



# Probabilistic classifiers

CE-477: Machine Learning - CS-828: Theory of Machine Learning  
Sharif University of Technology  
Fall 2024

Fatemeh Seyyedsalehi

# Topics

- ▶ Probabilistic approach
  - ▶ Bayes decision theory
  - ▶ Generative models
    - ▶ Gaussian Bayes classifier
    - ▶ Naïve Bayes
  - ▶ Discriminative models
    - ▶ Logistic regression

# Classification problem: probabilistic view

- ▶ Each feature as a random variable
- ▶ Class label also as a random variable
- ▶ We observe the feature values for a random sample and we intend to find its class label
  - ▶ Evidence: feature vector  $x$
  - ▶ Query: class label

# Definitions

- ▶ Posterior probability:  $p(\mathcal{C}_k|\mathbf{x})$
- ▶ Likelihood or class conditional probability:  $p(\mathbf{x}|\mathcal{C}_k)$
- ▶ Prior probability:  $p(\mathcal{C}_k)$

$p(\mathbf{x})$ : pdf of feature vector  $\mathbf{x}$  ( $p(\mathbf{x}) = \sum_{k=1}^K p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)$ )

$p(\mathbf{x}|\mathcal{C}_k)$ : pdf of feature vector  $\mathbf{x}$  for samples of class  $\mathcal{C}_k$

$p(\mathcal{C}_k)$ : probability of the label be  $\mathcal{C}_k$

# Bayes decision rule

$K = 2$

If  $P(\mathcal{C}_1|\mathbf{x}) > P(\mathcal{C}_2|\mathbf{x})$  decide  $\mathcal{C}_1$   
otherwise decide  $\mathcal{C}_2$

$$p(error|\mathbf{x}) = \begin{cases} p(\mathcal{C}_2|\mathbf{x}) & \text{if we decide } \mathcal{C}_1 \\ P(\mathcal{C}_1|\mathbf{x}) & \text{if we decide } \mathcal{C}_2 \end{cases}$$

- If we use Bayes decision rule:

$$P(error|\mathbf{x}) = \min\{P(\mathcal{C}_1|\mathbf{x}), P(\mathcal{C}_2|\mathbf{x})\}$$

Using Bayes rule, for each  $\mathbf{x}$ ,  $P(error|\mathbf{x})$  is as small as possible and thus this rule minimizes the probability of error

# Optimal classifier

- ▶ The optimal decision is the one that minimizes the expected number of mistakes
- ▶ We show that Bayes classifier is an optimal classifier

# Bayes decision rule

## Minimizing misclassification rate

► Decision regions:  $\mathcal{R}_k = \{\mathbf{x} | \alpha(\mathbf{x}) = k\}$

$K = 2$

► All points in  $\mathcal{R}_k$  are assigned to class  $\mathcal{C}_k$

$$p(\text{error}) = E_{\mathbf{x}, y}[I(\alpha(\mathbf{x}) \neq y)]$$

$$= p(\mathbf{x} \in \mathcal{R}_1, \mathcal{C}_2) + p(\mathbf{x} \in \mathcal{R}_2, \mathcal{C}_1)$$

$$= \int_{\mathcal{R}_1} p(\mathbf{x}, \mathcal{C}_2) d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathbf{x}, \mathcal{C}_1) d\mathbf{x}$$

$$= \int_{\mathcal{R}_1} p(\mathcal{C}_2 | \mathbf{x}) p(\mathbf{x}) d\mathbf{x} + \int_{\mathcal{R}_2} p(\mathcal{C}_1 | \mathbf{x}) p(\mathbf{x}) d\mathbf{x}$$

Choose class with highest  $p(\mathcal{C}_k | \mathbf{x})$  as  $\alpha(\mathbf{x})$

# Bayes minimum error

- ▶ Bayes minimum error classifier:

$$\min_{\alpha(\cdot)} E_{\mathbf{x},y}[I(\alpha(\mathbf{x}) \neq y)] \quad \text{Zero-one loss}$$

- ▶ If we know the probabilities in advance then the above optimization problem will be solved easily.
  - ▶  $\alpha(\mathbf{x}) = \operatorname{argmax}_y p(y|\mathbf{x})$
- ▶ In practice, we can estimate  $p(y|\mathbf{x})$  based on a set of training samples  $\mathcal{D}$



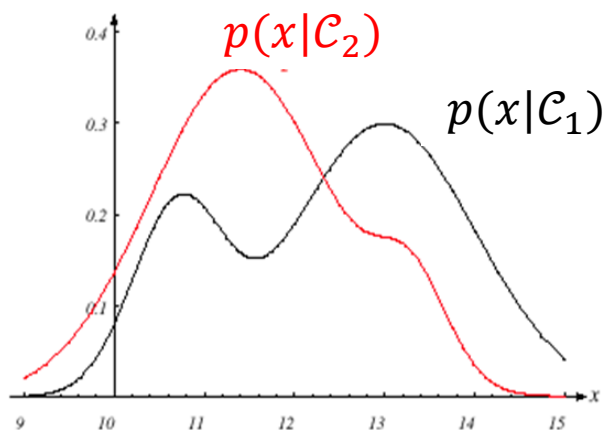
# Bayes theorem

- ▶ Bayes' theorem
- $$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{p(\mathbf{x})}$$
- Diagram illustrating the components of Bayes' theorem:
- Posterior:  $p(\mathcal{C}_k|\mathbf{x})$
  - Likelihood:  $p(\mathbf{x}|\mathcal{C}_k)$
  - Prior:  $p(\mathcal{C}_k)$
- ▶ Posterior probability:  $p(\mathcal{C}_k|\mathbf{x})$
  - ▶ Likelihood or class conditional probability:  $p(\mathbf{x}|\mathcal{C}_k)$
  - ▶ Prior probability:  $p(\mathcal{C}_k)$

$p(\mathbf{x})$ : pdf of feature vector  $\mathbf{x}$  ( $p(\mathbf{x}) = \sum_{k=1}^K p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)$ )  
 $p(\mathbf{x}|\mathcal{C}_k)$ : pdf of feature vector  $\mathbf{x}$  for samples of class  $\mathcal{C}_k$   
 $p(\mathcal{C}_k)$ : probability of the label be  $\mathcal{C}_k$

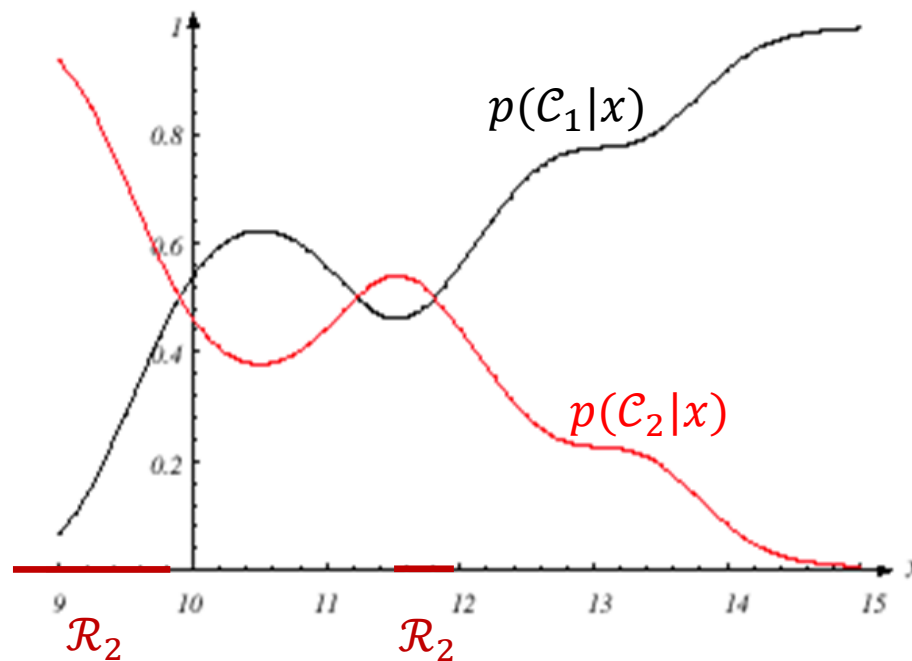
# Bayes decision rule: example

- Bayes decision: Choose the class with highest  $p(\mathcal{C}_k|\mathbf{x})$



$$p(\mathcal{C}_1) = \frac{2}{3}$$



$$p(\mathcal{C}_2) = \frac{1}{3}$$



$$p(\mathcal{C}_k|x) = \frac{p(x|\mathcal{C}_k)p(\mathcal{C}_k)}{p(x)}$$

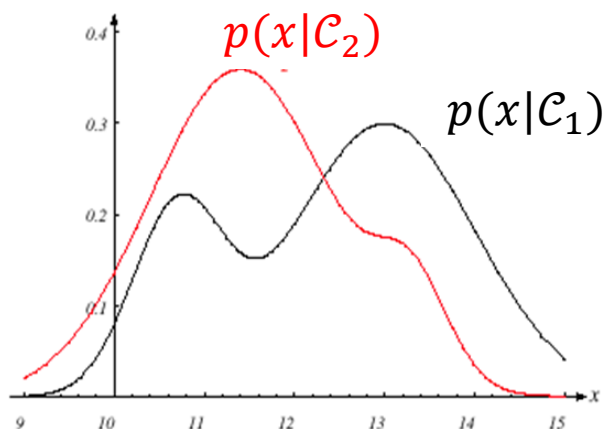
$$p(x) = p(\mathcal{C}_1)p(x|\mathcal{C}_1) + p(\mathcal{C}_2)p(x|\mathcal{C}_2)$$

