

# Machine Learning in Robotics

## Lecture 5: Density Estimation - MLE

**Prof. Dongheui Lee**

*Institute of Automatic Control Engineering  
Technische Universität München*

dhlee@tum.de

# Summary of last lecture : Unsupervised Learning

- Iterative Optimization
- K-means, variations of K-means, Binary Splitting, LBG algorithm

1. Initialization: Choose  $K$  random vectors  $\mathbf{Y} = \{\mathbf{y}^{(1)}, \mathbf{y}^{(2)}, \dots, \mathbf{y}^{(K)}\}$

2. E-step: **Label each data point.**

$$\mathbf{X}_i = \left\{ \mathbf{x}^{(j)} \mid d(\mathbf{x}^{(j)}, \mathbf{y}^{(i)}) \leq d(\mathbf{x}^{(j)}, \mathbf{y}^{(k)}), j = 1, \dots, n; k = 1, \dots, K \right\}$$

3. M-step: From the current clusters, their **mean vectors are updated**

$$\mathbf{y}^{(i)} = \frac{1}{N_i} \sum_{\mathbf{x} \in \omega_i} \mathbf{x} = \mathcal{C}(\mathbf{X}_i)$$

4. Calculate the total distortion

$$J = \sum_{i=1}^K \sum_{\mathbf{x} \in \omega_i} (\mathbf{x} - \mathbf{y}^{(i)})^2$$

5. Evaluate the convergence. If converged, stop. Else, go to step 2.

# Today Lecture Outline

- Maximum Likelihood Estimation
- Gaussian Mixture Model
- GMM Learning
- General Expectation Maximization Algorithm

# Density Estimation - Motivation

- In our previous lecture on decision theory, we saw that the optimal classifier could be expressed as a family of discriminant functions

$$g_i(\mathbf{x}) = p(\omega_i|\mathbf{x}) = \frac{p(\mathbf{x}|\omega_i)p(\omega_i)}{p(\mathbf{x})}$$

- Decision rule was

choose  $\omega_i$  if  $g_i(\mathbf{x}) > g_j(\mathbf{x}), \forall j \neq i$

- We need to estimate both **prior**  $p(\omega_i)$  and **likelihood**  $p(\mathbf{x}|\omega_i)$
- During next lectures, techniques to estimate the likelihood density function  $p(\mathbf{x}|\omega_i)$  will be introduced

# Approaches for Density Estimation

## Parametric Approach

- Assumes a particular form for the density (e.g., Gaussian) and estimates the parameters of the function (e.g., mean and covariance)
- called Parameter Estimation
  - Maximum Likelihood
  - Bayesian Estimation

## Non-Parametric Approach

- No functional form for the density function is assumed. The density estimate is driven entirely by the data
- called Parameter Estimation
  - Kernel Density Estimation
  - Nearest Neighbor Rule

# Maximum Likelihood Estimation (MLE)

- Consider a pdf  $p(\mathbf{x}|\theta)$  which depends on a set of parameters  $\theta = [\theta_1 \dots \theta_p]^T$
- If examples in the data set are drawn independently from  $p(\mathbf{x}|\theta)$ , the joint probability density of the entire set is given by

$$L(\theta) = p(\mathbf{X}|\theta) = \prod_{i=1} p(\mathbf{x}^{(i)}|\theta)$$

- We seek the set of parameters  $\theta_{ML}$  that maximize the likelihood, and call it the Maximum Likelihood estimate of parameters  $\theta_{ML} = \operatorname{argmax}_{\theta} p(\mathbf{X}|\theta)$
- It is easier to work with the logarithm of the likelihood

$$l(\theta) = \ln p(\mathbf{X}|\theta) = \sum_{i=1}^n \ln p(\mathbf{x}^{(i)}|\theta)$$

$$\nabla_{\theta} l(\theta) = \mathbf{0} \implies \theta_{ML}$$

# Maximum Likelihood Estimation (MLE)

- Has good convergence properties as the sample size increases
- Simpler than any other alternative techniques
- Why to use the logarithm of the likelihood?
- Summary
  - We define  $l(\theta)$  as the log-likelihood function  $l(\theta) = \ln p(\mathbf{X}|\theta)$
  - Solve

$$\nabla_{\theta} l(\theta) = \left[ \frac{\partial l}{\partial \theta_1}, \dots, \frac{\partial l}{\partial \theta_p} \right]^T = \mathbf{0}$$

## MLE Example 1: Gaussian with unknown mean

- The estimation of the parameters of a Gaussian distribution  
 $p(\mathbf{x}^{(i)}|\boldsymbol{\mu}) \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$
- Find the ML estimates
  - Parameters:  $\boldsymbol{\theta} = \boldsymbol{\mu}$

