# Probabilistic classification

Machine Learning

Hamid R Rabiee – Zahra Dehghanian
Spring 2025

Sharif University
of Technology

# Topics

- Probabilistic approach
  - Bayes decision theory
  - Generative models
    - Gaussian Bayes classifier
    - Naïve Bayes

**Probabilistic classification**

Sharif University
of Technology

# Classification problem: probabilistic view

- Given: Training set
  - labeled set of $N$ input-output pairs $D = \left\{ \left( \boldsymbol{x}^{(i)}, y^{(i)} \right) \right\}_{i=1}^{N}$
  - $y \in \{1, \dots, K\}$

- Goal: Given an input $\boldsymbol{x}$, assign it to one of $K$ classes

- Examples:
  - Spam filter
  - Handwritten digit recognition
  - …

**Probabilistic classification**

Sharif University
of Technology

# Definitions

- Posterior probability: $p(\mathcal{C}_k|\boldsymbol{x})$

- Likelihood or class conditional probability: $p(\boldsymbol{x}|\mathcal{C}_k)$

- Prior probability: $p(\mathcal{C}_k)$

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

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

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

**Probabilistic classification**

Sharif University
of Technology

# Bayes decision rule

- If $P(C_1|x) > P(C_2|x)$ decide $C_1$
  otherwise decide $C_2$

$$p(error|x) = \begin{cases} p(C_2|x) & \text{if we decide } C_1 \\ P(C_1|x) & \text{if we decide } C_2 \end{cases}$$

▸ If we use Bayes decision rule:

$$P(error|x) = \min\{P(C_1|x), P(C_2|x)\}$$

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

**Probabilistic classification**

Sharif University
of Technology

# Optimal classifier

- The optimal decision is the one that minimizes the expected number of mistakes

- We show that Bayes classifier is an optimal classifier

**Probabilistic classification**

Sharif University
of Technology

# Bayes decision rule
## Minimizing misclassification rate

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

$K = 2$

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

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

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

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

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

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

**Probabilistic classification**

Sharif University
of Technology

# Bayes minimum error

- Bayes minimum error classifier:

$$\min_{\alpha(.)} E_{x,y}\left[I(\alpha(x) \neq y)\right]$$    Zero-one loss

  - If we know the probabilities in advance then the above optimization problem will be solved easily.
    - $\alpha(x) = \underset{y}{\arg\max}\, p(y|x)$

- In practice, we can estimate $p(y|x)$ based on a set of training samples $\mathcal{D}$

**Probabilistic classification**

# Bayes theorem

- 

▸ Bayes' theorem

Posterior

Likelihood

Prior

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

▸ Posterior probability: $p(\mathcal{C}_k|\boldsymbol{x})$

▸ Likelihood or class conditional probability: $p(\boldsymbol{x}|\mathcal{C}_k)$

▸ Prior probability: $p(\mathcal{C}_k)$

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

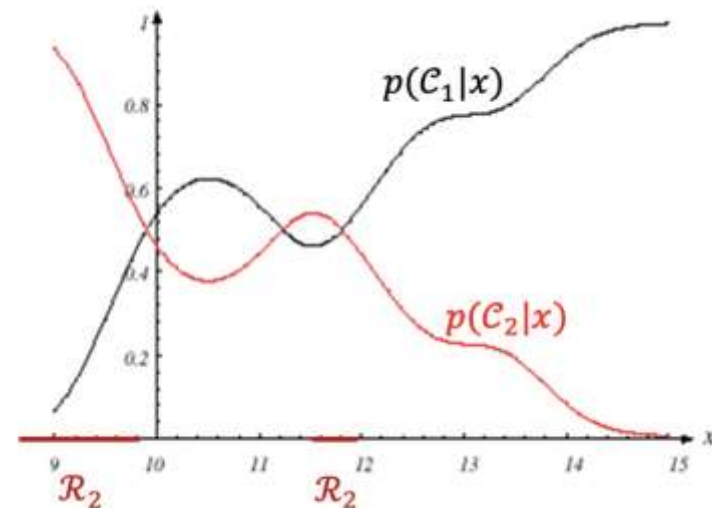**Probabilistic classification**

Sharif University
of Technology

# Bayes decision rule: example

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



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

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

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

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

**Probabilistic classification**

Sharif University
of Technology

# Bayesian decision rule

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

  otherwise decide $\mathcal{C}_2$

  Equivalent

- If $\dfrac{p(\boldsymbol{x}|\mathcal{C}_1)P(\mathcal{C}_1)}{p(\boldsymbol{x})} > \dfrac{p(\boldsymbol{x}|\mathcal{C}_2)P(\mathcal{C}_2)}{p(\boldsymbol{x})}$ decide $\mathcal{C}_1$