# Bayesian decision rule

- ▶ If  $P(\mathcal{C}_1|\mathbf{x}) > P(\mathcal{C}_2|\mathbf{x})$  decide  $\mathcal{C}_1$   
otherwise decide  $\mathcal{C}_2$  Equivalent
- ▶ If  $\frac{p(\mathbf{x}|\mathcal{C}_1)P(\mathcal{C}_1)}{p(\mathbf{x})} > \frac{p(\mathbf{x}|\mathcal{C}_2)P(\mathcal{C}_2)}{p(\mathbf{x})}$  decide  $\mathcal{C}_1$   
otherwise decide  $\mathcal{C}_2$  Equivalent
- ▶ If  $p(\mathbf{x}|\mathcal{C}_1)P(\mathcal{C}_1) > p(\mathbf{x}|\mathcal{C}_2)P(\mathcal{C}_2)$  decide  $\mathcal{C}_1$   
otherwise decide  $\mathcal{C}_2$

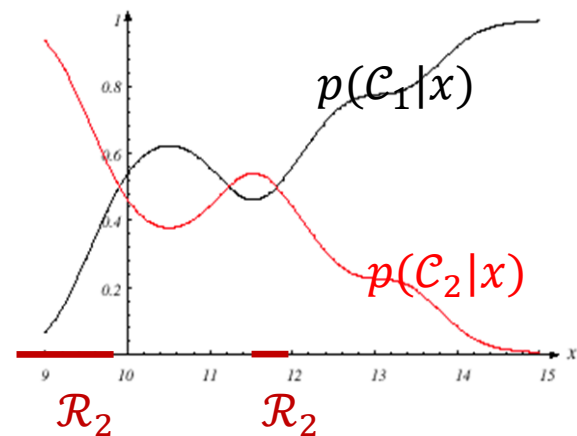
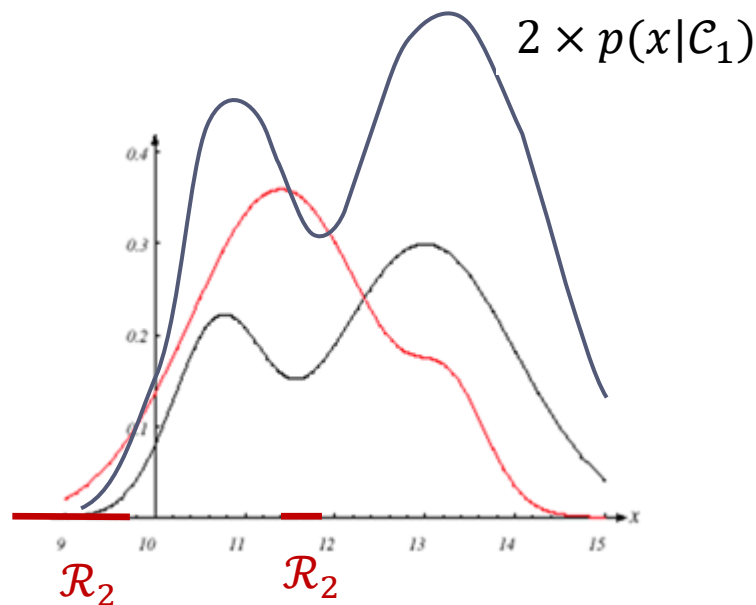
# Bayes decision rule: example

- Bayes decision: Choose the class with highest  $p(\mathcal{C}_k|\mathbf{x})$



$$p(\mathcal{C}_1) = \frac{2}{3}$$

$$p(\mathcal{C}_2) = \frac{1}{3}$$

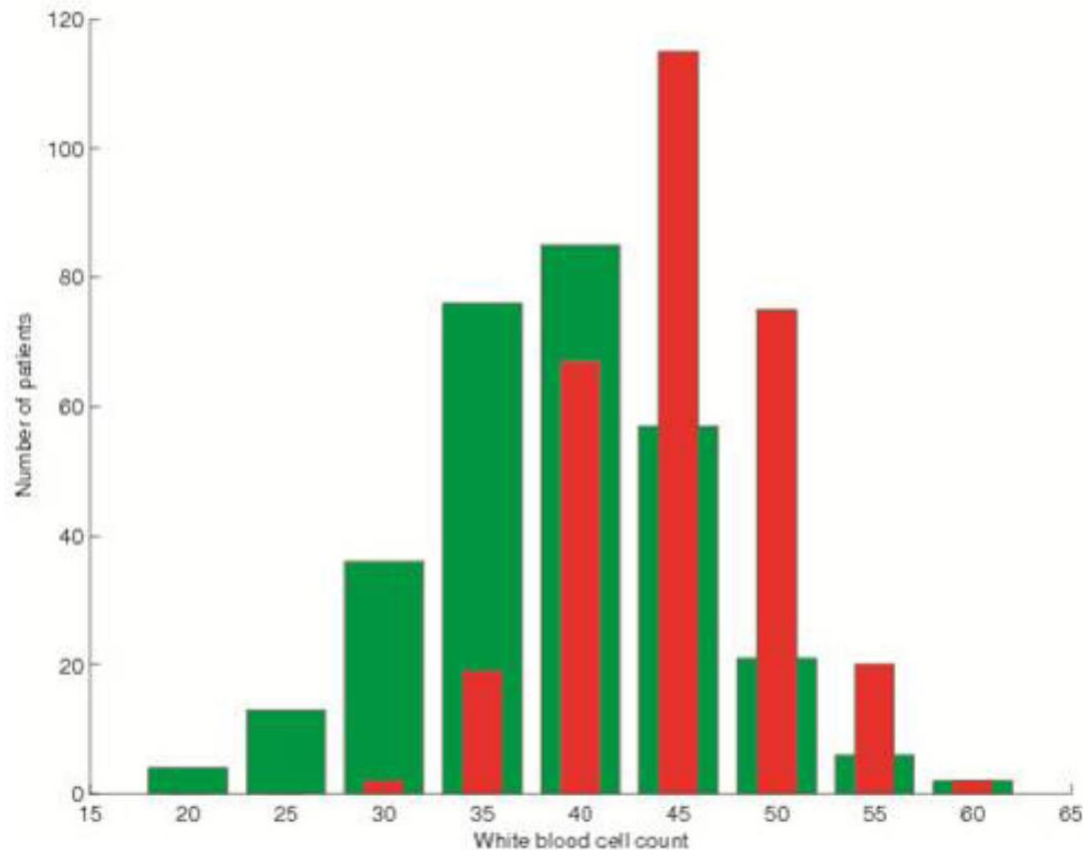


# Bayes Classifier

- ▶ Simple Bayes classifier: estimate posterior probability of each class
- ▶ What should the decision criterion be?
  - ▶ Choose class with highest  $p(\mathcal{C}_k|\mathbf{x})$
- ▶ The optimal decision is the one that minimizes the expected number of mistakes