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

$$\ln p(\mathbf{x}^{(i)}|\boldsymbol{\mu}) = -\frac{1}{2} \ln((2\pi)^m |\boldsymbol{\Sigma}|) - \frac{1}{2} (\mathbf{x}^{(i)} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu})$$

- ML estimates must satisfy

$$\nabla_{\boldsymbol{\theta}} \ln p(\mathbf{X}|\boldsymbol{\theta}) = \sum_{i=1}^n \boldsymbol{\Sigma}^{-1} (\mathbf{x}^{(i)} - \boldsymbol{\mu}) = \mathbf{0}$$

$$\boldsymbol{\theta}_{ML} = \hat{\boldsymbol{\mu}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)}$$



## Example 2: Univariate Gaussian with unknown $\mu, \sigma$

- The estimation of the parameters of a Gaussian distribution
- Find the ML estimates for the univariate case
  - Parameters:  $\theta_1 = \mu, \theta_2 = \sigma$
  - Log-likelihood

$$\ln p(x^{(i)}|\theta) = -\frac{1}{2} \ln(2\pi\theta_2^2) - \frac{1}{2\theta_2^2}(x^{(i)} - \theta_1)^2$$

- Gradient

$$\nabla_{\theta} \ln p(x^{(i)}|\theta) = \begin{bmatrix} \frac{1}{\theta_2^2}(x^{(i)} - \theta_1) \\ -\frac{1}{\theta_2} + \frac{(x^{(i)} - \theta_1)^2}{\theta_2^3} \end{bmatrix}$$

- ML estimates

$$\theta_{1,ML} = \hat{\mu} = \frac{1}{n} \sum_{i=1}^n x^{(i)} \quad , \quad \theta_{2,ML}^2 = \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x^{(i)} - \hat{\mu})^2$$

theta\_1: 对i从1到n累加得到,  $\hat{x}(i)$  累加为1.

总共n个theta都相等, 因此乘以系数n.

### Example 3: Multivariate Gaussian with unknown $\mu$ and $\Sigma$

多变量:  $\mu, \Sigma$  是多维向量和矩阵.

Find the ML estimates for the multivariate Gaussian

- Parameters:  $\theta_1 = \mu$ ,  $\theta_2 = \Sigma$
- ML (Maximum-likelihood) estimates for  $\mu$  and  $\sigma$

$$\theta_{1,ML} = \hat{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}^{(i)}$$

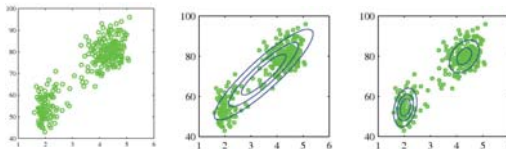
$$\theta_{2,ML} = \hat{\Sigma} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}^{(i)} - \hat{\mu})(\mathbf{x}^{(i)} - \hat{\mu})^T$$

# Mixture models

Consider the problem of modeling a pdf given a dataset of examples

$$X = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(n)}\}$$

- If the form of the underlying pdf is known (e.g. Single Gaussian distribution), the problem could be solved using the Maximum Likelihood Estimation method



Old Faithful data from Bishop2006

- Now we will consider an alternative density estimation method which is modeling the pdf with a mixture of parametric densities. In particular, we will focus on **mixture models of Gaussian densities**

$$p(\mathbf{x}|\boldsymbol{\theta}) = \sum_{j=1}^K p(\mathbf{x}|\boldsymbol{\theta}_j)p(\omega_j) = \sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$

$$p(\omega_j) = \pi_j$$

# Gaussian Mixture Model (GMM)

- Mixture of Gaussians
  - A superposition of K Gaussian densities  $p(\mathbf{x}) = \sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$
  - Parameters  $\pi_j, \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j$
  - Properties: Asymmetry, multi-modality  $0 \leq \pi_j \leq 1, \quad \sum_{j=1}^K \pi_j = 1$
- Previously, we estimated parameters for a single Gaussian distribution by MLE
- Log-likelihood function

$$l(\boldsymbol{\theta}) = \ln p(\mathbf{X} | \boldsymbol{\pi}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{i=1}^n \left[ \ln \left[ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}^{(i)} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \right] \right]$$

# Gaussian Mixture Model (GMM)

Log-likelihood function

$$l(\theta) = \ln p(X|\pi, \mu, \Sigma) = \sum_{i=1}^n \left[ \ln \left[ \sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}^{(i)} | \mu_k, \Sigma_k) \right] \right]$$