  otherwise decide $\mathcal{C}_2$

  Equivalent

- If $p(\boldsymbol{x}|\mathcal{C}_1)P(\mathcal{C}_1) > p(\boldsymbol{x}|\mathcal{C}_2)P(\mathcal{C}_2)$ decide $\mathcal{C}_1$

  otherwise decide $\mathcal{C}_2$

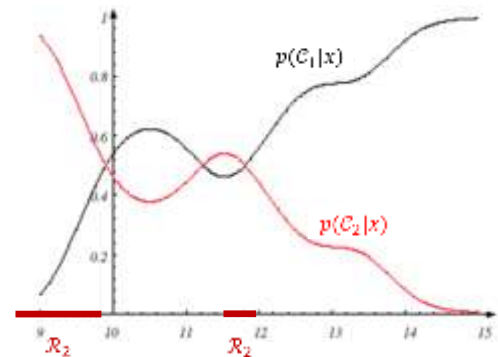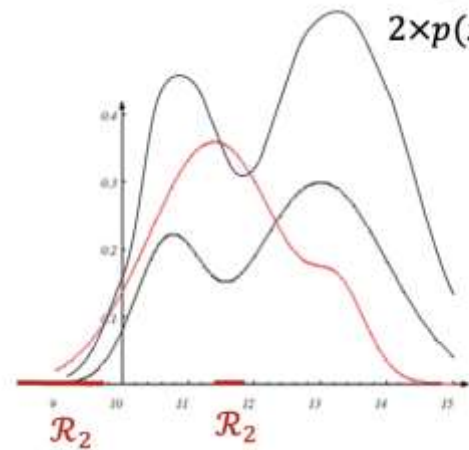**Probabilistic classification**

Sharif University
of Technology

# Bayes decision rule: example

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



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

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

**Probabilistic classification**

Sharif University
of Technology

# Bayes Classier

- Simple Bayes classifier: estimate posterior probability of each class

- What should the decision criterion be?
  - Choose class with highest $p(\mathcal{C}_k|\boldsymbol{x})$

- The optimal decision is the one that minimizes the expected number of mistakes

**Probabilistic classification**

Sharif University
of Technology

# Diabetes example

- white blood cell count



This example has been adopted from Sanja Fidler's slides, University of Toronto, CSC411

**Probabilistic classification**

Sharif University
of Technology

# 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

**Probabilistic classification**

Sharif University
of Technology

# Diabetes example

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



— $p(x|y = 0)$ (no diabetes)
— $p(x|y = 1)$ (diabetes)

This example has been adopted from Sanja Fidler's slides, University of Toronto, CSC411

**Probabilistic classification**

Sharif University
of Technology

# Estimate probability densities from data

- If we assume Gaussian distributions for $p(x|y = 0)$ and $p(x|y = 1)$

- Recall that for samples $\{x^{(1)}, \ldots, x^{(N)}\}$, if we assume a Gaussian distribution, the MLE estimates will be

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

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

**Probabilistic classification**

Sharif University
of Technology

# Diabetes example



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

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

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

This example has been adopted from Sanja Fidler's slides, University of Toronto, CSC411

**Probabilistic classification**

Sharif University
of Technology

# Diabetes example

- Add a second observation: Plasma glucose value



This example has been adopted from Sanja Fidler's slides, University of Toronto, CSC411

**Probabilistic classification**

Sharif University
of Technology

# Naïve Bayes classifier

- Generative methods
    - High number of parameters

- Assumption: Conditional independence

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

**Probabilistic classification**

# Naïve Bayes classifier

- In the decision phase, it finds the label of $x$ according to:

