

COSC 522 – Machine Learning

Discriminant Functions

Hairong Qi, Gonzalez Family Professor
Electrical Engineering and Computer Science
University of Tennessee, Knoxville
<https://www.eecs.utk.edu/people/hairong-qi/>
Email: hqi@utk.edu

Course Website: <http://web.eecs.utk.edu/~hqi/cosc522/>

Recap from Previous Lecture

- Definition of supervised learning (vs. unsupervised learning)
- The difference between the training set and the test set
- The difference between classification and regression
- Definition of “features”, “samples”, and “dimension”
- From histogram to probability density function (pdf)
- In Bayes' Formula, what is conditional pdf? Prior probability? Posterior probability?
- What does the normalization factor (or evidence) do?
- What is Bayesian decision rule? or MPP?
- What are decision regions?
- How to calculate conditional probability of error and overall probability of error?
- What are cost function (or objective function) and optimization method used in MPP?

terminologies

the Formula

decision rule

Recap

$$P(\omega_j | x) = \frac{p(x | \omega_j) P(\omega_j)}{p(x)}$$



Maximum
Posterior
Probability

For a given x , if $P(\omega_1 | x) > P(\omega_2 | x)$,
then x belongs to class 1, otherwise, 2.

Overall
probability
of error

$$P(\text{error}) = \int_{\mathfrak{R}_1} P(\omega_2 | x) p(x) dx + \int_{\mathfrak{R}_2} P(\omega_1 | x) p(x) dx$$

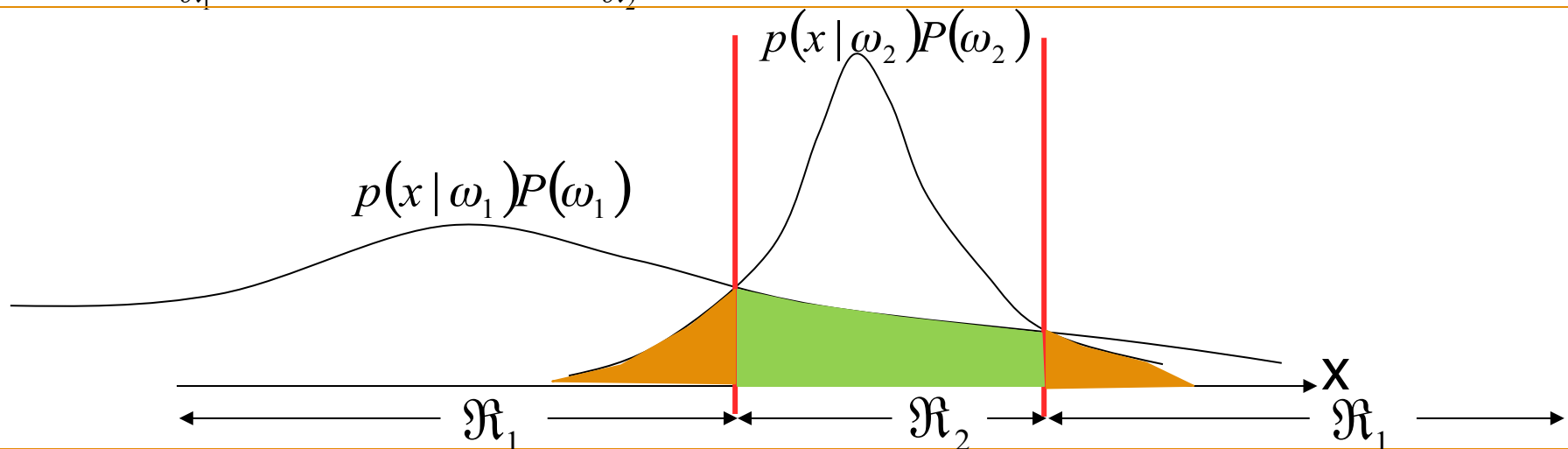
Decision Rule → Decision Region →
Conditional Probability of Error → Overall Probability of Error

$$P(\text{error} | x) = \begin{cases} P(\omega_1 | x) & \text{if we decide } \omega_2 \\ P(\omega_2 | x) & \text{if we decide } \omega_1 \end{cases} = \min[P(\omega_1 | x), P(\omega_2 | x)]$$

