

ECE595 / STAT598: Machine Learning I

Lecture 11.1: Maximum-Likelihood Estimation - Basic Principles

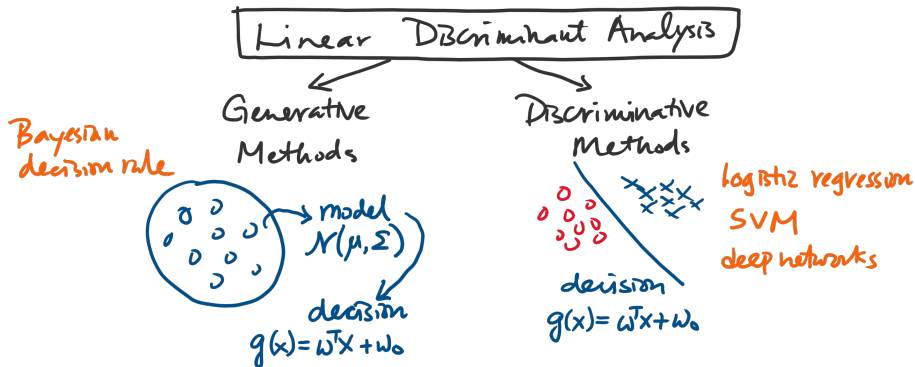
Spring 2020

Stanley Chan

School of Electrical and Computer Engineering
Purdue University



Overview



- In linear discriminant analysis (LDA), there are generally two types of approaches
- **Generative approach:** Estimate model, then define the classifier
- **Discriminative approach:** Directly define the classifier

Outline

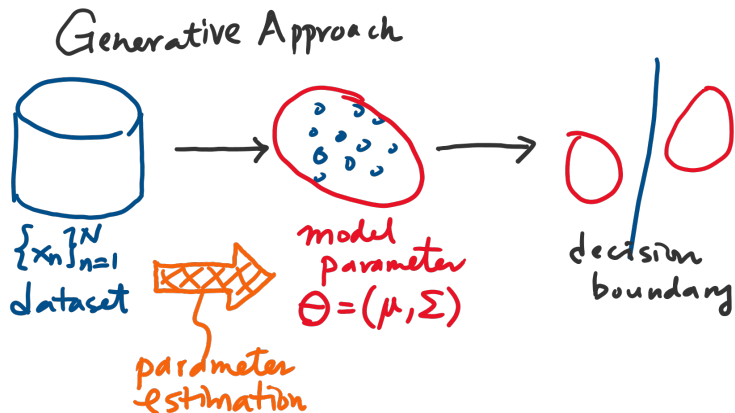
Generative Approaches

- Lecture 9 Bayesian Decision Rules
- Lecture 10 Evaluating Performance
- **Lecture 11 Parameter Estimation**
- Lecture 12 Bayesian Prior
- Lecture 13 Connecting Bayesian and Linear Regression

Today's Lecture

- **Basic Principles**
 - **Likelihood Function**
 - **Maximum Likelihood Estimate**
 - **1D Illustration**
 - **Gaussian Distributions**
- **Examples**
 - Non-Gaussian Distributions
 - Biased and Unbiased Estimators
 - From MLE to MAP

What is Parameter Estimation?



- The goal of parameter estimation is to determine $\Theta = (\mu, \Sigma)$ from dataset
- This is *the step* where you use data

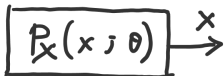
MLE and MAP

There are two typical ways of estimating parameters.

The Generative Process



Bayesian
(MAP estimation)



Frequentist
(ML estimation)

- Maximum-likelihood estimation (MLE): θ is deterministic.
- Maximum-a-posteriori estimation (MAP): θ is random and has a prior distribution.

Maximum Likelihood Estimation

Given the dataset $\mathcal{D} = \{\mathbf{x}_n\}_{n=1}^N$, how to estimate the model parameters?