$$\operatorname*{argmax}_{k=1,\ldots,K} p(C_k|\boldsymbol{x})$$

$$\operatorname*{argmax}_{k=1,\ldots,K} p(C_k) \prod_{i=1}^{n} p(x_i|C_k)$$

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

$$p(C_k|\boldsymbol{x}) \propto p(C_k) \prod_{i=1}^{n} p(x_i|C_k)$$

**Probabilistic classification**

Sharif University
of Technology

# Naïve Bayes: discrete example

- $p(h) = 0.3$

- $p(d|h) = \frac{1}{3}$
- $p(s|h) = \frac{2}{3}$

- $p(d|\bar{h}) = \frac{2}{7}$
- $p(s|\bar{h}) = \frac{2}{7}$

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

| 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 |

**Probabilistic classification**

# Naïve Bayes: discrete example

- $p(h) = 0.3$

- $p(d|h) = \frac{1}{3}$

- $p(s|h) = \frac{2}{3}$

- $p(d|\bar{h}) = \frac{2}{7}$

- $p(s|\bar{h}) = \frac{2}{7}$

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

| 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 $x = [d, \bar{s}]$ (a person that has diabetes but does not smoke):
  - $p(h|x) \propto p(h)p(d|h)p(\bar{s}|h) = 1/14$
  - $p(\bar{h}|x) \propto p(\bar{h})p(d|\bar{h})p(\bar{s}|\bar{h}) = 1/6$
  - Thus decide $H = No$

**Probabilistic classification**

# Naïve Bayes classifier

- Finds $d$ univariate distributions $p(x_1|C_k), \cdots, p(x_d|C_k)$ instead of finding one multi-variate distribution $p(\boldsymbol{x}|C_k)$

  - Example 1: For Gaussian class-conditional density $p(\boldsymbol{x}|C_k)$, 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 $p(\boldsymbol{x}|C_k)$, it finds $d$ (mean parameters on different dimensions) instead of $2^d - 1$ parameters

- It first estimates the class conditional densities $p(x_1|C_k), \cdots, p(x_d|C_k)$ and the prior probability $p(C_k)$ for each class ($k = 1, \ldots, K$) based on the training set.

**Probabilistic classification**

Sharif University of Technology

# Multivariate Gaussian

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

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

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

**Probabilistic classification**

# Multivariate Gaussian

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

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

$$k = 1,2$$

- Prior distribution $p(y)$:
  - $p(y = 1) = \pi, \qquad p(y = 0) = 1 - \pi$

**Probabilistic classification**

# Multivariate Gaussian

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

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

- $\pi = \dfrac{N_1}{N}$

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

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

$$N_2 = N - N_1$$

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

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

**Probabilistic classification**

Sharif University
of Technology

# Decision boundary for Gaussian Bayes classifier

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

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

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

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

**Probabilistic classification**

Sharif University
of Technology

# Decision boundary for Gaussian Bayes classifier

- $$p(C_1|\boldsymbol{x}) = p(C_2|\boldsymbol{x})$$

$$p(C_k|\boldsymbol{x}) = \frac{p(\boldsymbol{x}|C_k)p(C_k)}{p(\boldsymbol{x})}$$

$$\ln p(C_1|\boldsymbol{x}) = \ln p(C_2|\boldsymbol{x})$$

$$\ln p(\boldsymbol{x}|C_1) + \ln p(C_1) - \ln p(\boldsymbol{x})$$
$$= \ln p(\boldsymbol{x}|C_2) + \ln p(C_2) - \ln p(\boldsymbol{x})$$

$$\ln p(\boldsymbol{x}|C_1) + \ln p(C_1) = \ln p(\boldsymbol{x}|C_2) + \ln p(C_2)$$

$$\ln p(\boldsymbol{x}|C_k)$$
$$= -\frac{d}{2}\ln 2\pi - \frac{1}{2}\ln|\boldsymbol{\Sigma}_k| - \frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1}(\boldsymbol{x} - \boldsymbol{\mu}_k)$$

**Probabilistic classification**

Sharif University
of Technology

# Decision boundary



likelihoods

discriminant:
$P(t_1 | \mathbf{x}) = 0.5$

posterior for $t_1$

**Probabilistic classification**

Sharif University
of Technology

# Shared covariance matrix

- When classes share a single covariance matrix $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2$

$$p(\boldsymbol{x}|C_k) = \frac{1}{(2\pi)^{d/2}|\Sigma|^{1/2}} \exp\{-\frac{1}{2}(\boldsymbol{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}^{-1}(\boldsymbol{x} - \boldsymbol{\mu}_k)\}$$

$$k = 1,2$$

- $p(C_1) = \pi, \qquad p(C_2) = 1 - \pi$

**Probabilistic classification**

Sharif University
of Technology

# Likelihood

- $$\prod_{n=1}^{N} p(\boldsymbol{x}^{(n)}, y^{(n)}; \pi, \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma})$$

  $$= \prod_{n=1}^{N} p(\boldsymbol{x}^{(n)}|y^{(n)}; \boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}) p(y^{(n)}; \pi)$$