Unconditional risk, unconditional probability of error

$$P(\text{error}) = \int_{-\infty}^{\infty} P(\text{error}, x) dx = \int_{-\infty}^{\infty} P(\text{error} | x) p(x) dx$$

$$P(\text{error}) = \int_{\mathcal{R}_1}^{-\infty} P(\omega_2 | x) p(x) dx + \int_{\mathcal{R}_2}^{-\infty} P(\omega_1 | x) p(x) dx = P(\text{error} | \omega_2) + P(\text{error} | \omega_1)$$



Questions

- What is a discriminant function?
- What is a multivariate Gaussian (or normal density function)?
- What is the covariance matrix and what is its dimension?
- What would the covariance matrix look like if the features are independent from each other?
- What would the covariance matrix look like if the features are independent from each other AND have the same spread in each dimension?
- What is minimum (Euclidean) distance classifier? Is it a linear or quadratic classifier (machine)? What does the decision boundary look like?
- What are the assumptions made when using a minimum (Euclidean) distance classifier?
- What is minimum (Mahalanobis) distance classifier? Is it a linear or quadratic classifier (machine)? What does the decision boundary look like?
- What are the assumptions made when using a minimum (Mahalanobis) distance classifier?
- What does the decision boundary look like for a quadratic classifier?
- What are the cost functions for the discriminant functions? And what is the optimization method used to find the best solution?

Multi-variate Gaussian

Linear and Quadratic Machines
and their assumptions

Discriminant Function

- ◆ One way to represent pattern classifier- use discriminant functions $g_i(x)$

$$g_i(x) = P(\omega_i|x)$$

$$g_i(x) = p(x|\omega_i)P(\omega_i)$$

$$g_i(x) = \ln p(x|\omega_i) + \ln P(\omega_i)$$

The classifier will assign a feature vector x to class ω_i if

$$g_i(x) > g_j(x)$$

- ◆ For two-class cases,

$$g(x) = g_1(x) - g_2(x) = P(\omega_1|x) - P(\omega_2|x)$$

Multivariate Normal Density

$$p(\vec{x}) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left[-\frac{1}{2} (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})\right]$$

\vec{x} : d - component column vector

$\vec{\mu}$: d - component mean vector

Σ : d - by - d covariance matrix

$|\Sigma|$: determinant

Σ^{-1} : inverse

$$\vec{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}, \vec{\mu} = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_d \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} \sigma_{11} & \cdots & \sigma_{1d} \\ \vdots & \ddots & \vdots \\ \sigma_{d1} & \cdots & \sigma_{dd} \end{bmatrix} = \begin{bmatrix} \sigma_1^2 & \cdots & \sigma_{1d} \\ \vdots & \ddots & \vdots \\ \sigma_{d1} & \cdots & \sigma_d^2 \end{bmatrix}$$

When $d = 1$, $p(x) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left[-\frac{1}{2} \frac{(x - \mu)^2}{\sigma^2}\right]$

Estimating Normal Densities

◆ Calculate μ , Σ

$$\vec{\mu}_i = \begin{bmatrix} \mu_{i1} = \frac{1}{n_i} \sum_{k=1}^{n_i} x_{k1} \\ \vdots \\ \mu_{id} = \frac{1}{n_i} \sum_{k=1}^{n_i} x_{kd} \end{bmatrix}$$

$$\Sigma_i = \begin{bmatrix} \sigma_{11} & \cdots & \sigma_{1d} \\ \vdots & \ddots & \vdots \\ \sigma_{d1} & \cdots & \sigma_{dd} \end{bmatrix} = \frac{1}{n_i - 1} \sum_{k=1}^{n_i} (\vec{x}_k - \vec{\mu}_i)(\vec{x}_k - \vec{\mu}_i)^T$$