Find the maximum of this function by differentiation  
for  $\Sigma_k = \sigma_k^2 \mathbf{I}$

$$\begin{aligned} \frac{\partial l}{\partial \mu_j} = 0 &\rightarrow \hat{\mu}_j = \frac{\sum_{i=1}^n p(\omega_j | \mathbf{x}^{(i)}, \theta) \mathbf{x}^{(i)}}{\sum_{i=1}^n p(\omega_j | \mathbf{x}^{(i)}, \theta)} \\ \frac{\partial l}{\partial \sigma_j} = 0 &\rightarrow \hat{\sigma}_j^2 = \frac{\sum_{i=1}^n p(\omega_j | \mathbf{x}^{(i)}, \theta) (\mathbf{x}^{(i)} - \mu_j)^2}{\sum_{i=1}^n p(\omega_j | \mathbf{x}^{(i)}, \theta)} \\ \frac{\partial l}{\partial \pi_j} = 0 &\rightarrow \hat{\pi}_j = \frac{\sum_{i=1}^n p(\omega_j | \mathbf{x}^{(i)}, \theta)}{\sum_{k=1}^K \sum_{i=1}^n p(\omega_k | \mathbf{x}^{(i)}, \theta)} \end{aligned}$$

# Gaussian Mixture Model (GMM)

- NOT a closed form analytical solution for GMM parameters
- Due to responsibility depends on the GMM parameters
- Highly non-linear coupled system of equations

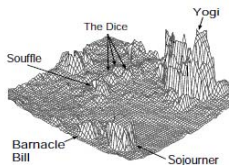
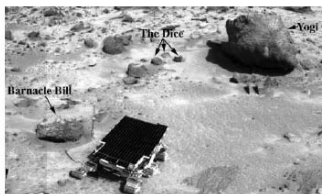
⇒ Iterative Numerical Optimization Technique is necessary. **EM algorithm**

# Maximum Likelihood technique for a rover localization

漫游者

**Task:** To perform rover localization by matching range maps.

**Motivation:** For greater autonomy in Mars rovers.



(Left) Annotated image, (Right) Terrain map generated from stereo image

全景的

- Global Map: panoramic imagery generated at the center of the area from lander.
- Local Map: Occupancy map of local terrain is generated using stereo vision on Sojourner.



Olson, Clark F., and Larry H. Matthies. *Maximum likelihood rover localization by matching range maps*. IEEE International Conference on Robotics and Automation. 1998.



# Maximum Likelihood technique for a rover localization

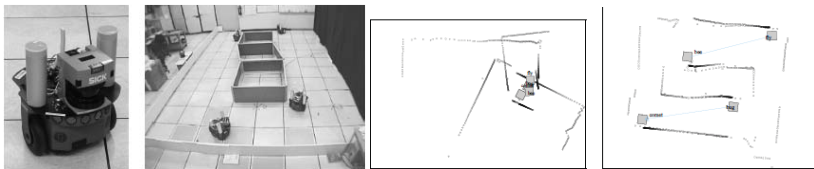
**Task:** To perform rover localization by matching range maps

- Rover position ( $x$  and  $y$  coordinates) denoted as  $X$ .
- Distances of nearby voxels ( $n$  occupied voxels) denoted as  $D_1, D_2, \dots, D_n$ . 立体像素
- The position  $X$  yielding the maximum likelihood value i.e.  $\ln L(X) = l(X) = \sum_{i=1}^n \ln p(D_i|X)$  is chosen to be the position of the rover.
- $p(D_i|X)$  is a normal distribution with a constant additive term. Normal distribution models difference occupied voxels in btw global map and local map
- Results: Comparison of rover position determined by a human operator and by the proposed localization method. Similar results.



# Localization for Mobile Robot Teams

**Task:** to infer the relative pose of every robot in the team without the use of GPS, external landmarks or instrumentation of the environment

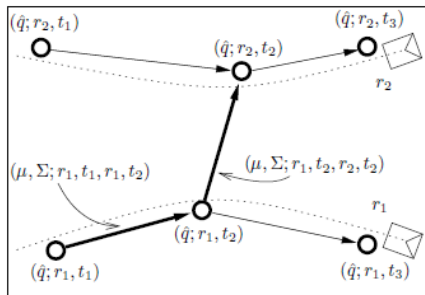


A. Howard, M. Mataric, G. Sukhatme. *Localization for Mobile Robot Teams Using Maximum Likelihood Estimation*. IROS 2002.

# Localization for Mobile Robot Teams

**Task:** to infer the relative pose of every robot in the team

- $r_1, r_2$  : robots 1 and 2.
- Nodes: robot pose estimates.
- Arcs: observations (motion observation and robot observation).



**Aim:** Maximize  $P(O|H)$  using Maximum Likelihood Estimation,  $O$ : observation set (motion sensor and robot sensor),  $H$ : robot pose estimation set.

# Localization for Mobile Robot Teams

**Task:** to infer the relative pose of every robot in the team

- experimental snapshots
- At  $t=1$ , the relative robot pose is completely unknown
- By time  $t=12$  sec, both robots following the outer wall have observed both robots following inner wall. The two robots on the outer wall can correctly determine each other pose even though they have never seen each other.
- errors: At  $t=0$ , the localization error is high. By the time  $t=20$ , the robot performs stable localization, The average range error is about 5.5cm.

