Sigma}^{-1} (\mu_k - \mu_\ell),$$
(4.9)

19

$$\log \frac{P(C_k|X)}{P(C_i|X)} = \log \frac{P(X|C_k)}{P(X|C_i)} + \log \frac{P(C_k)}{P(C_i)}$$

$$= \log \frac{\pi_k}{\pi_\ell} - \frac{1}{2} (\mu_k + \mu_\ell)^T \mathbf{\Sigma}^{-1} (\mu_k - \mu_\ell) + x^T \mathbf{\Sigma}^{-1} (\mu_k - \mu_\ell),$$
(4.9)

The above is derived from the following:

$$-\frac{1}{2}(x-\mu_k)^T \Sigma^{-1}(x-\mu_k) \, = \, x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k - \frac{1}{2} x^T \Sigma^{-1} x$$

$$\log \frac{P(C_k|X)}{P(C_l|X)} = \log \frac{P(X|C_k)}{P(X|C_l)} + \log \frac{P(C_k)}{P(C_l)}$$

$$= \underbrace{\log \frac{\pi_k}{\pi_\ell} - \frac{1}{2} (\mu_k + \mu_\ell)^T \Sigma^{-1} (\mu_k - \mu_\ell)}_{+x^T \Sigma^{-1} (\mu_k - \mu_\ell)}$$

$$+ \underbrace{x^T \Sigma^{-1} (\mu_k - \mu_\ell)}_{\alpha}$$

$$\Rightarrow \quad x^T \alpha + b = 0 \Rightarrow \quad \alpha \quad \text{linear line}$$

$$\text{decision boundary}$$

IS

LDA Classification Rule (also called as

Linear discriminant function:)

$$\underset{k}{\operatorname{argmax}} P(C_{k} | X) = \underset{k}{\operatorname{argmax}} P(X, C_{k}) = \underset{k}{\operatorname{argmax}} P(X | C_{k}) P(C_{k})$$

$$= \arg \max_{k} \left[-\log((2\pi)^{p/2} |\Sigma|^{1/2}) - \frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) + \log(\pi_k) \right]$$

$$= \arg \max_{k} \left[-\frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) + \log(\pi_k) \right]$$

- Note

Linear Discriminant Function for LDA

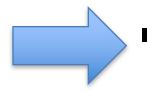
$$-\frac{1}{2}(x-\mu_k)^T \Sigma^{-1}(x-\mu_k) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k - \frac{1}{2} x^T \Sigma^{-1} x$$

Today: Generative Bayes Classifiers



- Gaussian distribution
- Naïve Gaussian BC
- Not-naïve Gaussian BC

LDA: Linear Discriminant Analysis



QDA: Quadratic Discriminant Analysis

(2) If covariance matrix are not the same e.g. → QDA (Quadratic Discriminant Analysis)

- Estimate the covariance matrix Σ_k separately for each class k, k = 1, 2, ..., K.
- Quadratic discriminant function:

$$\delta_k(x) = -\frac{1}{2} \log |\Sigma_k| - \frac{1}{2} (x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k) + \log \pi_k.$$

Classification rule:

$$\hat{G}(x) = \arg\max_{k} \delta_k(x)$$
.

- Decision boundaries are quadratic equations in x.
- QDA fits the data better than LDA, but has more parameters to estimate.

(2) If covariance matrix are not the same e.g. -> QDA (Quadratic Discriminant Analysis)

- \triangleright Estimate the covariance matrix Σ_k separately for each class k, k = 1, 2, ..., K.
- Quadratic discriminant function:

Quadratic discriminant function:
$$(\sum_{k} \sum_{k} \sum_{k}$$

Classification rule:

$$\hat{G}(x) = \arg\max_{k} \delta_{k}(x) .$$



QDA fits the data better than LDA, but has more parameters to estimate.

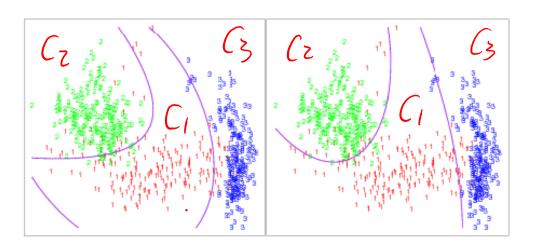




QDA vs. LDA on Expanded Basis

- ▶ Expand input space to include X_1X_2 , X_1^2 , and X_2^2 .
- ▶ Input is five dimensional: $X = (X_1, X_2, X_1X_2, X_1^2, X_2^2)$.