Covariance

For d sets of variates denoted $\{x_1\}, \dots, \{x_p\}, \dots, \{x_q\}, \dots, \{x_d\}$
 the covariance $\sigma_{pq} = \text{cov}(x_p, x_q)$ of x_p and x_q is defined by

$$\text{cov}(x_p, x_q) = E[(x_p - \mu_p)(x_q - \mu_q)]$$

$$= E[x_p x_q] - E[x_p \mu_q] - E[\mu_p x_q] + E[\mu_p \mu_q]$$

$$= E[x_p x_q] - \mu_q E[x_p] - \mu_p E[x_q] + \mu_p \mu_q$$

$$= E[x_p x_q] - \mu_q \mu_p - \mu_p \mu_q + \mu_p \mu_q$$

$$= E[x_p x_q] - \mu_q \mu_p$$

When $p = q$, $\sigma_{pp} = \text{cov}(x_p, x_p) = E[x_p x_p] - \mu_p \mu_p$

$$= E[x_p^2] - (E[x_p])^2$$

$$= \sigma_p^2$$

Discriminant Function for Normal Density

$$p(\vec{x} | w) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left[-\frac{1}{2} (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu}) \right]$$

$$\begin{aligned} g_i(\vec{x}) &= \ln p(\vec{x} | \omega_i) + \ln P(\omega_i) \\ &= -\frac{1}{2} (\vec{x} - \vec{\mu}_i)^T \Sigma_i^{-1} (\vec{x} - \vec{\mu}_i) - \frac{d}{2} \ln(2\pi) - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i) \\ &= -\frac{1}{2} (\vec{x} - \vec{\mu}_i)^T \Sigma_i^{-1} (\vec{x} - \vec{\mu}_i) - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i) \end{aligned}$$

Questions

- What is a discriminant function?
- What is a multivariate Gaussian (or normal density function)?
- What is the covariance matrix and what is its dimension?
- What would the covariance matrix look like if the features are independent from each other?
- What would the covariance matrix look like if the features are independent from each other AND have the same spread in each dimension?
- What is minimum (Euclidean) distance classifier? Is it a linear or quadratic classifier (machine)? What does the decision boundary look like?
- What are the assumptions made when using a minimum (Euclidean) distance classifier?
- What is minimum (Mahalanobis) distance classifier? Is it a linear or quadratic classifier (machine)? What does the decision boundary look like?
- What are the assumptions made when using a minimum (Mahalanobis) distance classifier?
- What does the decision boundary look like for a quadratic classifier?
- What are the cost functions for the discriminant functions? And what is the optimization method used to find the best solution?

Multi-variate Gaussian

Linear and Quadratic Machines
and their assumptions

Case 1: $\Sigma_i = \sigma^2 I$

- ◆ The features are statistically independent, and have the same variance
- ◆ Geometrically, the samples fall in equal-size hyperspherical clusters
- ◆ Decision boundary: hyperplane of $d-1$ dimension

$$\Sigma = \begin{bmatrix} \sigma^2 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \sigma^2 \end{bmatrix}, |\Sigma| = \sigma^{2d}, \Sigma^{-1} = \begin{bmatrix} \frac{1}{\sigma^2} & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \frac{1}{\sigma^2} \end{bmatrix}$$

Linear Discriminant Function and Linear Machine

$\|\vec{x} - \vec{\mu}_i\|$: the Euclidean norm (distance)

$$\|\vec{x} - \vec{\mu}_i\|^2 = (\vec{x} - \vec{\mu}_i)^T (\vec{x} - \vec{\mu}_i)$$

$$\begin{aligned} g_i(\vec{x}) &= -\frac{\|\vec{x} - \vec{\mu}_i\|^2}{2\sigma^2} + \ln P(\omega_i) \\ &= -\frac{\vec{x}^T \vec{x} - 2\vec{\mu}_i^T \vec{x} + \vec{\mu}_i^T \vec{\mu}_i}{2\sigma^2} + \ln P(\omega_i) \end{aligned}$$

$$g_i(\vec{x}) = \frac{\vec{\mu}_i^T}{\sigma^2} \vec{x} - \frac{\vec{\mu}_i^T \vec{\mu}_i}{2\sigma^2} + \ln P(\omega_i)$$

