

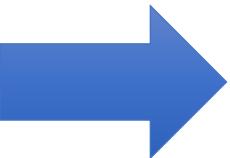
UVA CS 4501: Machine Learning

Lecture 17b: Gaussian BC and Generative vs. Discriminative Classifier

Dr. Yanjun Qi

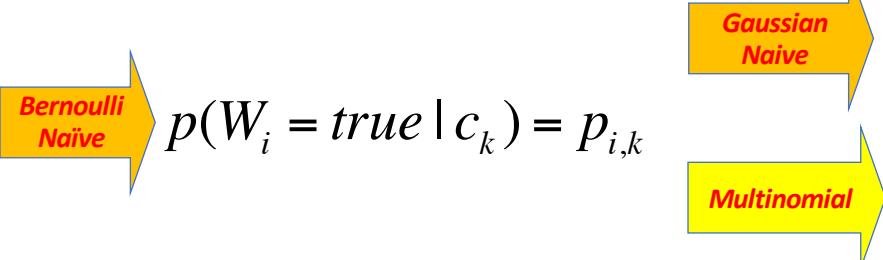
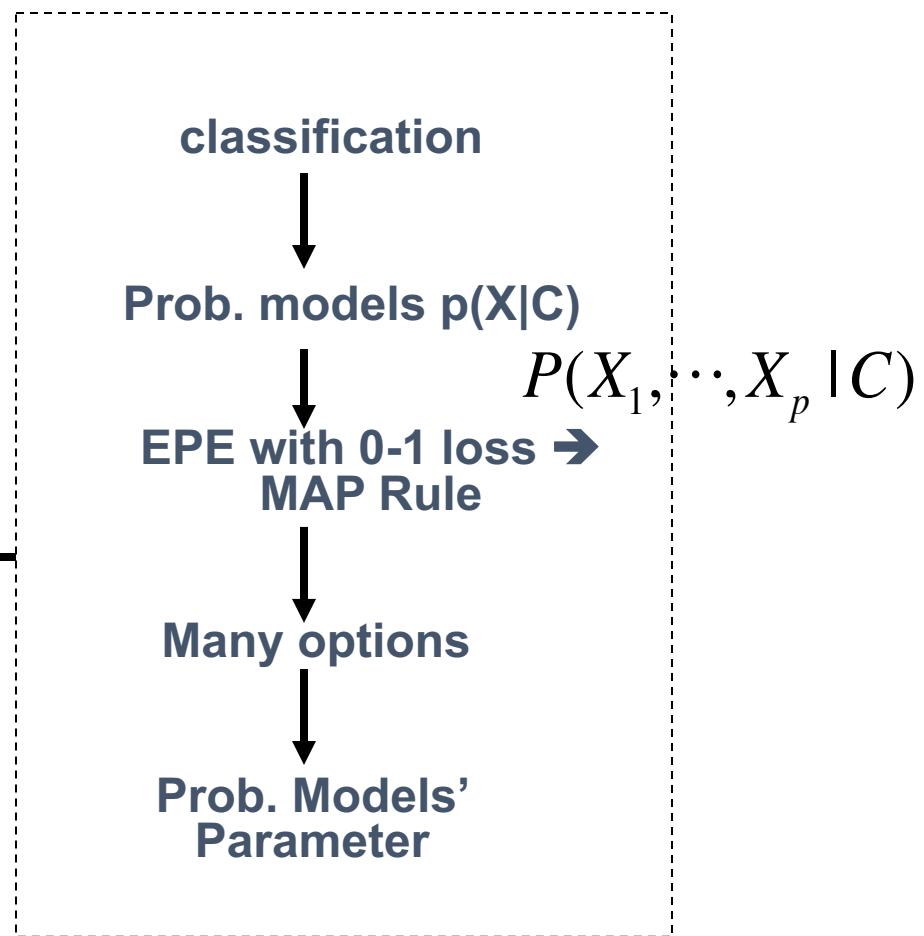
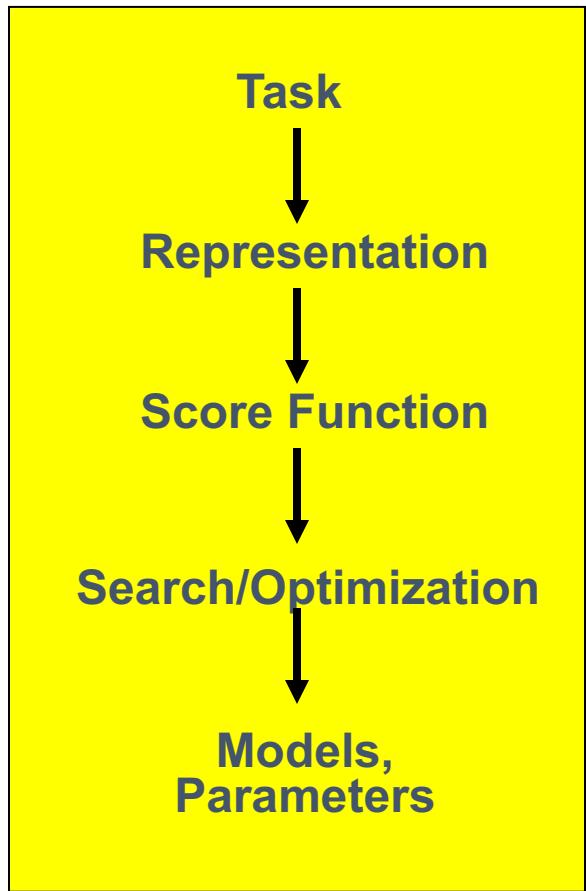
University of Virginia
Department of Computer Science

Today: More Generative Bayes Classifiers

- 
- ✓ Generative Bayes Classifier
 - ✓ Naïve Bayes Classifier
 - ✓ Gaussian Bayes Classifiers
 - Gaussian distribution
 - Naïve Gaussian BC
 - Not-naïve Gaussian BC → LDA, QDA
 - ✓ Discriminative vs. Generative

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

Generative Bayes Classifier



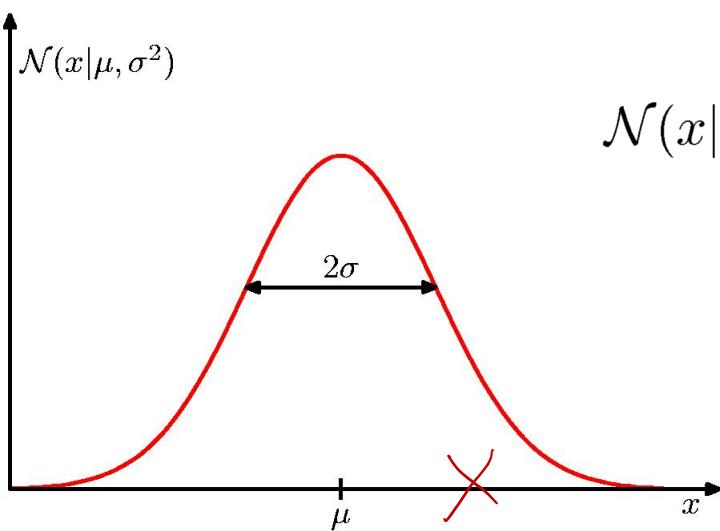
$$\hat{P}(X_j | C = c_k) = \frac{1}{\sqrt{2\pi}\sigma_{jk}} \exp\left(-\frac{(X_j - \mu_{jk})^2}{2\sigma_{jk}^2}\right)$$

$$P(W_1 = n_1, \dots, W_v = n_v | c_k) = \frac{N!}{n_{1k}! n_{2k}! \dots n_{vk}!} \theta_{1k}^{n_{1k}} \theta_{2k}^{n_{2k}} \dots \theta_{vk}^{n_{vk}}$$

Review: Continuous Random Variables

- Probability density function (pdf) instead of probability mass function (pmf)
 - For discrete RV: Probability mass function (pmf): $P(X = x_i)$
- A pdf (prob. Density func.) is any function $f(x)$ that describes the probability density in terms of the input variable x .