# Diabetes example

- ▶ white blood cell count

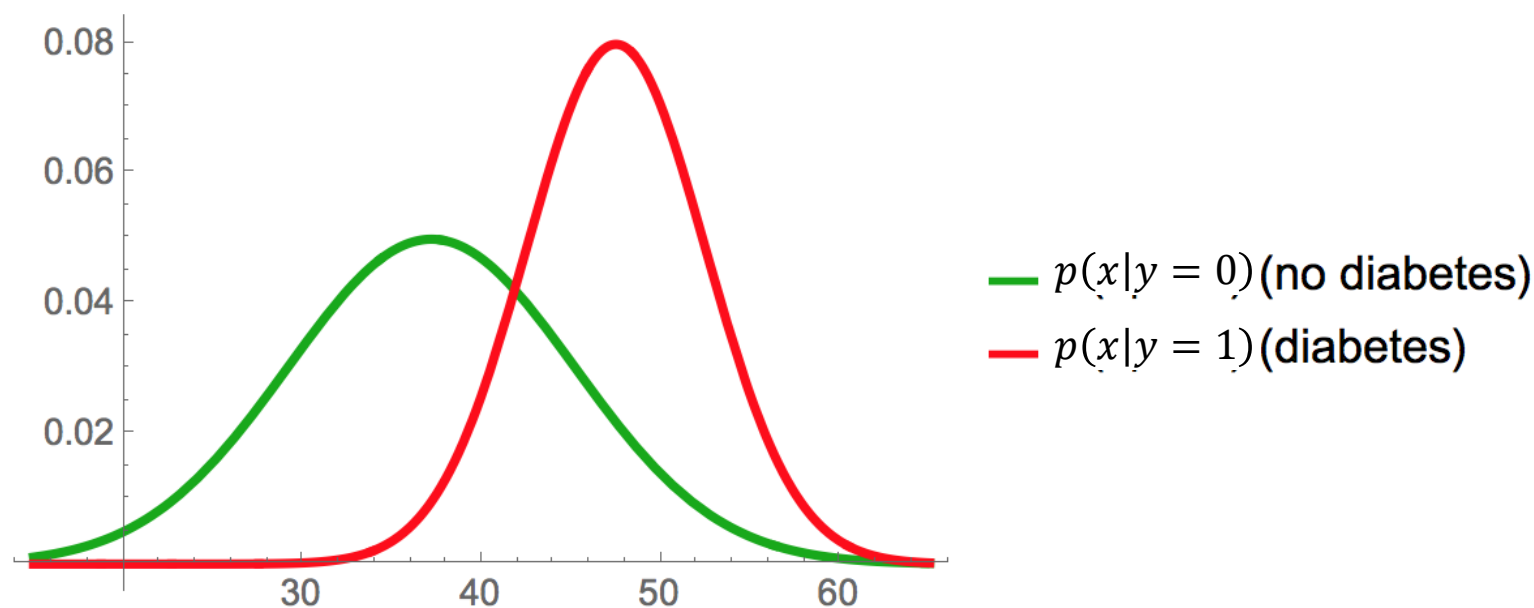


# Diabetes example

- ▶ Doctor has a prior  $p(y = 1) = 0.2$ 
  - ▶ Prior: In the absence of any observation, what do I know about the probability of the classes?
- ▶ A patient comes in with white blood cell count  $x$
- ▶ Does the patient have diabetes  $p(y = 1|x)$ ?
  - ▶ given a new observation, we still need to compute the posterior

# Diabetes example

$$p(x = 40|y = 0)P(y = 0) >? p(x = 40|y = 1)P(y = 1)$$



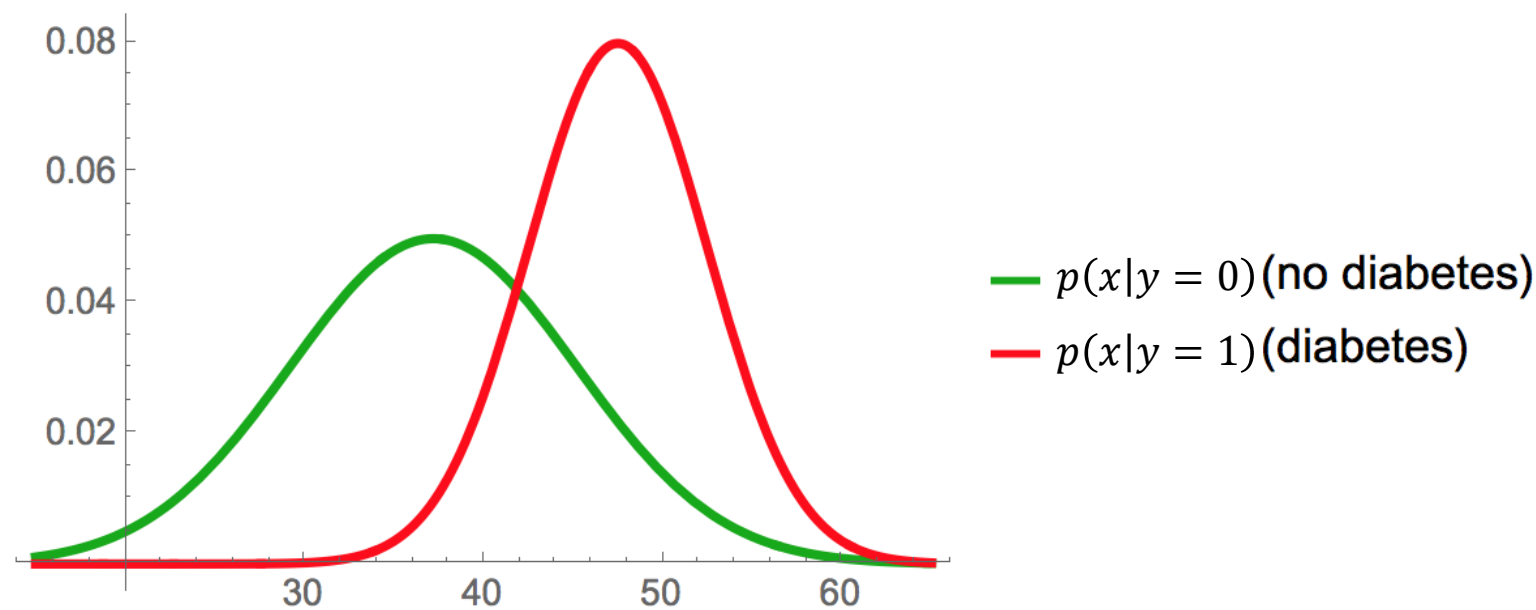


# Estimate probability densities from data

- ▶ If we assume Gaussian distributions for  $p(x|y = 0)$  and  $p(x|y = 1)$
- ▶ Recall that for samples  $\{x^{(1)}, \dots, x^{(N)}\}$ , if we assume a Gaussian distribution, the MLE estimates will be

$$\begin{aligned}\mu &= \frac{1}{N} \sum_{n=1}^N x^{(n)} \\ \sigma^2 &= \frac{1}{N} \sum_{n=1}^N (x^{(n)} - \mu)^2\end{aligned}$$

# Diabetes example



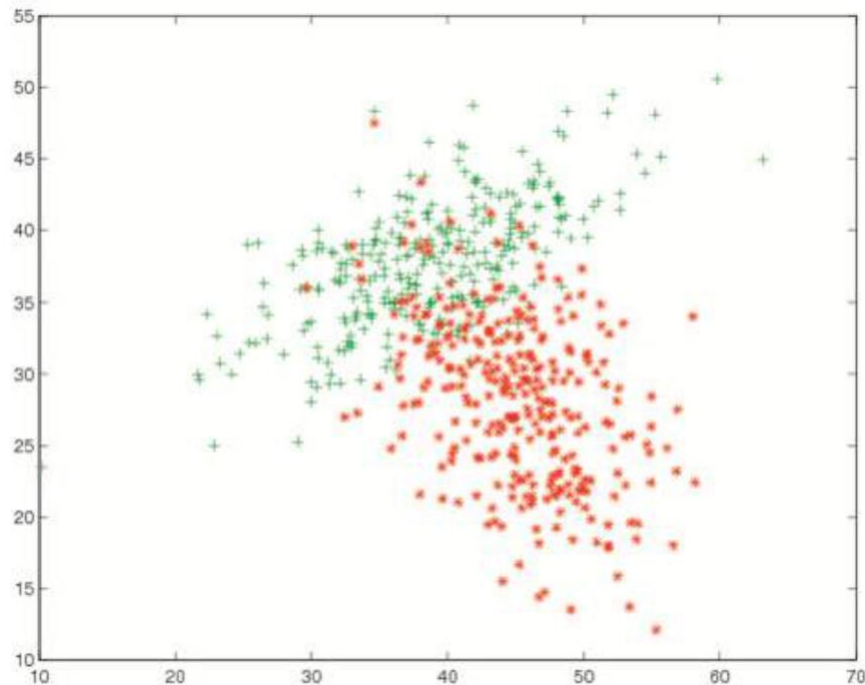
$$p(x|y = 1) = N(\mu_1, \sigma_1^2)$$

$$\mu_1 = \frac{\sum_{n: y^{(n)}=1} x^{(n)}}{\sum_{n: y^{(n)}=1} 1} = \frac{\sum_{n: y^{(n)}=1} x^{(n)}}{N_1}$$

$$\sigma_1^2 = \frac{\sum_{n: y^{(n)}=1} (x^{(n)} - \mu_1)^2}{N_1}$$

# Diabetes example

- Add a second observation: Plasma glucose value



## Generative approach for this example

- ▶ Multivariate Gaussian distributions for  $p(\mathbf{x}|\mathcal{C}_k)$ :

$$p(\mathbf{x}|y = k) = \frac{1}{(2\pi)^{d/2} |\Sigma_k|^{1/2}} \exp\left\{-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k)\right\}$$

$$k = 1, 2$$

- ▶ Prior distribution  $p(y)$ :