**Probabilistic classification**

Sharif University
of Technology

# Shared covariance matrix

- Maximum likelihood estimation ($D = \{(\boldsymbol{x}^{(i)}, y^{(i)})\}_{i=1}^{n}$):

$$\pi = \frac{N_1}{N}$$

$$\boldsymbol{\mu}_1 = \frac{\sum_{n=1}^{N} y^{(n)} \boldsymbol{x}^{(n)}}{N_1}$$

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

$$\boldsymbol{\Sigma} = \frac{1}{N} \left( \sum_{n \in C_1} (\boldsymbol{x}^{(n)} - \boldsymbol{\mu}_1)(\boldsymbol{x}^{(n)} - \boldsymbol{\mu}_1)^T + \sum_{n \in C_2} (\boldsymbol{x}^{(n)} - \boldsymbol{\mu}_2)(\boldsymbol{x}^{(n)} - \boldsymbol{\mu}_2)^T \right)$$

**Probabilistic classification**

Sharif University
of Technology

# Decision boundary when shared covariance matrix

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

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

**Probabilistic classification**

Sharif University
of Technology

# Multi-class Bayes decision rule

- Multi-class problem: Probability of error of Bayesian decision rule
  - Simpler to compute the probability of correct decision

$$P(error) = 1 - P(correct)$$

$$P(Correct) = \sum_{i=1}^{K} \int_{\mathcal{R}_i} p(\boldsymbol{x}, \mathcal{C}_i) \, d\boldsymbol{x}$$

$$= \sum_{i=1}^{K} \int_{\mathcal{R}_i} p(\mathcal{C}_i|\boldsymbol{x}) p(\boldsymbol{x}) \, d\boldsymbol{x}$$

$\mathcal{R}_i$: the subset of feature space assigned to the class $\mathcal{C}_i$ using the classifier

**Probabilistic classification**

Sharif University
of Technology

# Bayes minimum error

- Bayes minimum error classifier:

$$\min_{\alpha(.)} E_{x,y}\left[I(\alpha(x) \neq y)\right] \qquad \text{Zero-one loss}$$

$$\alpha(x) = \operatorname*{argmax}_{y} p(y|x)$$

**Probabilistic classification**

Sharif University
of Technology

# Minimizing Bayes risk (expected loss)

$$E_{x,y}[L(\alpha(x), y)]$$

$$= \int \sum_{j=1}^{K} L(\alpha(x), C_j) p(x, C_j) dx$$

**Probabilistic classification**

Sharif University
of Technology

# Minimizing Bayes risk (expected loss)

$$E_{x,y}[L(\alpha(x), y)]$$

$$= \int \sum_{j=1}^{K} L(\alpha(x), \mathcal{C}_j) p(x, \mathcal{C}_j) dx$$

$$= \int p(x) \underbrace{\sum_{j=1}^{K} L(\alpha(x), \mathcal{C}_j) p(\mathcal{C}_j|x)}_{} dx$$

for each $x$ minimize it that is called conditional risk

**Probabilistic classification**

# Minimizing Bayes risk (expected loss)

$$E_{x,y}[L(\alpha(x), y)]$$

$$= \int \sum_{j=1}^{K} L(\alpha(x), C_j) p(x, C_j) dx$$

$$= \int p(x) \underbrace{\sum_{j=1}^{K} L(\alpha(x), C_j) p(C_j|x)} dx$$

for each $x$ minimize it that is called conditional risk

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

$$\hat{\alpha}(x) = \operatorname*{argmin}_{i=1,\ldots,K} \sum_{j=1}^{K} L_{ij} p(C_j|x)$$

The loss of assigning a sample to $C_i$ where the correct class is $C_j$

**Probabilistic classification**

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

- Problem definition for this special case:

  - If action $\alpha(\boldsymbol{x}) = i$ 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:

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

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

**Probabilistic classification**

Sharif University of Technology

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

- Problem definition for this special case:
  - If action $\alpha(\boldsymbol{x}) = i$ 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:

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

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

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

$$= \underset{i=1,\ldots,K}{\operatorname{argmin}} \, 1 - p(\mathcal{C}_i|\boldsymbol{x}) = \underset{i=1,\ldots,K}{\operatorname{argmax}} \, p(\mathcal{C}_i|\boldsymbol{x})$$

**Probabilistic classification**

# Probabilistic classifiers

- How can we find the probabilities required in the Bayes decision rule?