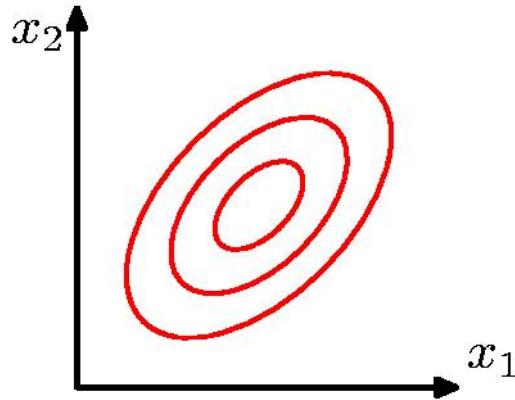
Single-Variate Gaussian Distribution



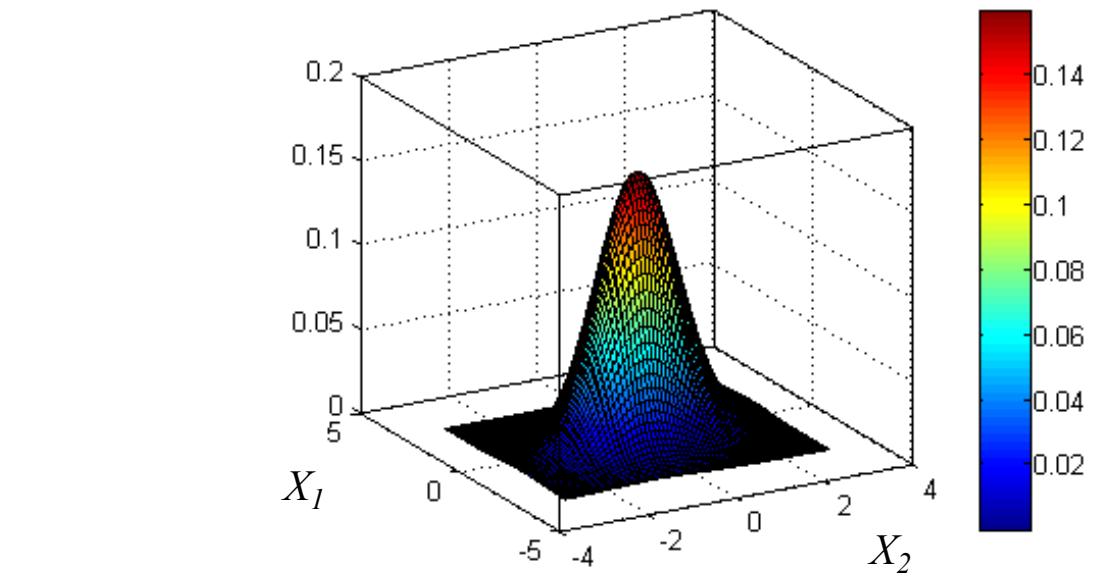
$$N(x|\mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2}(x - \mu)^2 \right\}$$

$$X \sim \mathcal{N}(\mu, \sigma^2)$$

Bi-Variate Gaussian Distribution



Bivariate
normal PDF:



- Mean of normal PDF is at peak value. Contours of equal PDF form ellipses.

- The covariance matrix captures linear dependencies among the variables

Multivariate Normal (Gaussian) PDFs

The only widely used continuous joint PDF is the multivariate normal (or Gaussian):

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

Where $|*$ represents determinant

Mean Covariance Matrix

- Mean of normal PDF is at peak value. Contours of equal PDF form ellipses.

- The covariance matrix captures linear dependencies among the variables

Example: the Bivariate Normal distribution

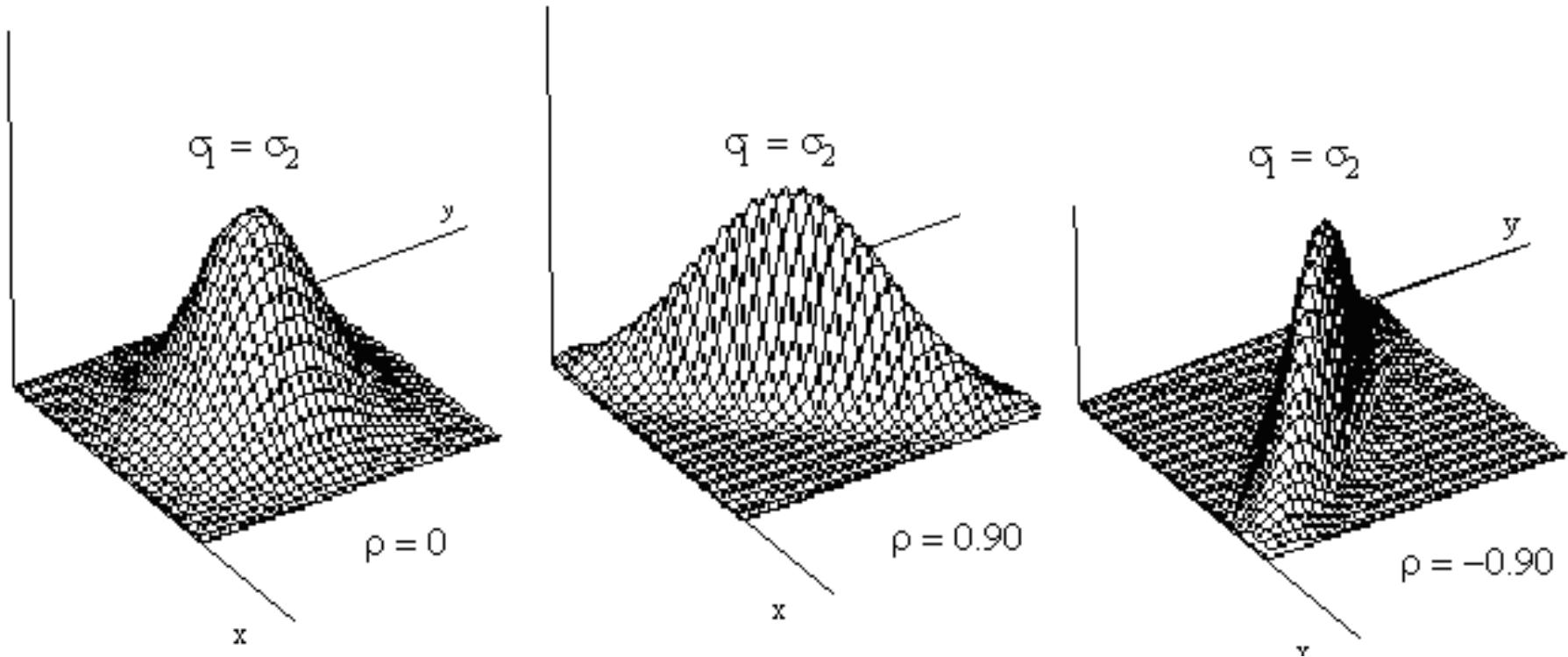
$$f(x_1, x_2) = \frac{1}{(2\pi)^{1/2} |\Sigma|} e^{-\frac{1}{2} (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})}$$

with $\vec{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}$ and

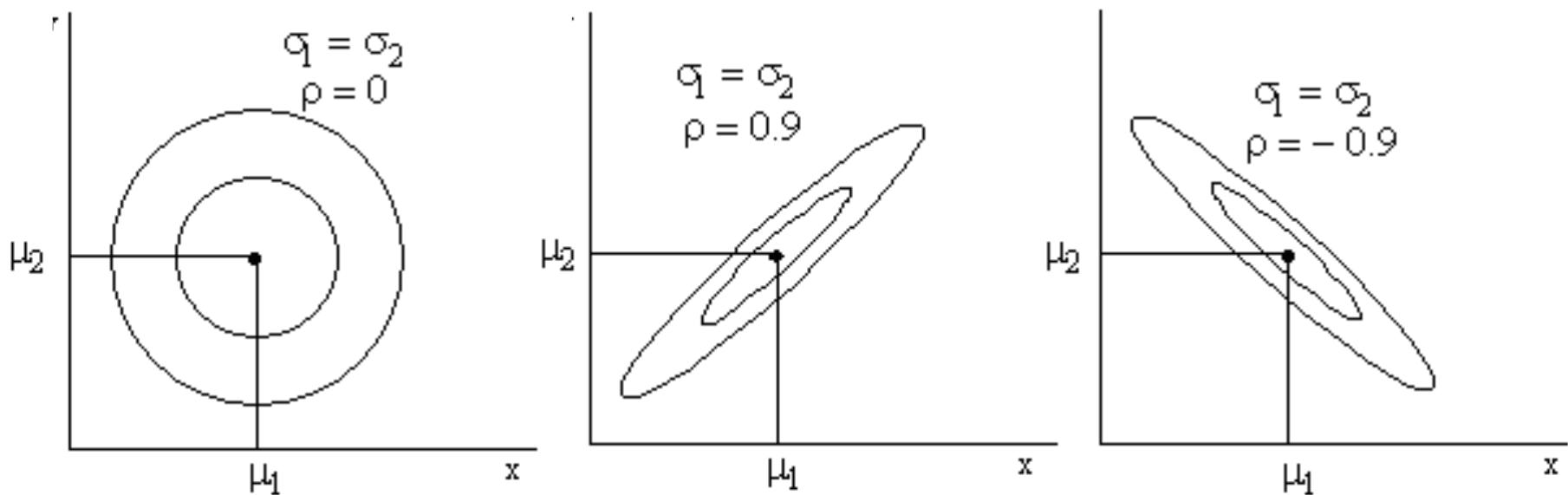
$$\Sigma_{2 \times 2} = \begin{bmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{12} & \sigma_{22} \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & \text{Cov}(x_1, x_2) \\ \underbrace{\rho \sigma_1 \sigma_2}_{\text{Cov}(x_1, x_2)} & \sigma_2^2 \end{bmatrix}_{2 \times 2}$$

$$|\Sigma| = \sigma_{11}\sigma_{22} - \sigma_{12}^2 = \sigma_1^2 \sigma_2^2 (1 - \rho^2)$$