- We are going to use Gaussian as an illustration.
- Denote θ as the model parameter.
- In Gaussian

$$\theta = \{\boldsymbol{\mu}, \boldsymbol{\Sigma}\}$$

- The likelihood for one data point \mathbf{x}_n is

$$p(\mathbf{x}_n | \overbrace{\theta}^{\{ \boldsymbol{\mu}, \boldsymbol{\Sigma} \}}) = \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \right\}$$

- θ is a deterministic quantity, not a random variable.
- θ does not have a distribution.
- θ is fixed but unknown.

Likelihood for the Entire Dataset

- Likelihood for the entire dataset $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ is

$$\begin{aligned} p(\mathcal{D} | \boldsymbol{\theta}) &= \prod_{n=1}^N \left\{ \frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \right\} \right\} \\ &= \left(\frac{1}{\sqrt{(2\pi)^d |\boldsymbol{\Sigma}|}} \right)^N \exp \left\{ \sum_{n=1}^N -\frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \right\} \end{aligned}$$

- The Negative Log-Likelihood is

$$\begin{aligned} -\log p(\mathcal{D} | \boldsymbol{\theta}) &= \frac{N}{2} \log |\boldsymbol{\Sigma}| + \frac{N}{2} \log (2\pi)^d \\ &\quad + \sum_{n=1}^N \left\{ \frac{1}{2} (\mathbf{x}_n - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}_n - \boldsymbol{\mu}) \right\}. \end{aligned}$$

Maximum Likelihood Estimation

- Goal: Find θ that maximizes the likelihood:

$$\begin{aligned}\hat{\theta} &= \operatorname{argmax}_{\theta} p(\mathcal{D} | \theta) \\ &= \operatorname{argmax}_{\theta} \prod_{n=1}^N \left\{ \frac{1}{\sqrt{(2\pi)^d |\Sigma|}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) \right\} \right\} \\ &= \operatorname{argmin}_{\theta} -\log(\dots) \\ &= \operatorname{argmin}_{\theta} \frac{N}{2} \log |\Sigma| + \frac{N}{2} \log(2\pi)^d \\ &\quad + \sum_{n=1}^N \left\{ \frac{1}{2} (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) \right\}.\end{aligned}$$

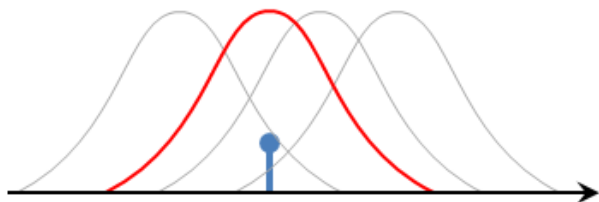
- This optimization is called the **maximum likelihood estimation** (MLE).

Illustrating MLE when $N = 1$. Known σ .

When $N = 1$: The MLE solution is

$$\begin{aligned}\hat{\mu} &= \operatorname{argmax}_{\mu} \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(x_1 - \mu)^2}{2\sigma^2}\right\} \\ &= \operatorname{argmin}_{\mu} (x_1 - \mu)^2 = x_1.\end{aligned}$$

- Which μ will give you the best Gaussian?
- When $\mu = x_1$, the probability of obtaining x_1 is the highest.

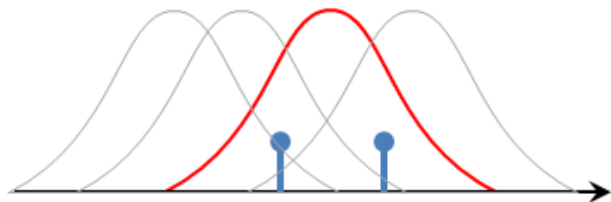


Illustrating MLE when $N = 2$. Known σ .

When $N = 2$: The MLE solution is

$$\begin{aligned}\hat{\mu} &= \operatorname{argmax}_{\mu} \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^2 \exp \left\{ -\frac{(x_1 - \mu)^2 + (x_2 - \mu)^2}{2\sigma^2} \right\} \\ &= \operatorname{argmin}_{\mu} (x_1 - \mu)^2 + (x_2 - \mu)^2 = \frac{x_1 + x_2}{2}.\end{aligned}$$

- Which μ will give you the best Gaussian?
- When $\mu = (x_1 + x_2)/2$, the prob. of obtaining x_1 and x_2 is highest.