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

**Probabilistic classification**

# Generative approach

- <u>Inference stage</u>
  - Determine class conditional densities $p(\boldsymbol{x}|\mathcal{C}_k)$ and priors $p(\mathcal{C}_k)$
  - Use the Bayes theorem to find $p(\mathcal{C}_k|\boldsymbol{x})$

- <u>Decision stage:</u> After learning the model (inference stage), make optimal class assignment for new input
  - if $p(\mathcal{C}_i|\boldsymbol{x}) > p(\mathcal{C}_j|\boldsymbol{x}) \quad \forall j \neq i$ then decide $\mathcal{C}_i$

**Probabilistic classification**

Sharif University
of Technology

# Probabilistic discriminant functions

- **Discriminant functions:** A popular way of representing a classifier
  - A discriminant function $f_i(x)$ for each class $C_i$ $(i = 1, \dots, K)$:
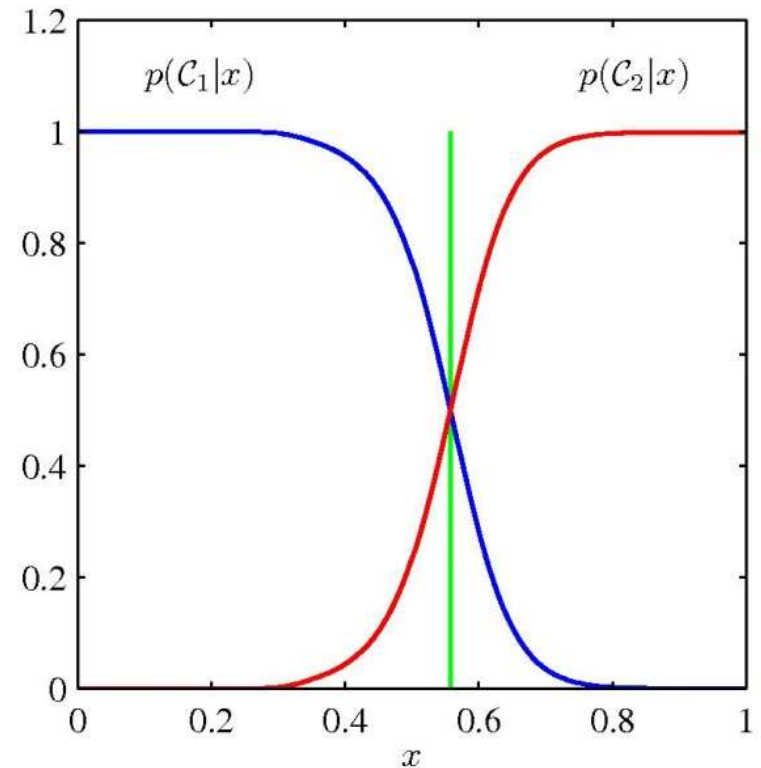    - $x$ is assigned to class $C_i$ if:

$$f_i(x) > f_j(x) \quad \forall j \neq i$$

- Representing Bayesian classifier using discriminant functions:
  - Classifier minimizing error rate: $f_i(x) = P(C_i|x)$
  - Classifier minimizing risk: $f_i(x) = -\sum_{j=1}^{K} L_{ij} p(C_j|x)$

**Probabilistic classification**

Sharif University
of Technology

# Discriminative vs. generative approach



[Bishop]

**Probabilistic classification**

Sharif University
of Technology

# Class conditional densities vs. posterior



[Bishop]

$$p(C_1|x) = \sigma(w^T x + w_0)$$

$$w = \Sigma^{-1}(\mu_1 - \mu_2)$$

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

$$w_0 = -\frac{1}{2}\mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2}\mu_2^T \Sigma^{-1} \mu_2 + \ln\frac{p(C_1)}{p(C_2)}$$

**Probabilistic classification**

Sharif University
of Technology

# Feed back

? https://forms.gle/vKRbyVVsWRKcZuqr8

**Regression: Probabilistic perspective**

Sharif University
of Technology

# Resources

- C. Bishop, "Pattern Recognition and Machine Learning", Chapter 4.2-4.3.

- Course CE-717, Dr. M.Soleymani

**Discriminative models**

Sharif University
of Technology