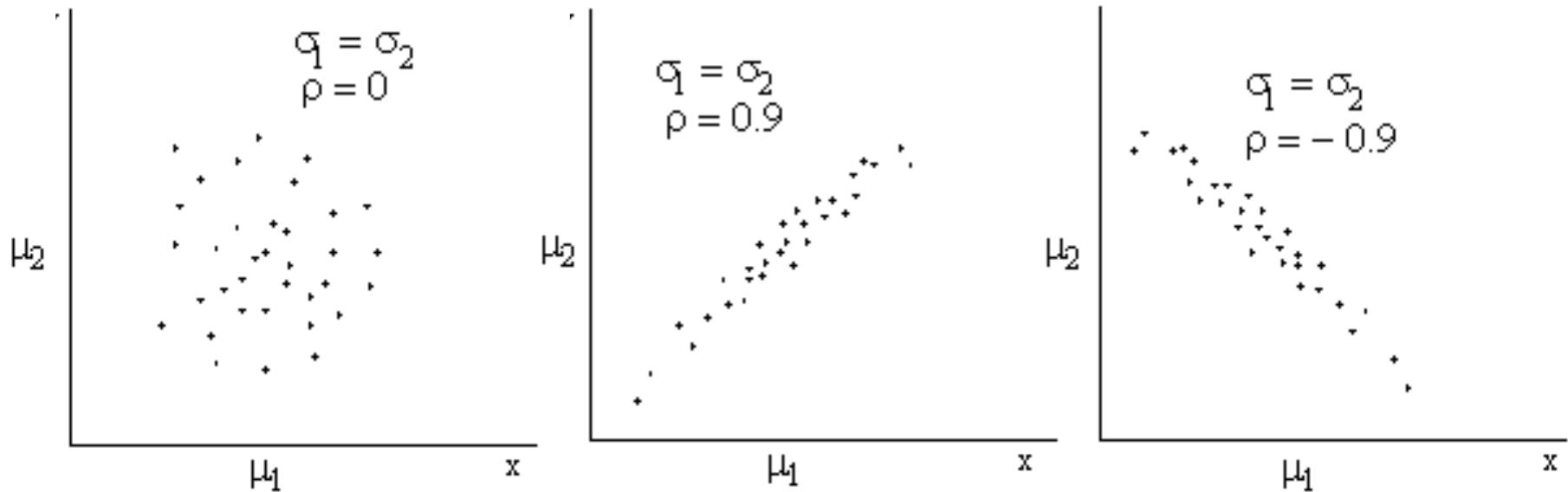
Surface Plots of the bivariate Normal distribution



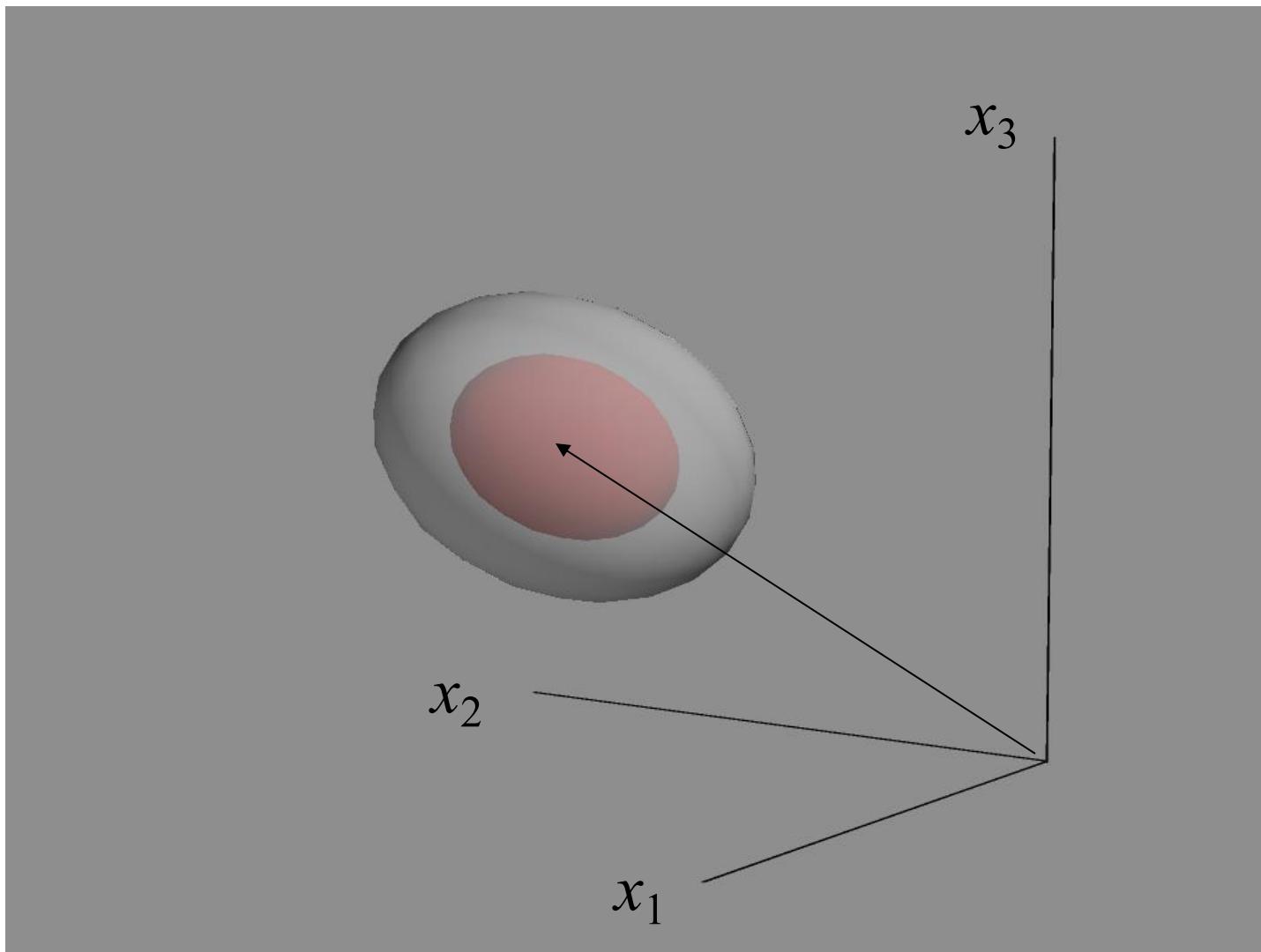
Contour Plots of the bivariate Normal distribution



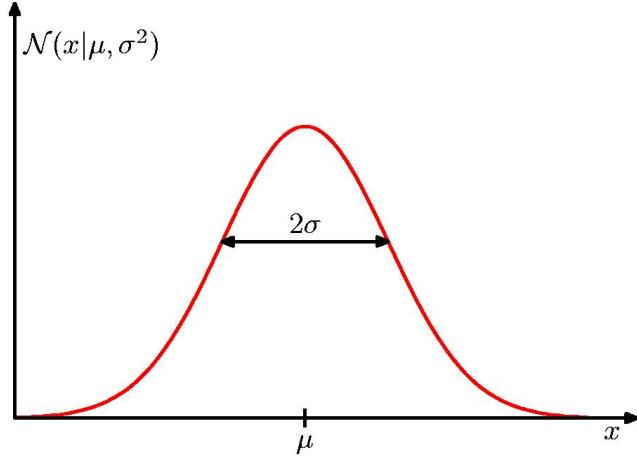
Scatter Plots of data from the bivariate Normal distribution



Trivariate Normal distribution



How to Estimate 1D Gaussian: MLE



- In the 1D Gaussian case, we simply set the mean and the variance to the **sample mean** and the **sample variance**:

$$\bar{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\bar{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{\mu})^2$$

How to Estimate p-D Gaussian: MLE

$$\langle X_1, X_2 \dots, X_p \rangle \sim N(\vec{\mu}, \Sigma)$$

$$\vec{\mu} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_p \end{bmatrix} \quad p \times 1$$

$$\mu_i = \frac{1}{n} \sum_{j=1}^N X_j^{(i)}$$

$\in \{1, 2, \dots, p\}$

i -th feature

j -th sample

$\in \{1, 2, \dots, N\}$

$$\Sigma = \begin{bmatrix} \text{Var}(X_1) & & \\ & \ddots & \\ & & \text{Var}(X_p) \end{bmatrix} \quad p \times p$$

i

j

$\text{Cov}(X_i, X_j)$

Gaussian Naïve Bayes Classifier

$$\operatorname{argmax}_C P(C | X) = \operatorname{argmax}_C P(X, C) = \operatorname{argmax}_C P(X | C)P(C)$$

Naïve
Bayes
Classifier

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

Gaussian Naïve Bayes Classifier

$$\operatorname{argmax}_C P(C | X) = \operatorname{argmax}_C P(X, C) = \operatorname{argmax}_C P(X | C)P(C)$$

Naïve
Bayes
Classifier

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

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

μ_{ji} : mean (average) 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 $C = c_i$

Gaussian Naïve Bayes Classifier

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

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

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

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

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

$$N(\mu_{ji}, \sigma_{ji}^2) \quad \left\{ \begin{array}{l} \mu_{ji} \\ \sigma_{ji} \end{array} \right. \quad \xrightarrow{j \in \{1, 2, \dots, p\}} \quad \text{MLE} \quad \left\{ \begin{array}{l} \text{sample mean} \\ \text{sample variance} \end{array} \right.$$

$j \in \{1, 2, \dots, p\}$

$i \in \{1, 2, \dots, L\}$

Gaussian Naïve Bayes Classifier

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

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

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

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

- **Learning Phase:** for $\mathbf{X} = (X_1, \dots, X_p)$, $C = c_1, \dots, c_L$
Output: $p \times L$ normal distributions and $P(C = c_i)$ $i = 1, \dots, L$
- **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 $\arg\max_i p(C=c_i) p(X'_1|c_i) \dots p(X'_p|c_i)$