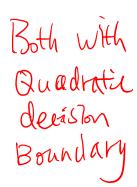
LPA With Q(X)



RDA

LDA with quadratic basis Versus QDA

Figure 4.6: Two methods for fitting quadratic boundaries. The left plot shows the quadratic decision boundaries for the data in Figure 4.1 (obtained using LDA in the five-dimensional space $x_1, x_2, x_{12}, x_1^2, x_2^2$). The right plot shows the quadratic decision boundaries found by QDA. The differences are small, as is usually the case.



(3) Regularized Discriminant Analysis

- A compromise between LDA and QDA.
- Shrink the separate covariances of QDA toward a common covariance as in LDA.
- Regularized covariance matrices:

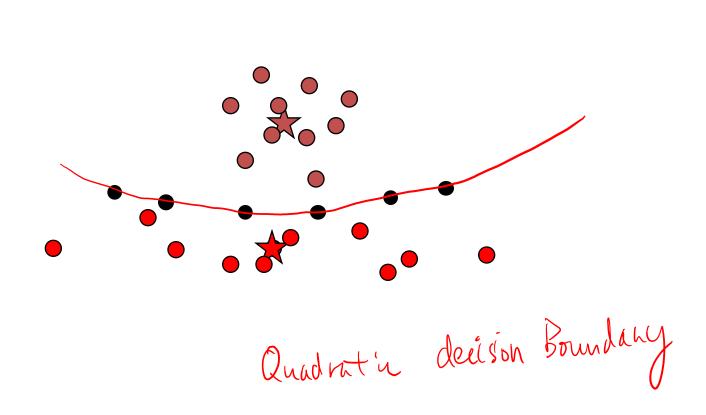
$$\hat{\Sigma}_k(\alpha) = \alpha \hat{\Sigma}_k + (1 - \alpha) \hat{\Sigma} .$$

- ▶ The quadratic discriminant function $\delta_k(x)$ is defined using the shrunken covariance matrices $\hat{\Sigma}_k(\alpha)$.
- ▶ The parameter α controls the complexity of the model.

An example: Gaussian Bayes Classifier

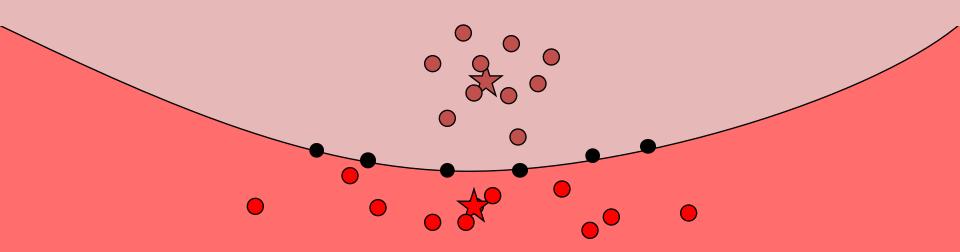
 $\Sigma_1 + \Sigma_7$ Naive BC diagonal Siz/I

Gaussian Bayes Classifier



Gaussian Bayes Classifier

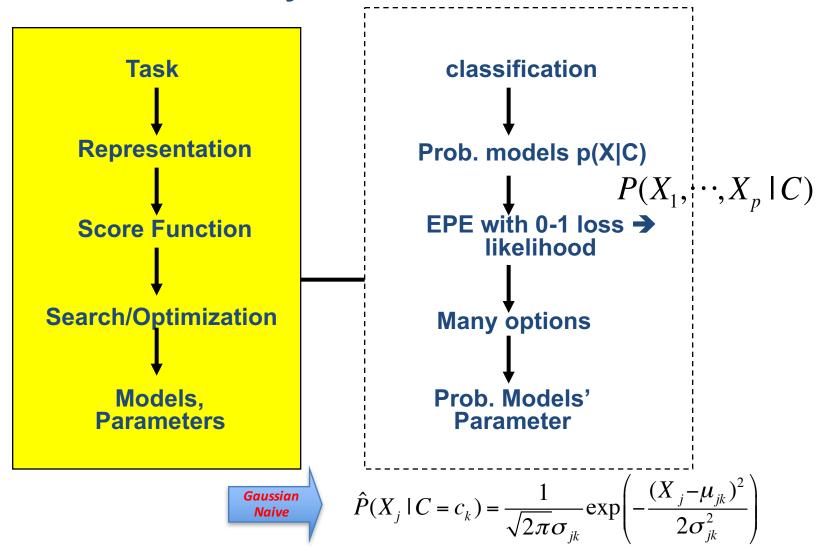
Blue Team



Red Team

Naïve Gaussian Bayes Classifier is not a linear classifier!

Generative Gaussian Bayes Classifier



 $\operatorname{argmax} P(C_k \mid X) = \operatorname{argmax} P(X, C) = \operatorname{argmax} P(X \mid C) P(C)$

References

- Prof. Andrew Moore's review tutorial
- ☐ Prof. Ke Chen NB slides
- ☐ Prof. Carlos Guestrin recitation slides