- ▶  $p(y = 1) = \pi, \quad p(y = 0) = 1 - \pi$

# MLE for multivariate Gaussian

- ▶ For samples  $\{x^{(1)}, \dots, x^{(N)}\}$ , if we assume a multivariate Gaussian distribution, the MLE estimates will be:

$$\boldsymbol{\mu} = \frac{\sum_{n=1}^N \mathbf{x}^{(n)}}{N}$$

$$\boldsymbol{\Sigma} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}^{(n)} - \boldsymbol{\mu})(\mathbf{x}^{(n)} - \boldsymbol{\mu})^T$$

# Generative approach: example

$$y \in \{0,1\}$$

Maximum likelihood estimation ( $D = \{(\mathbf{x}^{(n)}, y^{(n)})\}_{n=1}^N$ ):

- ▶  $\pi = \frac{N_1}{N}$

- ▶  $\boldsymbol{\mu}_1 = \frac{\sum_{n=1}^N y^{(n)} \mathbf{x}^{(n)}}{N_1}, \boldsymbol{\mu}_2 = \frac{\sum_{n=1}^N (1-y^{(n)}) \mathbf{x}^{(n)}}{N_2}$

$$N_1 = \sum_{n=1}^N y^{(n)}$$

- ▶  $\boldsymbol{\Sigma}_1 = \frac{1}{N_1} \sum_{n=1}^N y^{(n)} (\mathbf{x}^{(n)} - \boldsymbol{\mu})(\mathbf{x}^{(n)} - \boldsymbol{\mu})^T$

$$N_2 = N - N_1$$

- ▶  $\boldsymbol{\Sigma}_2 = \frac{1}{N_2} \sum_{n=1}^N (1-y^{(n)}) (\mathbf{x}^{(n)} - \boldsymbol{\mu})(\mathbf{x}^{(n)} - \boldsymbol{\mu})^T$

## Decision boundary for Gaussian Bayes classifier

$$p(\mathcal{C}_1|\mathbf{x}) = p(\mathcal{C}_2|\mathbf{x})$$

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{p(\mathbf{x})}$$

$$\ln p(\mathcal{C}_1|\mathbf{x}) = \ln p(\mathcal{C}_2|\mathbf{x})$$

$$\begin{aligned} \ln p(\mathbf{x}|\mathcal{C}_1) + \ln p(\mathcal{C}_1) - \ln p(\mathbf{x}) \\ = \ln p(\mathbf{x}|\mathcal{C}_2) + \ln p(\mathcal{C}_2) - \ln p(\mathbf{x}) \end{aligned}$$

# Decision boundary for Gaussian Bayes classifier

$$p(\mathcal{C}_1|\mathbf{x}) = p(\mathcal{C}_2|\mathbf{x})$$

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{p(\mathbf{x}|\mathcal{C}_k)p(\mathcal{C}_k)}{p(\mathbf{x})}$$

$$\ln p(\mathcal{C}_1|\mathbf{x}) = \ln p(\mathcal{C}_2|\mathbf{x})$$

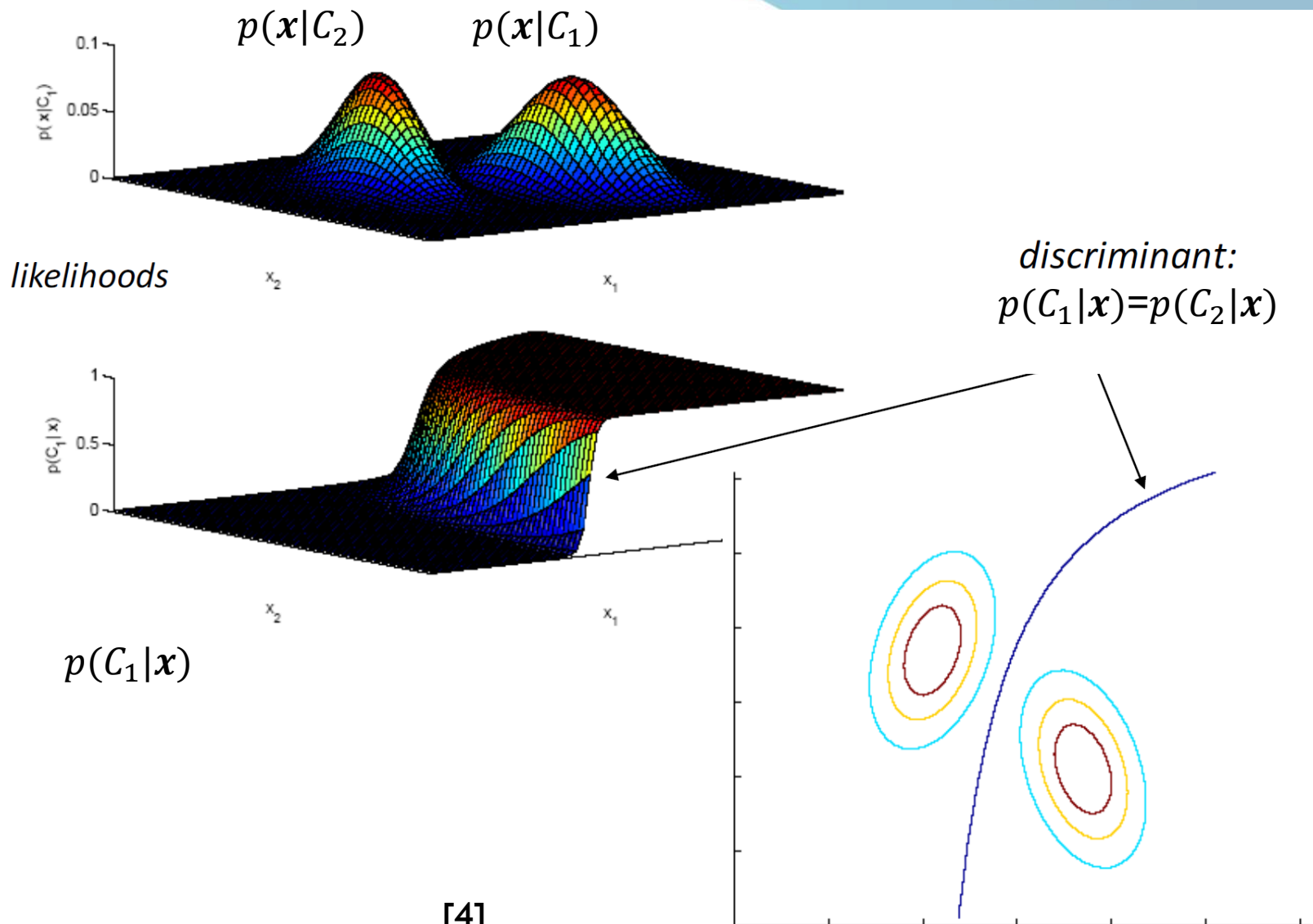
$$\begin{aligned} \ln p(\mathbf{x}|\mathcal{C}_1) + \ln p(\mathcal{C}_1) - \ln p(\mathbf{x}) \\ = \ln p(\mathbf{x}|\mathcal{C}_2) + \ln p(\mathcal{C}_2) - \ln p(\mathbf{x}) \end{aligned}$$

$$\ln p(\mathbf{x}|\mathcal{C}_1) + \ln p(\mathcal{C}_1) = \ln p(\mathbf{x}|\mathcal{C}_2) + \ln p(\mathcal{C}_2)$$

$$\begin{aligned} \ln p(\mathbf{x}|\mathcal{C}_k) \\ = -\frac{d}{2} \ln 2\pi - \frac{1}{2} \ln |\boldsymbol{\Sigma}_k| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1} (\mathbf{x} - \boldsymbol{\mu}_k) \end{aligned}$$



# Decision boundary



# Gaussian discriminant analysis

- ▶ When the likelihood probability is distributed according to a multivariate normal distribution.
- ▶ The general model is:

$$C_k \sim \text{Bernouli}(\phi) = \phi^y (1 - \phi)^{(1-y)}$$

$$\begin{aligned} p(\mathbf{x}|y = k) &\sim N(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \\ &= \frac{1}{(2\pi)^{d/2} |\boldsymbol{\Sigma}_k|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)\right\} \end{aligned}$$

# Gaussian discriminant analysis

- ▶ The log-likelihood of the data:

- ▶  $k = 2$

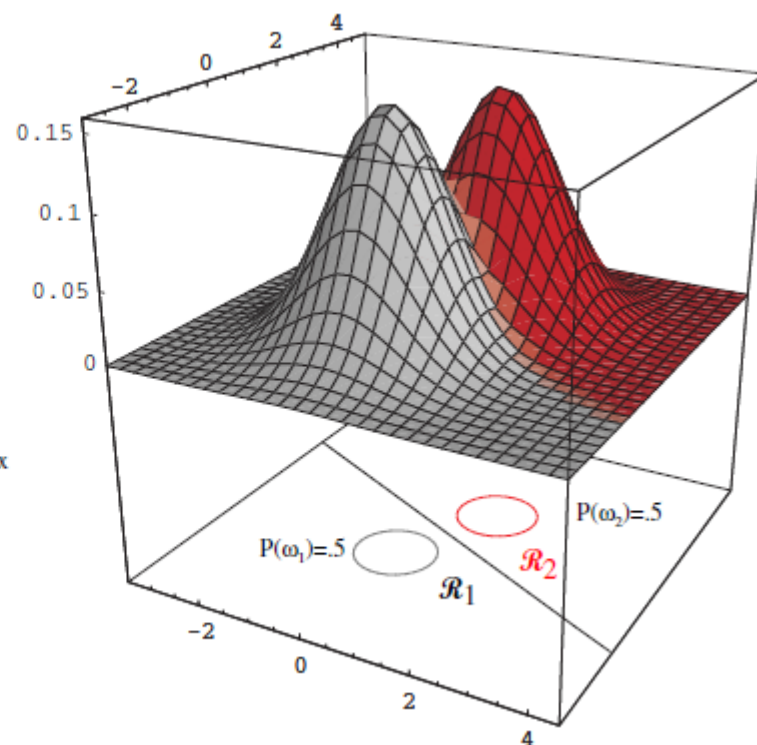
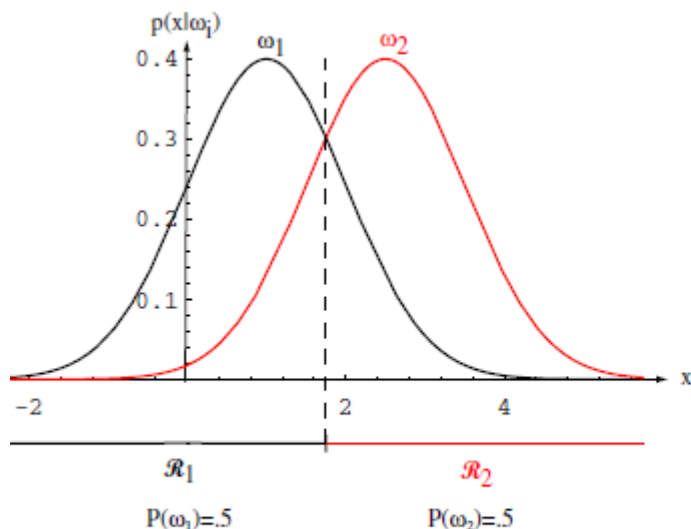
$$\ell(\phi, \mu_0, \mu_1, \Sigma_1, \Sigma_2) = \log \prod_{i=1}^N p(\mathbf{x}^{(i)}, \mathbf{y}^{(i)}; \phi, \mu_0, \mu_1, \Sigma_1, \Sigma_2)$$

- ▶ By maximizing with respect to the parameters

$$\begin{aligned}\phi &= \frac{N_1}{N} \\ \mu_1 &= \frac{\sum_{i=1}^{N_1} y^{(i)} \mathbf{x}^{(i)}}{N_1}, \mu_2 = \frac{\sum_{i=1}^{N_2} (1 - y^{(i)}) \mathbf{x}^{(i)}}{N_2} \\ \Sigma_1 &= \frac{1}{N_1} \sum_{i=1}^{N_1} y^{(i)} (\mathbf{x}^{(i)} - \mu) (\mathbf{x}^{(i)} - \mu)^T \\ \Sigma_2 &= \frac{1}{N_2} \sum_{i=1}^{N_2} (1 - y^{(i)}) (\mathbf{x}^{(i)} - \mu) (\mathbf{x}^{(i)} - \mu)^T\end{aligned}$$

# Gaussian discriminant analysis

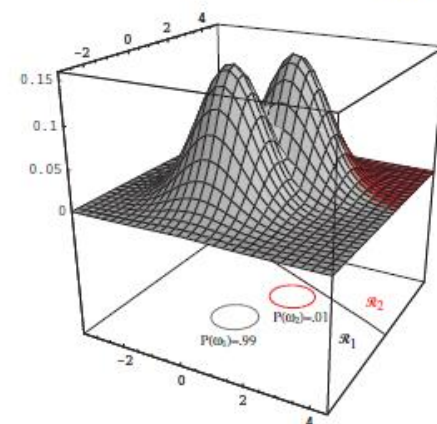
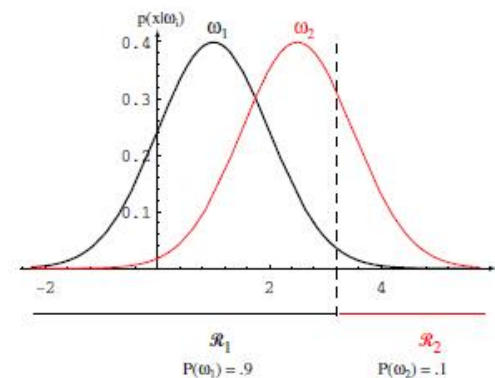
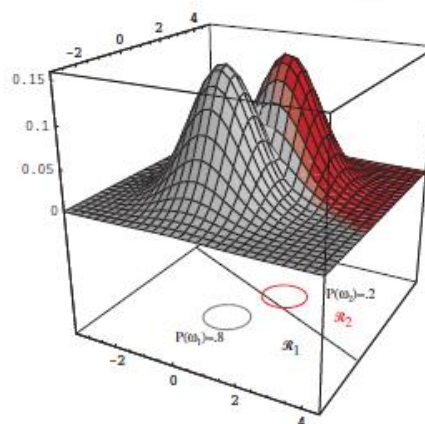
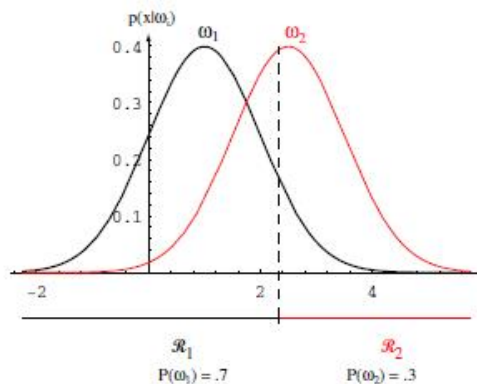
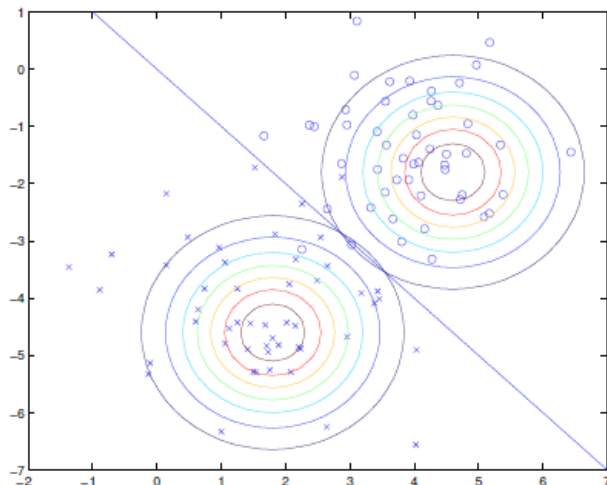
- ▶ A special case:  $\Sigma_k = \sigma^2 I$ 
  - ▶ The decision boundary is a Hyperplane which is orthogonal to the line between the means.