Naïve

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

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

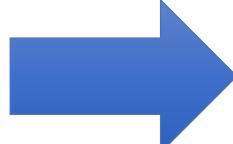
$\sum |C_i| = \begin{bmatrix} \sigma_{11} & & & \\ & \sigma_{22} & & \\ & & \ddots & \\ & & & \sigma_{pp} \end{bmatrix}$

Diagonal Matrix

$$\sum c_k = \Lambda c_k$$

Each class' covariance matrix is diagonal

Today: More Generative Bayes Classifiers

- ✓ Generative Bayes Classifier
 - ✓ Naïve Bayes Classifier
 - ✓ Gaussian Bayes Classifiers
 - Gaussian distribution
 - Naïve Gaussian BC
 - Not-naïve Gaussian BC → LDA, QDA
 - ✓ Discriminative vs. Generative
- 

Not Naïve Gaussian means ?

$$\text{Total \# param} \Rightarrow L \times \{P + P \times P\}$$

μ/C
 Σ/C

Not
Naïve

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

Total # param $\Rightarrow L \times (P + P)$

Naïve

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

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

$\sum |C_i| = \begin{bmatrix} \sigma_{11} & & \\ & \ddots & \\ & & \sigma_{pp} \end{bmatrix}$

Diagonal Matrix

$$\sum c_k = \Lambda c_k$$

Each class' covariance matrix is diagonal

Today : Generative Bayes Classifiers

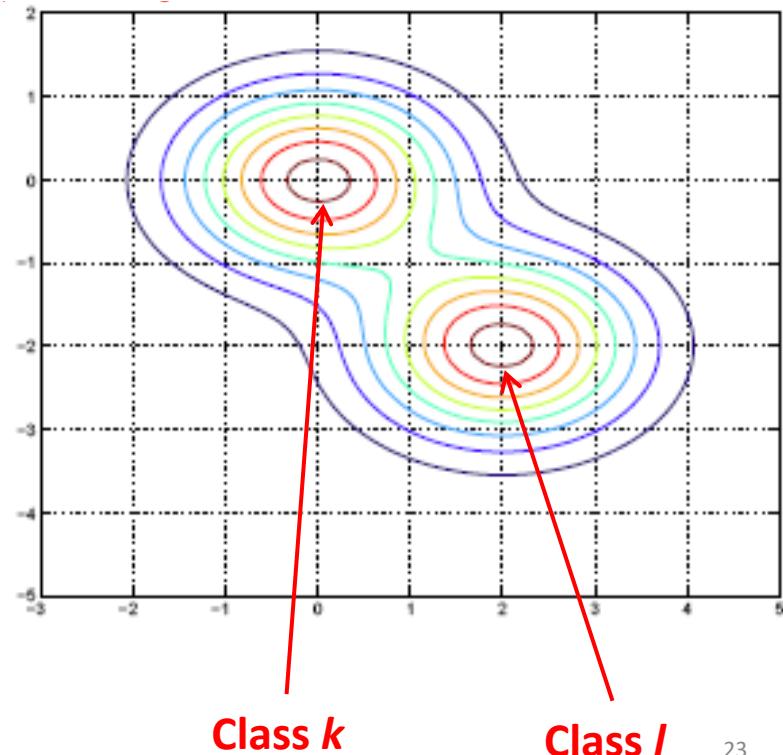
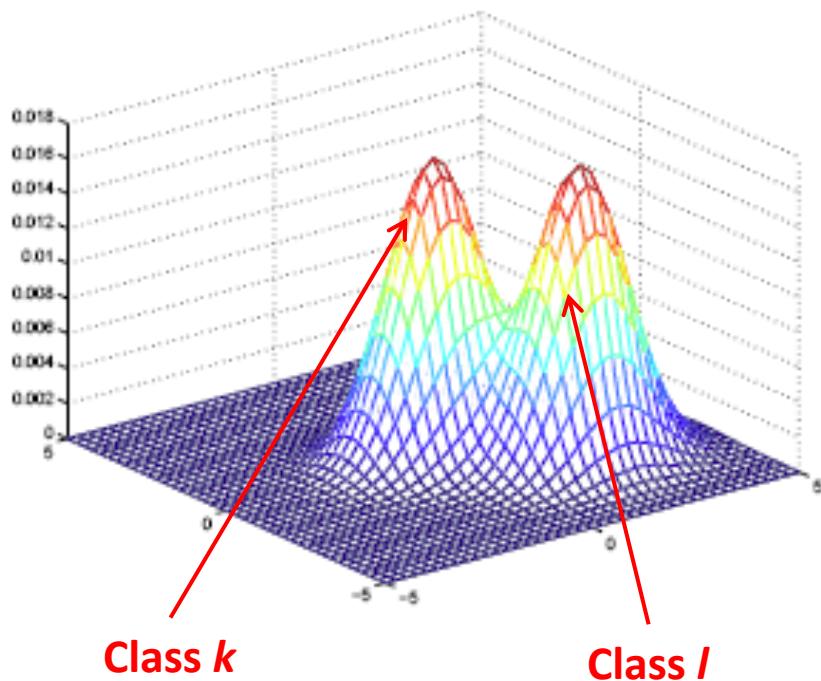
- ✓ Gaussian Bayes Classifiers
 - Gaussian distribution
 - Naïve Gaussian BC
 - Not-naïve Gaussian BC
- ■ LDA: Linear Discriminant Analysis
- QDA: Quadratic Discriminant Analysis

(1) covariance matrix are the same across classes
→ LDA (Linear Discriminant Analysis)

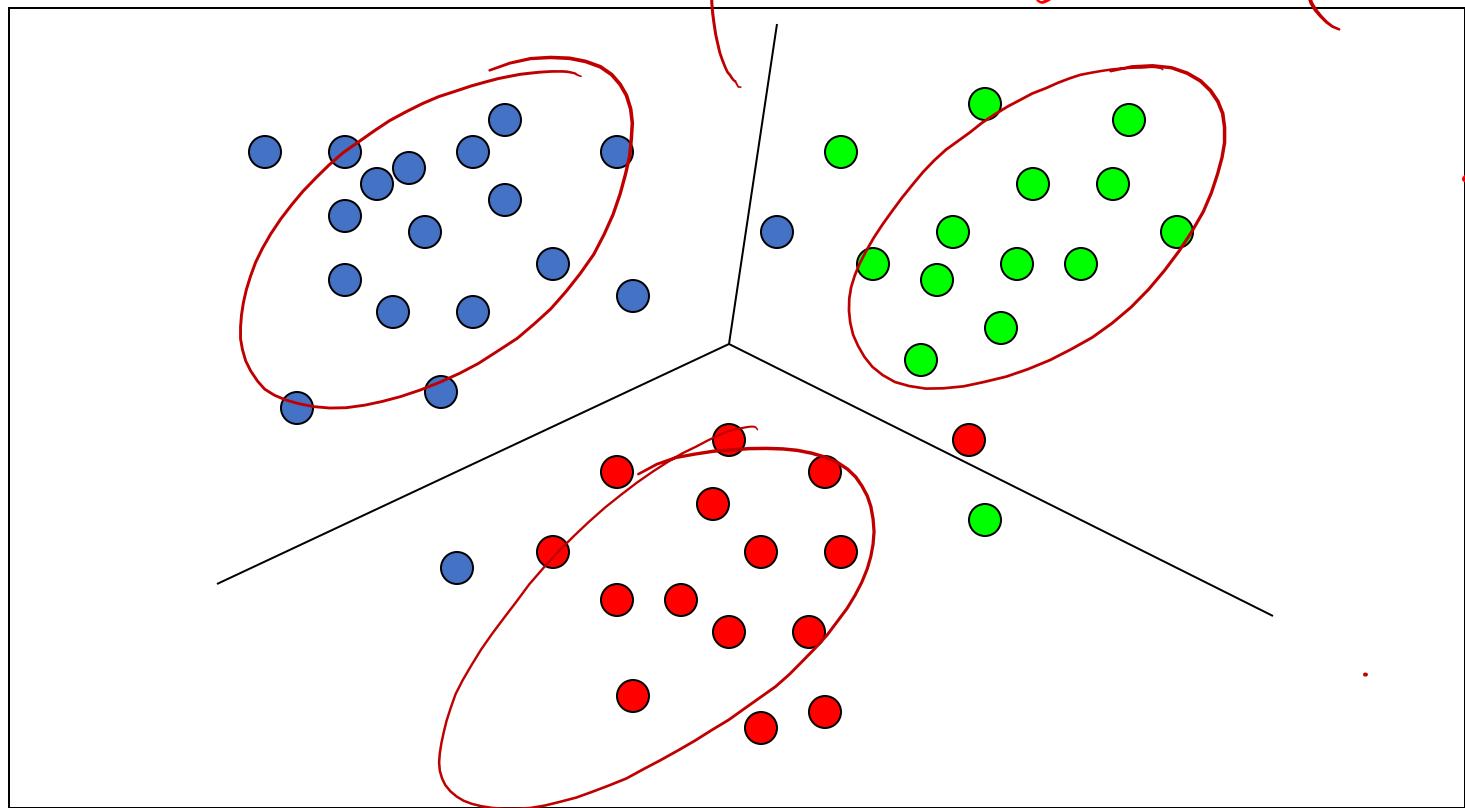
Linear Discriminant Analysis : $\Sigma_k = \Sigma$, $\forall k$ PXP

Each class' covariance matrix is the same

The Gaussian Distribution are shifted versions of each other



Visualization (three classes)