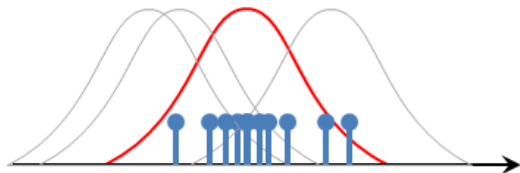


Illustrating MLE when $N =$ arbitrary integer

The MLE solution is

$$\begin{aligned}\hat{\mu} &= \operatorname{argmax}_{\mu} \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^2 \exp \left\{ - \sum_{n=1}^N \frac{(x_n - \mu)^2}{2\sigma^2} \right\} \\ &= \operatorname{argmin}_{\mu} \sum_{n=1}^N (x_n - \mu)^2 = \frac{1}{N} \sum_{n=1}^N x_n.\end{aligned}$$

- Which μ will give you the best Gaussian?
- When $\mu = \frac{1}{N} \sum_{n=1}^N x_n$, the prob. of obtaining $\{x_n\}$ is highest.



Estimation in High-dimension

- Assume Σ is known and fixed.
- Thus, $\theta = \mu$. Estimate μ

$$\begin{aligned}\hat{\mu} &= \operatorname{argmin}_{\mu} \cancel{\frac{N}{2} \log |\Sigma|} + \cancel{\frac{N}{2} \log (2\pi)^d} \\ &\quad + \sum_{n=1}^N \left\{ \frac{1}{2} (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) \right\} \\ &= \operatorname{argmin}_{\mu} \sum_{n=1}^N \left\{ (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) \right\}\end{aligned}$$

- Take derivative, setting to zero:

$$\nabla_{\mu} \left\{ \sum_{n=1}^N (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) \right\} = 2 \sum_{n=1}^N \Sigma^{-1} (\mathbf{x}_n - \mu) = \mathbf{0}.$$

Estimation in High-dimension

- Let us do some algebra

$$\sum_{n=1}^N \Sigma^{-1}(\mathbf{x}_n - \boldsymbol{\mu}) = \mathbf{0} \implies \sum_{n=1}^N \mathbf{x}_n = \sum_{n=1}^N \boldsymbol{\mu}$$

- Then we can show that the MLE solution is

$$\hat{\boldsymbol{\mu}} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n.$$

- This is just the **empirical average** of the entire dataset!
- You can show that if $\mathbb{E}[\mathbf{x}_n] = \boldsymbol{\mu}$ for all n , then

$$\mathbb{E}[\hat{\boldsymbol{\mu}}] = \frac{1}{N} \sum_{n=1}^N \mathbb{E}[\mathbf{x}_n] = \boldsymbol{\mu}.$$

- We say that $\hat{\boldsymbol{\mu}}$ is a **unbiased estimator** of $\boldsymbol{\mu}$ since $\mathbb{E}[\hat{\boldsymbol{\mu}}] = \boldsymbol{\mu}$.

When both μ and Σ are Unknown

What will be the MLE when both μ and Σ are unknown?

$$\begin{aligned}(\hat{\mu}, \hat{\Sigma}) = \operatorname{argmin}_{\mu, \Sigma} & \frac{N}{2} \log |\Sigma| + \frac{N}{2} \log(2\pi)^d \\ & + \underbrace{\sum_{n=1}^N \left\{ \frac{1}{2} (\mathbf{x}_n - \mu)^T \Sigma^{-1} (\mathbf{x}_n - \mu) \right\}}_{\varphi(\mu, \Sigma)}.\end{aligned}$$

You need to take derivative with respect to μ and Σ , and solve

$$\begin{aligned}\nabla_{\mu} \varphi(\mu, \Sigma) &= \mathbf{0} \\ \nabla_{\Sigma} \varphi(\mu, \Sigma) &= \mathbf{0}\end{aligned}$$

With some (tedious) matrix calculus, we can show that

$$\hat{\mu} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n, \quad \text{and} \quad \hat{\Sigma} = \frac{1}{N} \sum_{n=1}^N (\mathbf{x}_n - \hat{\mu})(\mathbf{x}_n - \hat{\mu})^T.$$

Exercise: Prove this result when \mathbf{x}_n is a 1D scalar.