[3]

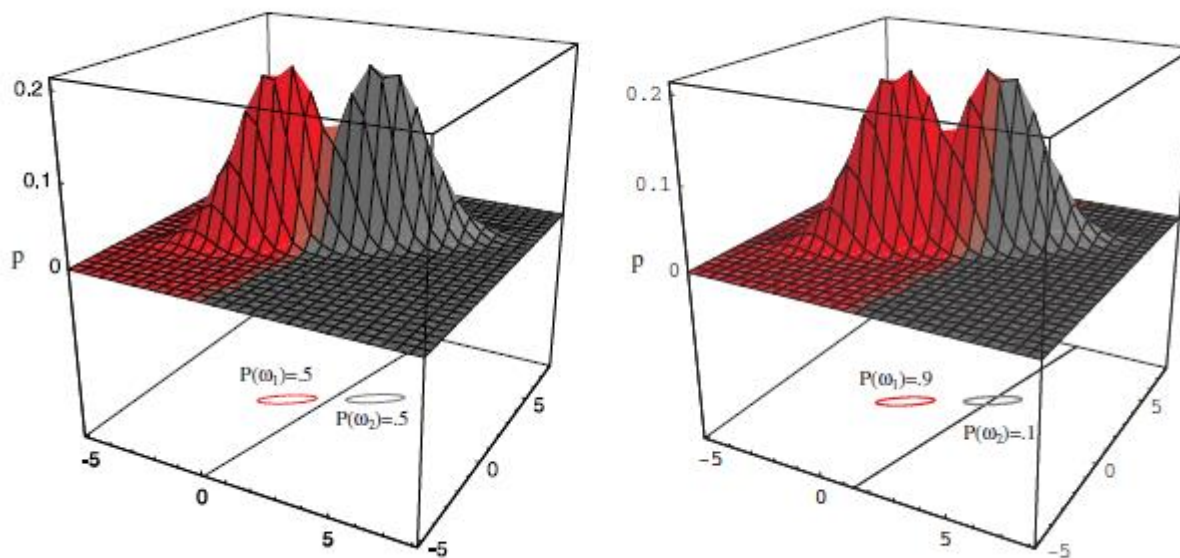
# Gaussian discriminant analysis

- ▶ A special case:  $\Sigma_k = \sigma^2 I$ 
  - ▶ Features are independent random variables

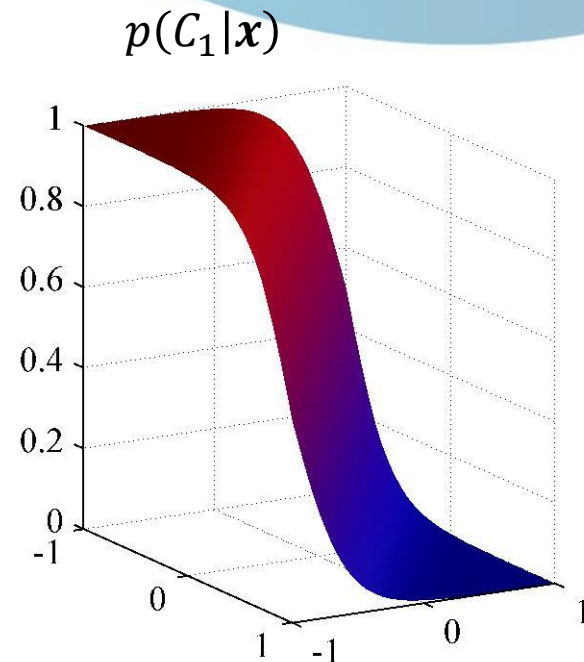
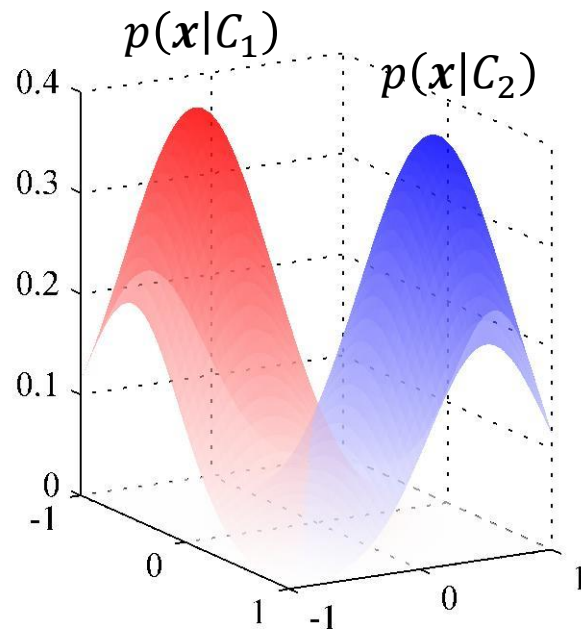


# Gaussian discriminant analysis

- ▶ A special case:  $\Sigma_k = \Sigma$ 
  - ▶ Shared covariance matrix
    - ▶ Equivalent to LDA decision boundary
  - ▶ The decision surface is a Hyperplane, but is not necessarily orthogonal to the line between the means.



# Class conditional densities vs. posterior



[1]

$$\sigma(z) = \frac{1}{1 + \exp(z)}$$

$$p(\mathcal{C}_1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

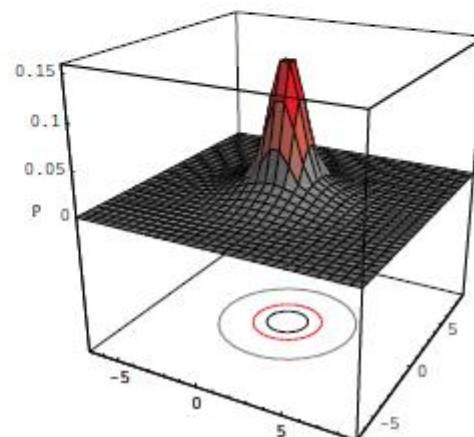
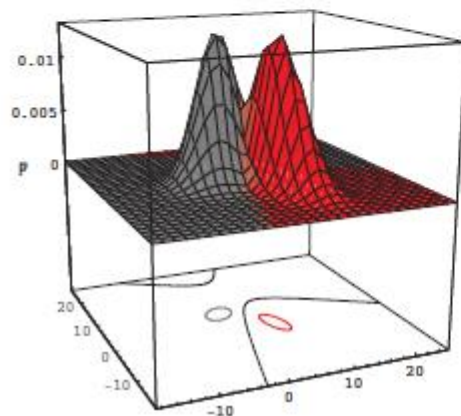
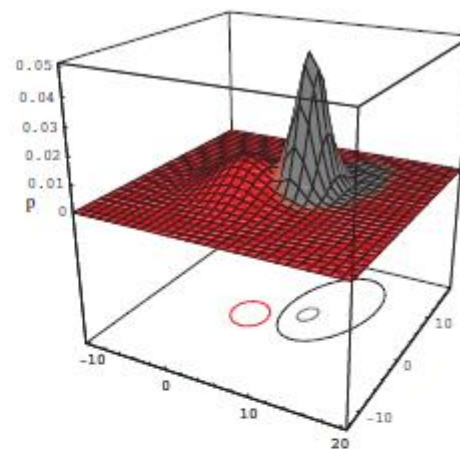
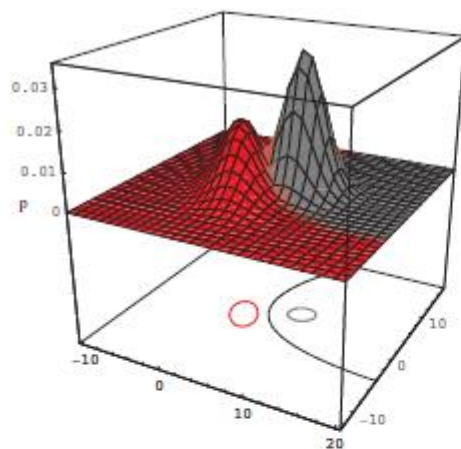
$$\mathbf{w} = \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

$$w_0 = -\frac{1}{2}\boldsymbol{\mu}_1^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_1 + \frac{1}{2}\boldsymbol{\mu}_2^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\mu}_2 + \ln \frac{p(\mathcal{C}_1)}{p(\mathcal{C}_2)}$$



# Gaussian discriminant analysis

- ▶ A special case:  $\Sigma_k = \text{arbitrary}$ 
  - ▶ Decision boundaries are hyperquadrics



[3]



# Naïve Bayes classifier

- ▶ Generative methods
  - ▶ High number of parameters
- ▶ Naïve Bayes assumption: Conditional independence of features

$$p(\mathbf{x}|C_k) = p(x_1|C_k) \times p(x_2|C_k) \times \cdots \times p(x_d|C_k)$$

# Naïve Bayes classifier

- In the decision phase, it finds the label of  $\mathbf{x}$  according to:

$$\operatorname{argmax}_{k=1,\dots,K} p(C_k | \mathbf{x})$$
$$\operatorname{argmax}_{k=1,\dots,K} p(C_k) \prod_{i=1}^N p(x_i | C_k)$$

$$p(\mathbf{x} | C_k) = p(x_1 | C_k) \times p(x_2 | C_k) \times \dots \times p(x_d | C_k)$$
$$p(C_k | \mathbf{x}) \propto p(C_k) \prod_{i=1}^n p(x_i | C_k)$$