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

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

Decision Boundary means those points

satisfying: $P(C_i | X) = P(C_j | X)$

$$\frac{P(C_i | X)}{P(C_j | X)} = 1$$

$$\Rightarrow \log \frac{P(C_i | X)}{P(C_j | X)} = 0$$

$$\operatorname{argmax}_k P(C_k | X) = \operatorname{argmax}_k P(X, C_k) = \operatorname{argmax}_k P(X | C_k) P(C_k)$$

$$= \operatorname{argmax}_k \log \{ P(X | C_k) P(C_k) \}$$

$$= \operatorname{argmax}_k \log P(X | C_k) + \log P(C_k) \Rightarrow \pi_k$$

Decision Boundary points

$$\log \frac{P(C_k | X)}{P(C_l | X)} = 0 = \log \frac{P(X | C_k)}{P(X | C_l)} + \log \frac{\pi_k}{\pi_l}$$

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

$$\begin{aligned}
& \log \frac{P(C_k | X)}{P(C_l | X)} = \log \frac{P(X | C_k)}{P(X | C_l)} + \log \frac{P(C_k)}{P(C_l)} \\
&= \log \frac{\pi_k}{\pi_\ell} - \frac{1}{2} (\mu_k + \mu_\ell)^T \Sigma^{-1} (\mu_k - \mu_\ell) \\
&\quad + x^T \Sigma^{-1} (\mu_k - \mu_\ell), \tag{4.9}
\end{aligned}$$

$$\begin{aligned}
& \log \frac{P(C_k | X)}{P(C_l | X)} = \log \frac{P(X | C_k)}{P(X | C_l)} + \log \frac{P(C_k)}{P(C_l)} \\
&= \log \frac{\pi_k}{\pi_\ell} - \frac{1}{2}(\mu_k + \mu_\ell)^T \Sigma^{-1} (\mu_k - \mu_\ell) \\
&\quad + x^T \Sigma^{-1} (\mu_k - \mu_\ell), \tag{4.9}
\end{aligned}$$

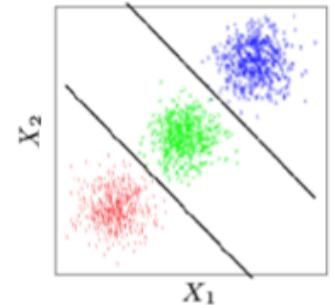
The above is derived from the following :

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

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

$$= \underbrace{\log \frac{\pi_k}{\pi_\ell} - \frac{1}{2}(\mu_k + \mu_\ell)^T \Sigma^{-1} (\mu_k - \mu_\ell)}_{+ \underbrace{x^T \Sigma^{-1} (\mu_k - \mu_\ell)}_a} \quad b \quad (4.9)$$

$\Rightarrow x^T a + b = 0 \Rightarrow$ a linear line
decision boundary



LDA Classification Rule (also called as Linear discriminant function:)

$$\operatorname{argmax}_k P(C_k | X) = \operatorname{argmax}_k P(X, C_k) = \operatorname{argmax}_k P(X | C_k) P(C_k)$$

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

$$= \operatorname{argmax}_k \boxed{-\frac{1}{2}(x - \mu_k)^T \Sigma^{-1} (x - \mu_k) + \log(\pi_k)}$$

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: More Generative Bayes Classifiers

- ✓ Generative Bayes Classifier
- ✓ Naïve Bayes Classifier
- ✓ Gaussian Bayes Classifiers
 - Gaussian distribution
 - Naïve Gaussian BC
 - Not-naïve Gaussian BC → LDA, QDA
- ✓ Discriminative vs. Generative

(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 ,
 $k = 1, 2, \dots, K$.
- ▶ *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)

- ▶ Estimate the covariance matrix Σ_k separately for each class k ,
 $k = 1, 2, \dots, K$.

- ▶ Quadratic discriminant function:

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

$$\{\Sigma_1, \Sigma_2, \dots, \Sigma_K, \mu_1, \mu_2, \dots, \mu_K, \pi_1, \pi_2, \dots, \pi_K\}$$

- ▶ Classification rule:

$$\delta_1(x) - \delta_2(x) = 0$$

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

Total # para

$$K \times (P + P^2)$$

$$\{\mu_k, \Sigma_k\}$$

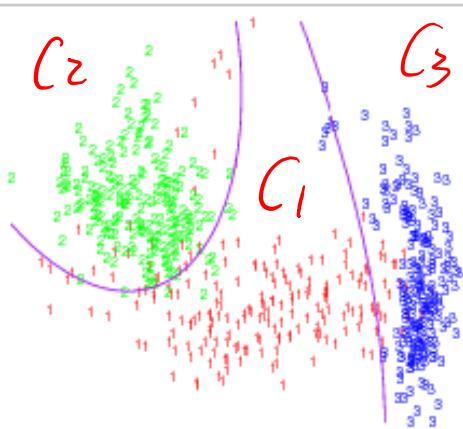
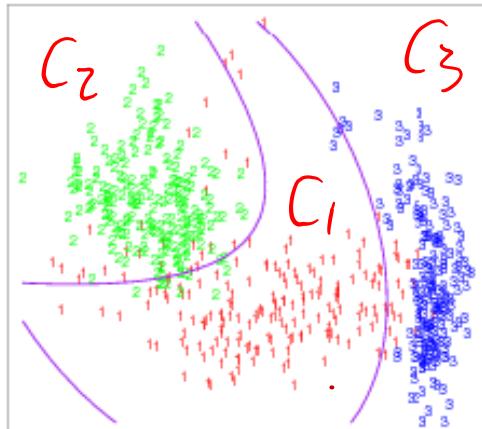
- ▶ Decision boundaries are quadratic equations in x .

- ▶ 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: $X = (X_1, X_2, X_1X_2, X_1^2, X_2^2)$.

LDA
With
 $Q(X)$



QDA

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.