Minimum-Distance Classifier

- ◆ When $P(\omega_i)$ are the same for all c classes, the discriminant function is actually measuring the minimum distance from each x to each of the c mean vectors

$$g_i(\vec{x}) = -\frac{\|\vec{x} - \vec{\mu}_i\|^2}{2\sigma^2}$$

Case 2: $\Sigma_i = \Sigma$

- ◆ The covariance matrices for all the classes are identical but not a scalar of identity matrix.
- ◆ Geometrically, the samples fall in hyperellipsoidal
- ◆ Decision boundary: hyperplane of d-1 dimension

$$g_i(\vec{x}) = \ln p(\vec{x} | \omega_i) + \ln P(\omega_i)$$

$$= -\frac{1}{2} (\vec{x} - \vec{\mu}_i)^T \Sigma_i^{-1} (\vec{x} - \vec{\mu}_i) + \ln P(\omega_i)$$

$$= \vec{\mu}_i^T (\Sigma^{-1})^T \vec{x} - \frac{1}{2} \vec{\mu}_i^T \Sigma^{-1} \vec{\mu}_i + \ln P(\omega_i)$$

Squared Mahalanobis distance

Case 3: $\Sigma_i = \text{arbitrary}$

- ◆ The covariance matrices are different from each category
- ◆ Quadratic classifier
- ◆ Decision boundary: hyperquadratic for 2-D Gaussian

$$\begin{aligned}g_i(\vec{x}) &= \ln p(\vec{x} | \omega_i) + \ln P(\omega_i) \\&= -\frac{1}{2}(\vec{x} - \vec{\mu}_i)^T \Sigma_i^{-1}(\vec{x} - \vec{\mu}_i) - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i) \\&= -\frac{1}{2} \vec{x}^T \Sigma_i^{-1} \vec{x} + \vec{\mu}_i^T (\Sigma_i^{-1}) \vec{x} - \frac{1}{2} \vec{\mu}_i^T \Sigma_i^{-1} \vec{\mu}_i - \frac{1}{2} \ln |\Sigma_i| + \ln P(\omega_i)\end{aligned}$$

Questions

- What is a discriminant function?
- What is a multivariate Gaussian (or normal density function)?
- What is the covariance matrix and what is its dimension?
- What would the covariance matrix look like if the features are independent from each other?
- What would the covariance matrix look like if the features are independent from each other AND have the same spread in each dimension?
- What is minimum (Euclidean) distance classifier? Is it a linear or quadratic classifier (machine)? What does the decision boundary look like?
- What are the assumptions made when using a minimum (Euclidean) distance classifier?
- What is minimum (Mahalanobis) distance classifier? Is it a linear or quadratic classifier (machine)? What does the decision boundary look like?
- What are the assumptions made when using a minimum (Mahalanobis) distance classifier?
- What does the decision boundary look like for a quadratic classifier?
- **What are the cost functions for the discriminant functions? And what is the optimization method used to find the best solution?**

Multi-variate Gaussian

Linear and Quadratic Machines
and their assumptions

Bayes Decision Rule

$$P(\omega_j | x) = \frac{p(x | \omega_j) P(\omega_j)}{p(x)}$$

Maximum
Posterior
Probability

For a given x , if $P(\omega_1 | x) > P(\omega_2 | x)$,
then x belongs to class 1, otherwise, 2.

Discriminant
Function

The classifier will assign a feature vector x to class ω_i if
 $g_i(x) > g_j(x)$

Case 1: Minimum Euclidean Distance (Linear Machine), $\Sigma_i = \sigma^2 I$

Case 2: Minimum Mahalanobis Distance (Linear Machine), $\Sigma_i = \Sigma$

Case 3: Quadratic classifier, $\Sigma_i = \text{arbitrary}$

All assuming Gaussian pdf