# Naïve Bayes: discrete example

▶  $p(h) = 0.3$

$$H = Yes \equiv h$$

$$H = No \equiv \bar{h}$$

▶  $p(d|h) = \frac{1}{3}$

▶  $p(s|h) = \frac{2}{3}$

▶  $p(d|\bar{h}) = \frac{2}{7}$

▶  $p(s|\bar{h}) = \frac{2}{7}$

Diabetes (D)	Smoke (S)	Heart Disease (H)
Y	N	Y
Y	N	N
N	Y	N
N	Y	N
N	N	N
N	Y	Y
N	N	N
N	Y	Y
N	N	N
Y	N	N

▶ Decision on  $\mathbf{x} = [d, \bar{s}]$  (a person that has diabetes but does not smoke):

▶  $p(h|\mathbf{x}) \propto p(h)p(d|h)p(\bar{s}|h) = 1/14$

▶  $p(\bar{h}|\mathbf{x}) \propto p(\bar{h})p(d|\bar{h})p(\bar{s}|\bar{h}) = 1/6$

▶ Thus decide  $H = No$

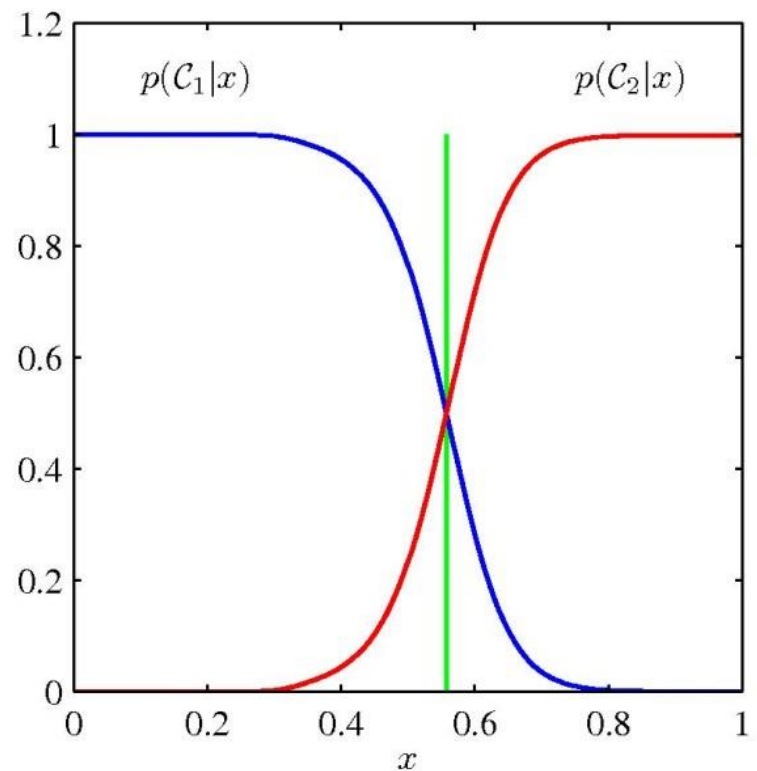
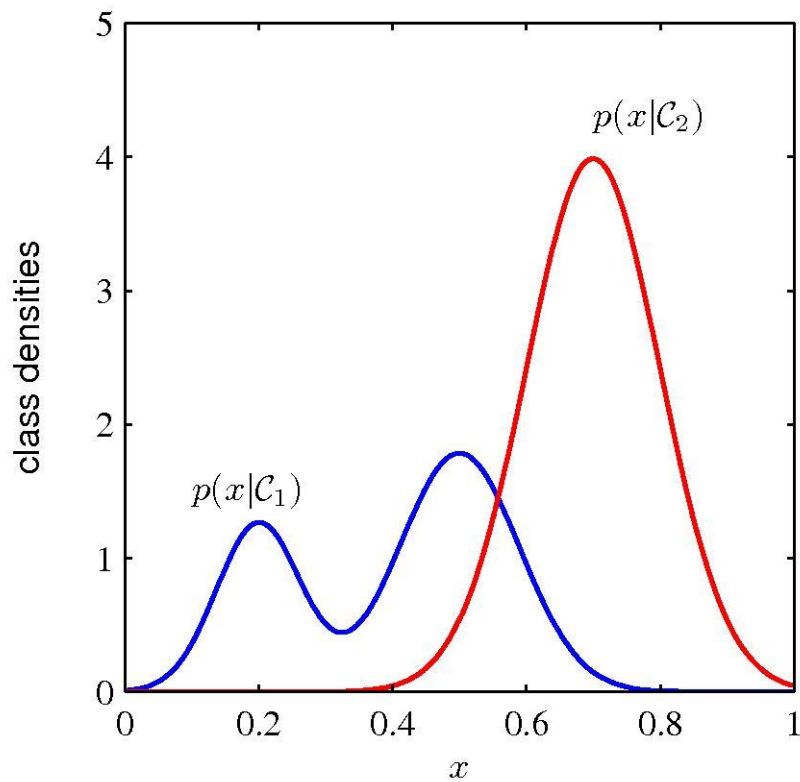
# Naïve Bayes classifier

- ▶ Finds  $d$  univariate distributions  $p(x_1|C_k), \dots, p(x_d|C_k)$  instead of finding one multi-variate distribution  $p(\mathbf{x}|C_k)$ 
  - ▶ Example 1: For Gaussian class-conditional density  $p(\mathbf{x}|C_k)$ , it finds  $d + d$  (mean and sigma parameters on different dimensions) instead of  $d + \frac{d(d+1)}{2}$  parameters
  - ▶ Example 2: For Bernoulli class-conditional density  $p(\mathbf{x}|C_k)$ , it finds  $d$  (mean parameters on different dimensions) instead of  $2^d - 1$  parameters
- ▶ It first estimates the class conditional densities  $p(x_1|C_k), \dots, p(x_d|C_k)$  and the prior probability  $p(C_k)$  for each class ( $k = 1, \dots, K$ ) based on the training set.

# Probabilistic classifiers

- ▶ Probabilistic classification approaches can be divided in two main categories:
  - ▶ **Generative**
    - ▶ Estimate pdf  $p(\mathbf{x}, \mathcal{C}_k)$  for each class  $\mathcal{C}_k$  and then use it to find  $p(\mathcal{C}_k|\mathbf{x})$ 
      - Or alternatively estimate both pdf  $p(\mathbf{x}|\mathcal{C}_k)$  and  $p(\mathcal{C}_k)$  to find  $p(\mathcal{C}_k|\mathbf{x})$
  - ▶ **Discriminative**
    - ▶ Directly estimate  $p(\mathcal{C}_k|\mathbf{x})$  for each class  $\mathcal{C}_k$

# Discriminative vs. generative approach



[1]

# Discriminative approach: logistic regression

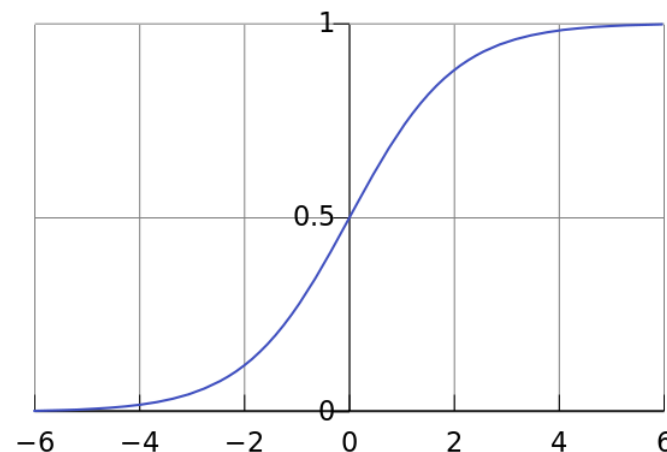
- ▶  $h(\mathbf{x}; \mathbf{w})$  predicts posterior probabilities  $P(y = 1 | \mathbf{x})$   $K = 2$

$$h(\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x})$$

$$\mathbf{x} = [1, x_1, \dots, x_d]$$
$$\mathbf{w} = [w_0, w_1, \dots, w_d]$$

- ▶ Sigmoid (logistic) function

$$\sigma(z) = \frac{1}{1 + e^{-z}}$$



# Logistic regression

- ▶  $h(\mathbf{x}; \mathbf{w})$ : probability that  $y = 1$  given  $\mathbf{x}$  (parameterized by  $\mathbf{w}$ )