Both with
Quadratic
decision
Boundary

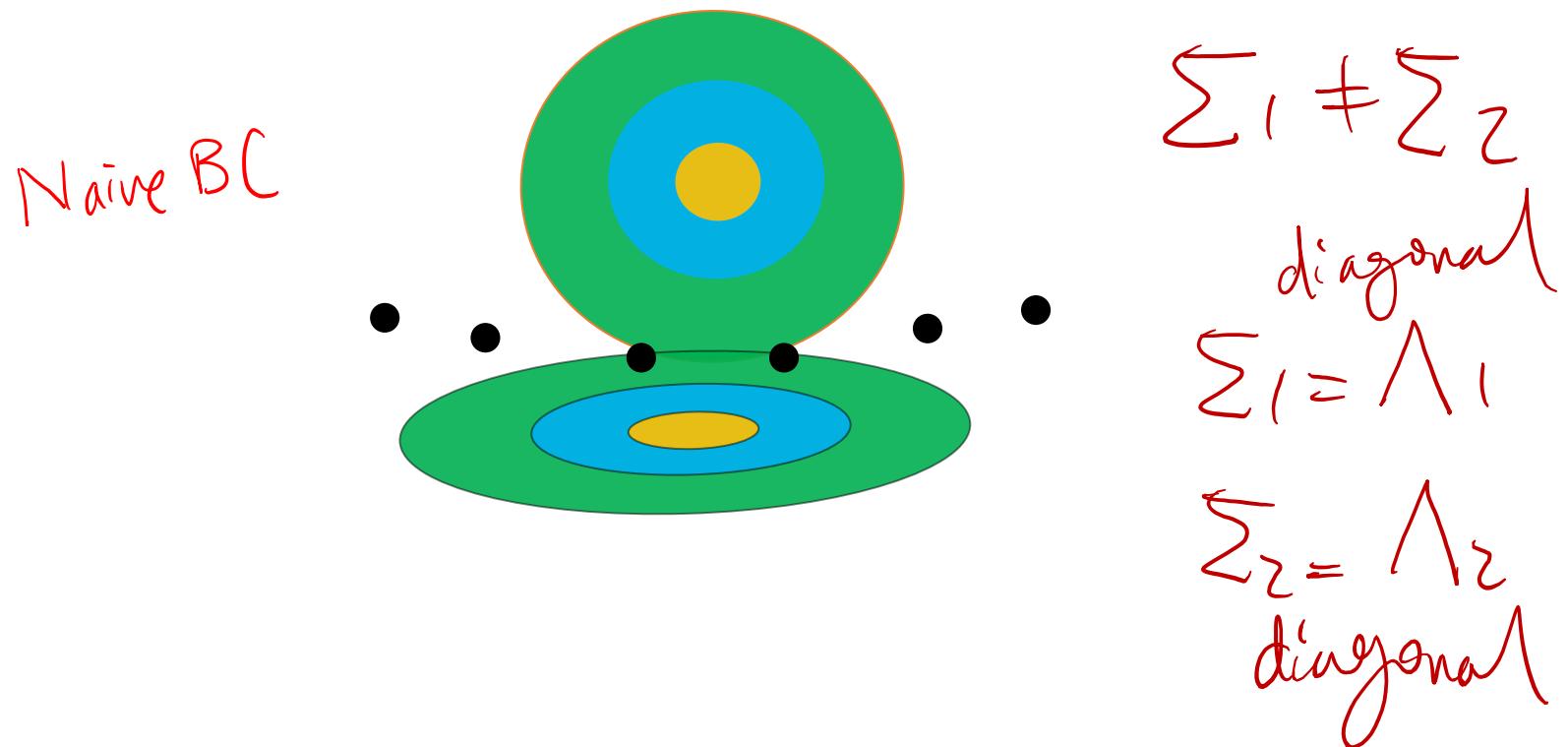
(3) Regularized Discriminant Analysis

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

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

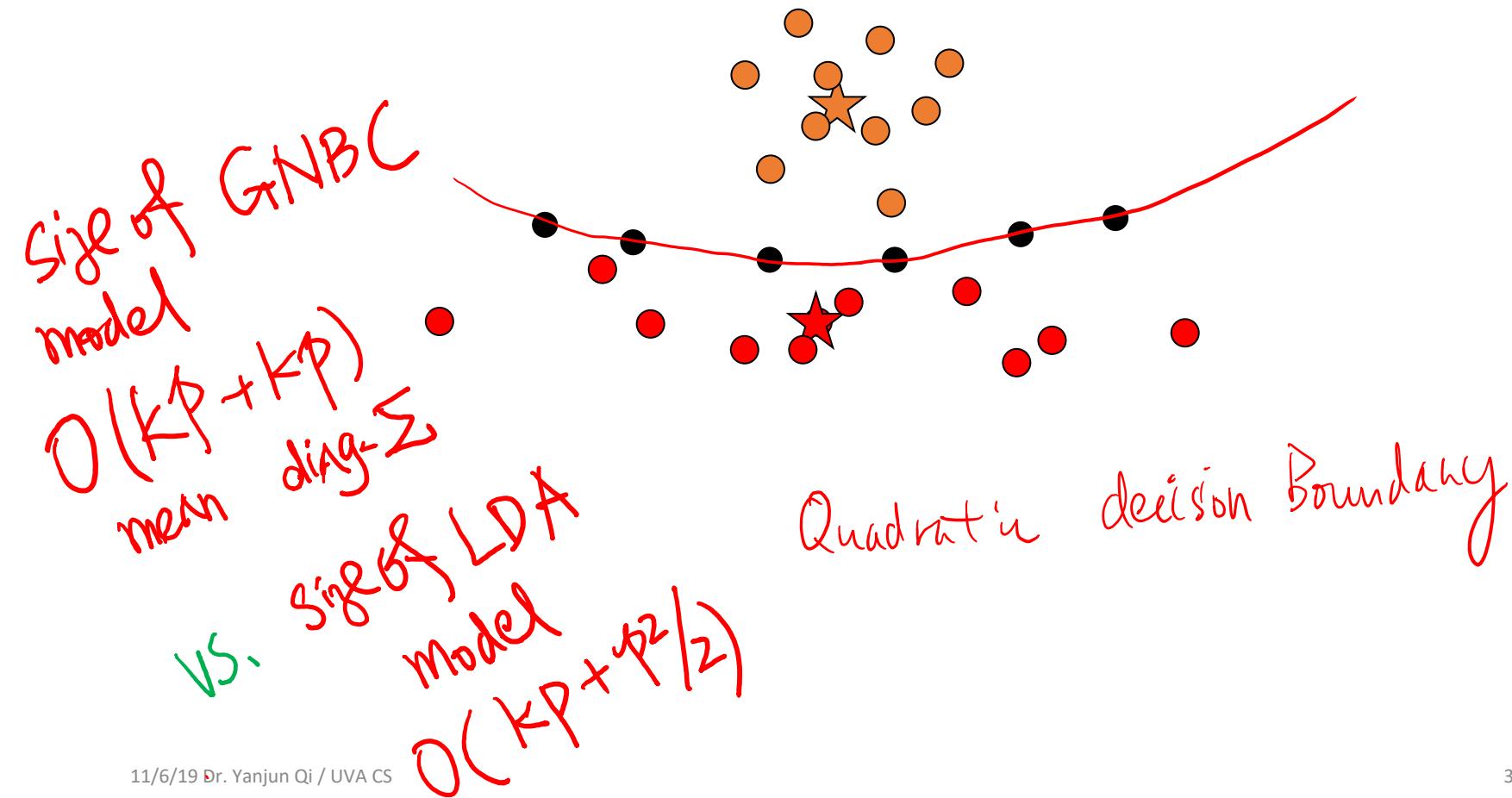
- ▶ The quadratic discriminant function $\delta_k(x)$ is defined using the shrunken covariance matrices $\hat{\Sigma}_k(\alpha)$.
- ▶ The parameter α controls the complexity of the model.

An example: Gaussian Bayes Classifier



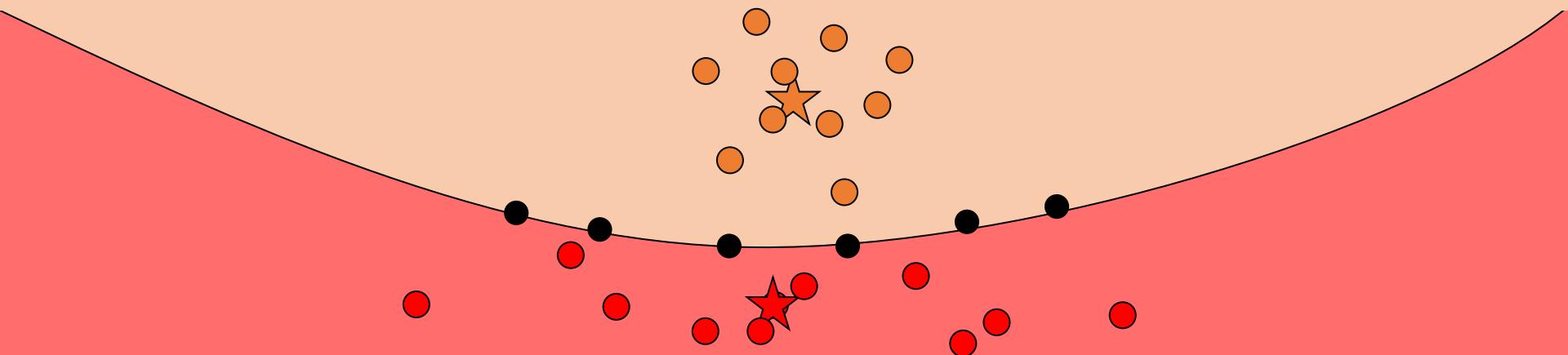
Gaussian Bayes Classifier

Naive
(GNBC)



Gaussian Bayes Classifier

Blue Team



Red Team

Naïve Gaussian Bayes Classifier is
not a linear classifier!

Today: More Generative Bayes Classifiers

- ✓ Generative Bayes Classifier
- ✓ Naïve Bayes Classifier
- ✓ Gaussian Bayes Classifiers
 - Gaussian distribution
 - Naïve Gaussian BC
 - Not-naïve Gaussian BC → LDA, QDA
- ✓ Discriminative vs. Generative

Discriminative vs. Generative

Generative approach

- Model the joint distribution $p(X, C)$ using
 $p(X | C = c_k)$ and $p(C = c_k)$

Discriminative approach

- Model the conditional distribution $p(c | X)$ directly

Class prior

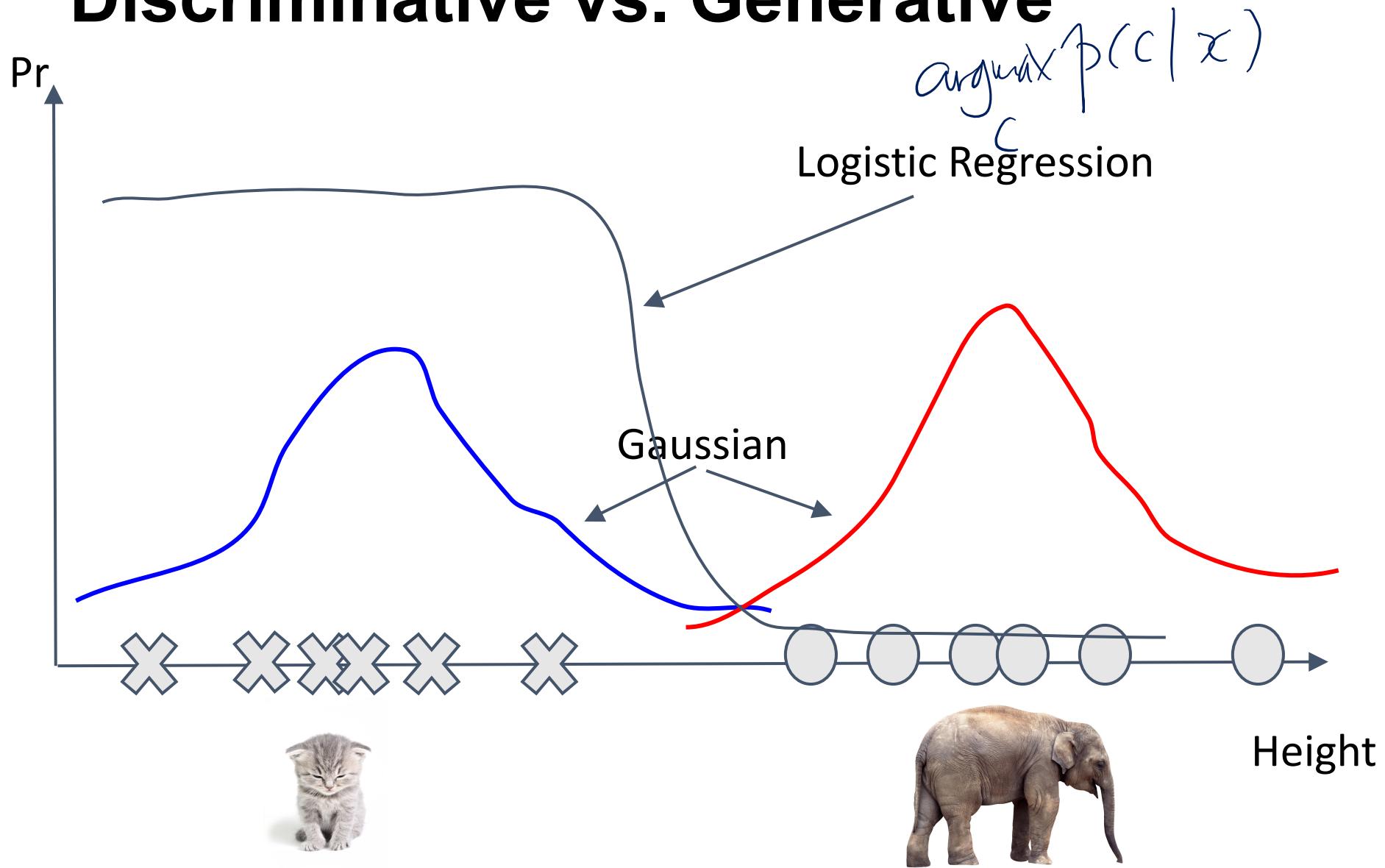
$$P(c=1 | x) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 * x)}}$$

e.g.,

Discriminative vs. Generative

- Empirically, **generative** classifiers approach their asymptotic error faster than discriminative ones
 - Good for small training set
 - Handle missing data well (EM)
- Empirically, **discriminative** classifiers have lower asymptotic error than generative ones
 - Good for larger training set

Discriminative vs. Generative



LDA vs. Logistic Regression

• LDA (Generative model)

- Assumes Gaussian class-conditional densities and a common covariance
- Model parameters are estimated by maximizing the [full log likelihood,]
parameters for each class are estimated independently of other classes,
 $K p + \frac{p(p+1)}{2} + (K - 1)$ parameters
- Makes use of marginal density information $\Pr(x)$
- Easier to train, low variance, more efficient if model is correct
- Higher asymptotic error, but converges faster

$$p(x_{p+1} | c_i)$$

$$\Rightarrow \text{mean } Kp + p^2 \underset{\text{Conv}}{\text{Conv}}$$

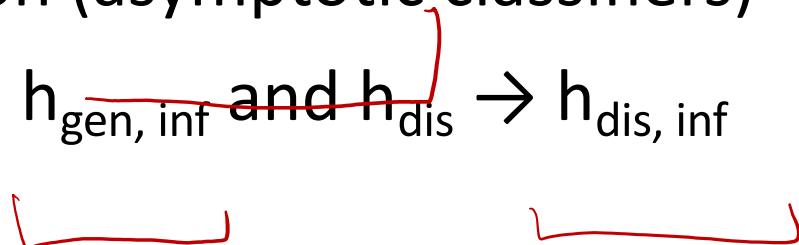
• Logistic Regression (Discriminative model)

- Assumes class-conditional densities are members of the (same) exponential family distribution $p(c_i | x)$
- Model parameters are estimated by maximizing the [conditional log likelihood]
simultaneous consideration of all other classes, $(K - 1)(p + 1)$ parameters
- Ignores marginal density information $\Pr(x)$
- Harder to train, robust to uncertainty about the data generation process
- Lower asymptotic error, but converges more slowly

$$\Rightarrow (K-1)(p+1)$$

Discriminative vs. Generative

- Definitions
 - h_{gen} and h_{dis} : generative and discriminative classifiers
 - $h_{gen, inf}$ and $h_{dis, inf}$: same classifiers but trained on the entire population (asymptotic classifiers)
 - $n \rightarrow \text{infinity}$, $h_{gen} \rightarrow h_{\underline{\text{gen, inf}}}$ and $h_{dis} \rightarrow h_{\underline{\text{dis, inf}}}$



Ng, Jordan., "On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes." *Advances in neural information processing systems* 14 (2002): 841.

Discriminative vs. Generative

Proposition 1:

$$\epsilon(h_{dis,inf}) \leq \epsilon(h_{gen,inf})$$

Proposition 2:

ϵ : ~~assymptotic error~~

Proposition 1 states that asymptotically, the error of the discriminative logistic regression is smaller than that of the generative naive Bayes. This is easily shown

- ϵ : generalization error

Logistic Regression vs. NBC

Discriminative classifier (Logistic Regression)

- Smaller asymptotic error
- Slow convergence $\sim O(p)$

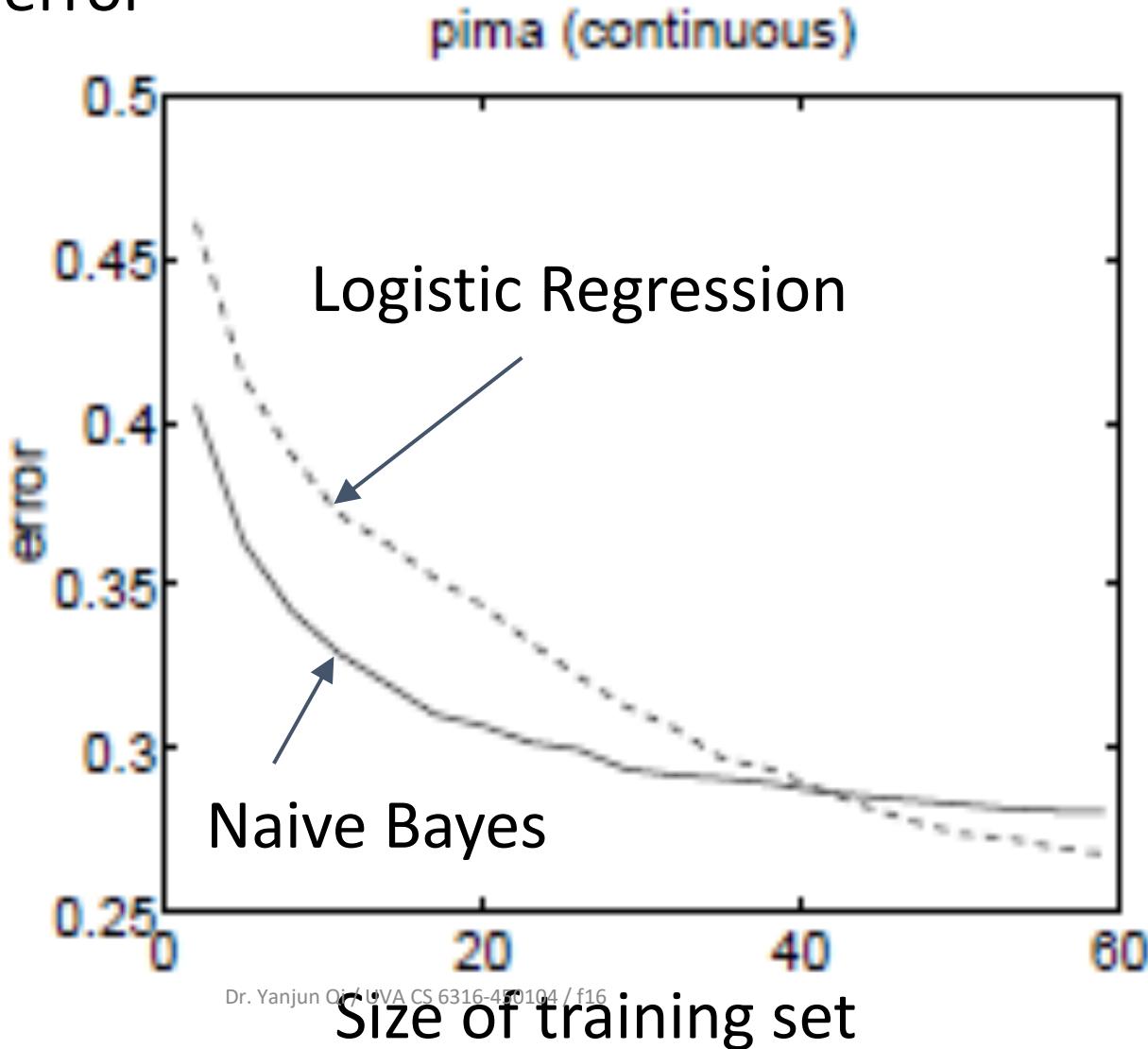
Generative classifier (Naive Bayes)

- Larger asymptotic error
- Can handle missing data (EM)
- Fast convergence $\sim O(\lg(p))$

In numerical analysis, the speed at which a convergent sequence approaches its limit is called the rate of convergence.