- ▶  $0 \leq h(\mathbf{x}; \mathbf{w}) \leq 1$

$$K = 2$$
$$y \in \{0,1\}$$

$$P(y = 1|\mathbf{x}; \mathbf{w}) = h(\mathbf{x}; \mathbf{w})$$

$$P(y = 0|\mathbf{x}; \mathbf{w}) = 1 - h(\mathbf{x}; \mathbf{w})$$

- ▶ Decision surface

$$h(\mathbf{x}; \mathbf{w}) = \sigma(\mathbf{w}^T \mathbf{x}) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x})}} = 0.5$$



# Logistic regression: ML estimation

- ▶ Maximum (conditional) log likelihood:

$$\hat{\mathbf{w}} = \operatorname{argmax}_{\mathbf{w}} \log \prod_{i=1}^n p(y^{(i)} | \mathbf{w}, \mathbf{x}^{(i)})$$

$$p(y^{(i)} | \mathbf{w}, \mathbf{x}^{(i)}) = h(\mathbf{x}^{(i)}; \mathbf{w})^{y^{(i)}} (1 - h(\mathbf{x}^{(i)}; \mathbf{w}))^{(1-y^{(i)})}$$

$$\log p(\mathbf{y} | \mathbf{X}, \mathbf{w}) = \sum_{i=1}^n \left[ y^{(i)} \log(h(\mathbf{x}^{(i)}; \mathbf{w})) + (1 - y^{(i)}) \log(1 - h(\mathbf{x}^{(i)}; \mathbf{w})) \right]$$

# Logistic regression: cost function

$$\hat{\mathbf{w}} = \underset{\mathbf{w}}{\operatorname{argmin}} J(\mathbf{w})$$

$$\begin{aligned} J(\mathbf{w}) &= - \sum_{i=1}^N \log p(y^{(i)} | \mathbf{w}, \mathbf{x}^{(i)}) \\ &= \sum_{i=1}^N -y^{(i)} \log(h(\mathbf{x}^{(i)}; \mathbf{w})) - (1 - y^{(i)}) \log(1 - h(\mathbf{x}^{(i)}; \mathbf{w})) \end{aligned}$$

- ▶ No closed form solution for

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = 0$$

- ▶ However  $J(\mathbf{w})$  is convex.
  - ▶ Global optimum can be found by gradient ascent

# Logistic regression: Gradient descent

- ▶ The gradient descent update

$$\mathbf{w}^{t+1} = \mathbf{w}^t - \eta \nabla_{\mathbf{w}} J(\mathbf{w}^t)$$

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \sum_{i=1}^N (h(\mathbf{x}^{(i)}; \mathbf{w}) - y^{(i)}) \mathbf{x}^{(i)}$$

- ▶ Is it similar to gradient of SSE for linear regression?

$$\nabla_{\mathbf{w}} J(\mathbf{w}) = \sum_{i=1}^N (\mathbf{w}^T \mathbf{x}^{(i)} - y^{(i)}) \mathbf{x}^{(i)}$$

# Posterior probabilities

- ▶ Two-class LR

$$p(\mathcal{C}_1|\mathbf{x}) = \sigma(a(\mathbf{x})) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x})}}$$

- ▶ Multi-class:  $p(\mathcal{C}_k|\mathbf{x})$  can be written as a soft-max function

$$p(\mathcal{C}_k|\mathbf{x}) = \frac{e^{-(\mathbf{w}_k^T \mathbf{x})}}{\sum_{j=1}^K e^{-(\mathbf{w}_j^T \mathbf{x})}}$$

- ▶ To make a prediction:

$$\alpha(\mathbf{x}) = \operatorname{argmax}_{k=1,\dots,K} h_k(\mathbf{x})$$

# Logistic regression: multi-class

## ► The gradient descent update

$$\widehat{\mathbf{W}} = \underset{\mathbf{W}}{\operatorname{argmin}} J(\mathbf{W}) \quad \mathbf{W} = [\mathbf{w}_1 \quad \cdots \quad \mathbf{w}_K]$$

$$\begin{aligned} J(\mathbf{W}) &= -\log \prod_{i=1}^n p(\mathbf{y}^{(i)} | \mathbf{x}^{(i)}, \mathbf{W}) \\ &= -\log \prod_{i=1}^n \prod_{k=1}^K h_k(\mathbf{x}^{(i)}; \mathbf{W})^{y_k^{(i)}} \\ &= -\sum_{i=1}^n \sum_{k=1}^K y_k^{(i)} \log(h_k(\mathbf{x}^{(i)}; \mathbf{W})) \end{aligned}$$

$\mathbf{y}$  is a vector of length  $K$  (1-of- $K$  coding)

e.g.,  $\mathbf{y} = [0, 0, 1, 0]^T$  when the target class is  $C_3$

# Logistic regression: multi-class

- ▶ The gradient descent update

$$\mathbf{w}_j^{t+1} = \mathbf{w}_j^t - \eta \nabla_{\mathbf{W}} J(\mathbf{W}^t)$$

$$\nabla_{\mathbf{w}_j} J(\mathbf{W}) = \sum_{i=1}^n \left( h_j(\mathbf{x}^{(i)}; \mathbf{W}) - y_j^{(i)} \right) \mathbf{x}^{(i)}$$

# LR vs. GDA

- ▶  $d$ -dimensional feature space

- ▶ Logistic regression:  $d + 1$  parameters

$$\mathbf{w} = (w_0, w_1, \dots, w_d)$$

- ▶ GDA with shared covariance matrix

- ▶  $2d$  parameters for means
    - ▶  $d(d + 1)/2$  parameters for shared covariance matrix
    - ▶ one parameter for class prior  $p(C_1)$

- ▶ LR is more robust, less sensitive to incorrect modeling assumptions

# Summary of alternatives

## ▶ Generative

- ▶ Most demanding, because it finds the joint distribution  $p(\mathbf{x}, \mathcal{C}_k)$
- ▶ Usually needs a large training set to find  $p(\mathbf{x}|\mathcal{C}_k)$
- ▶ Can find  $p(\mathbf{x}) \Rightarrow$  Outlier detection

## ▶ Discriminative

- ▶ Specifies what is really needed (i.e.,  $p(\mathcal{C}_k|\mathbf{x})$ )
- ▶ More computationally efficient



# Generalization of Bayes decision rule

## Minimizing Bayes risk (expected loss)

$$\begin{aligned} & E_{\mathbf{x},y}[L(\alpha(\mathbf{x}), y)] \\ &= \int \sum_{j=1}^K L(\alpha(\mathbf{x}), \mathcal{C}_j) p(\mathbf{x}, \mathcal{C}_j) d\mathbf{x} \\ &= \int p(\mathbf{x}) \underbrace{\sum_{j=1}^K L(\alpha(\mathbf{x}), \mathcal{C}_j) p(\mathcal{C}_j|\mathbf{x})}_{\text{conditional risk}} d\mathbf{x} \end{aligned}$$

for each  $\mathbf{x}$  minimize it that is called conditional risk

- Bayes minimum loss (risk) decision rule:  $\hat{\alpha}(\mathbf{x})$

$$\hat{\alpha}(\mathbf{x}) = \operatorname{argmin}_{i=1,\dots,K} \sum_{j=1}^K \underset{\downarrow}{L_{ij}} p(\mathcal{C}_j|\mathbf{x})$$

The loss of assigning a sample to  $\mathcal{C}_i$  where the correct class is  $\mathcal{C}_j$

# Minimizing expected loss: special case (loss = misclassification rate)

- ▶ Problem definition for this special case:

- ▶ If action  $\alpha(\mathbf{x}) = i$  is taken and the true category is  $\mathcal{C}_j$ , then the decision is correct if  $i = j$  and otherwise it is incorrect.

- ▶ Zero-one loss function:

$$L_{ij} = 1 - \delta_{ij} = \begin{cases} 0 & i = j \\ 1 & o.w. \end{cases}$$

$$\hat{\alpha}(\mathbf{x}) = \operatorname{argmin}_{i=1,\dots,K} \sum_{j=1}^K L_{ij} p(\mathcal{C}_j | \mathbf{x})$$

$$= \operatorname{argmin}_{i=1,\dots,K} 0 \times p(\mathcal{C}_i | \mathbf{x}) + \sum_{j \neq i} p(\mathcal{C}_j | \mathbf{x})$$

$$= \operatorname{argmin}_{i=1,\dots,K} 1 - p(\mathcal{C}_i | \mathbf{x}) = \operatorname{argmax}_{i=1,\dots,K} p(\mathcal{C}_i | \mathbf{x})$$

# Resources

- ▶ [1] C. Bishop, “Pattern Recognition and Machine Learning”, Chapter 4.2-4.3.
- ▶ [2]: Andrew Ng, Machine learning, Stanford
- ▶ [3]: Pattern classification, Duda, Hart & Stork, 2002
- ▶ [4]: Mahdieh Soleymani, Machine learning, Sharif university of technology