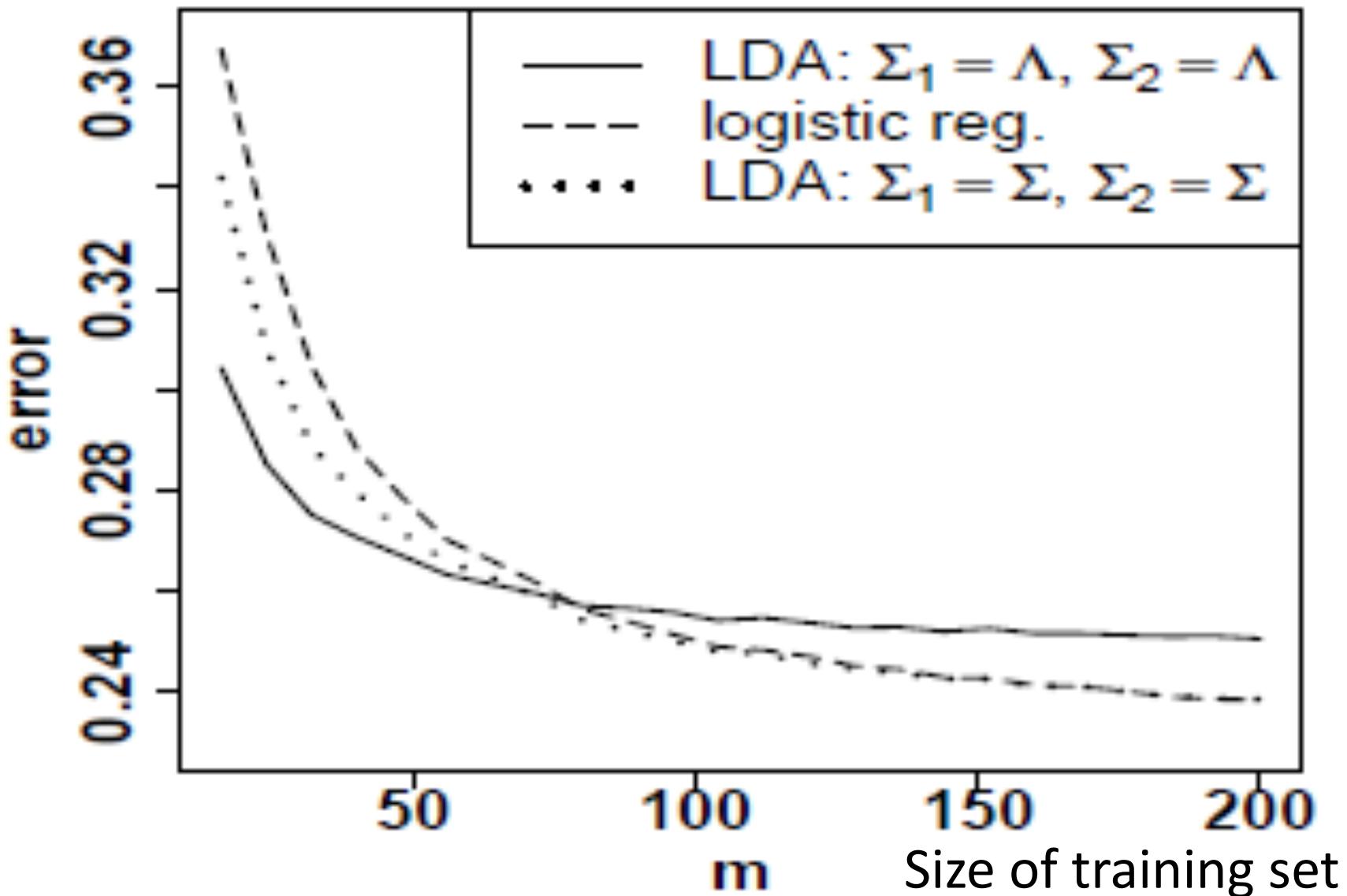
Ng, Jordan,. "On discriminative vs. generative classifiers: A comparison of logistic regression and naive bayes." *Advances in neural information processing systems* 14 (2002): 841.

generalization error



generalization error

pima



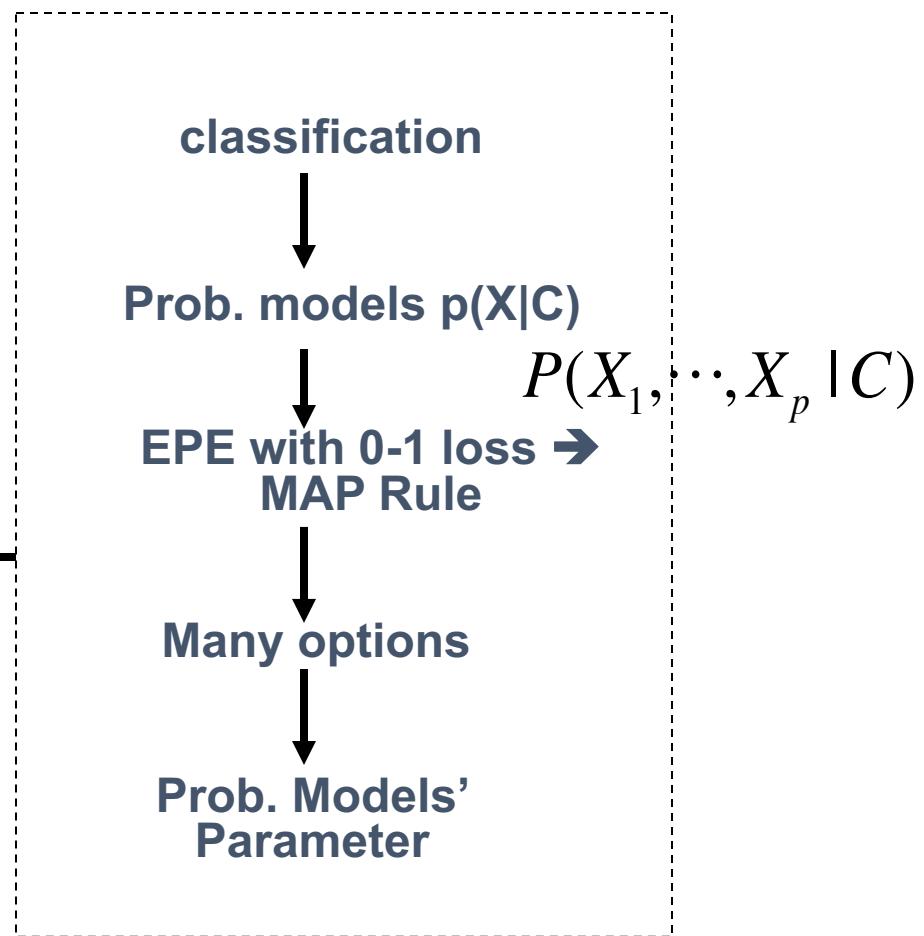
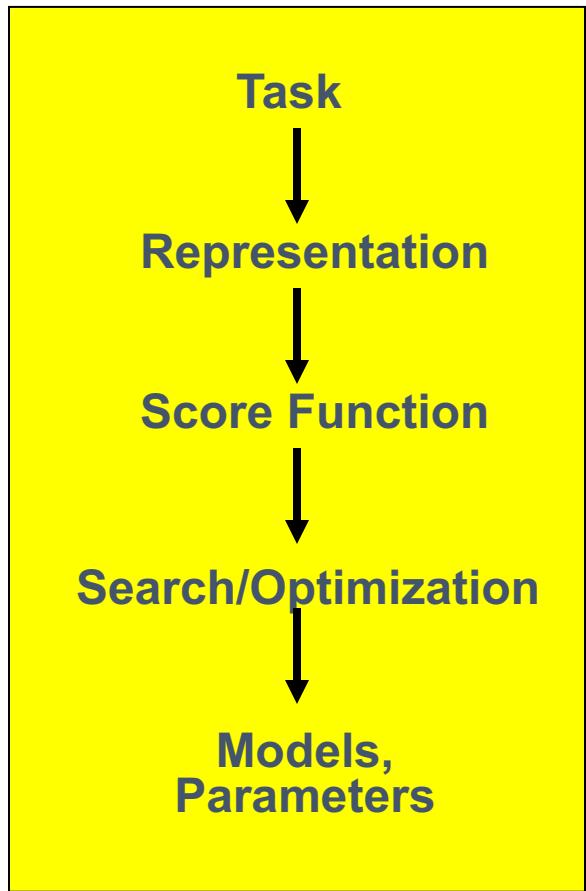
Xue, Jing-Hao, and D. Michael Titterington. "Comment on "On discriminative vs. generative classifiers: A comparison of logistic regression and naive Bayes"." *Neural processing letters* 28.3 (2008): 169-187.

Discriminative vs. Generative

- Empirically, **generative** classifiers approach their asymptotic error faster than discriminative ones
 - Good for small training set
 - Handle missing data well (EM)
- Empirically, **discriminative** classifiers have lower asymptotic error than generative ones
 - Good for larger training set

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

Generative Bayes Classifier



Bernoulli Naïve

$$p(W_i = \text{true} | c_k) = p_{i,k}$$

Gaussian Naïve

$$\hat{P}(X_j | C = c_k) = \frac{1}{\sqrt{2\pi}\sigma_{jk}} \exp\left(-\frac{(X_j - \mu_{jk})^2}{2\sigma_{jk}^2}\right)$$

Multinomial

$$P(W_1 = n_1, \dots, W_v = n_v | c_k) = \frac{N!}{n_{1k}! n_{2k}! \dots n_{vk}!} \theta_{1k}^{n_{1k}} \theta_{2k}^{n_{2k}} \dots \theta_{vk}^{n_{vk}}$$

References

- Prof. Tan, Steinbach, Kumar's "Introduction to Data Mining" slide
- Prof. Andrew Moore's slides
- Prof. Eric Xing's slides
- Prof. Ke Chen NB slides
- Hastie, Trevor, et al. *The elements of statistical learning*. Vol. 2. No. 1. New York: Springer, 2009.