

Parametric Estimation: Estimators, MLE & Mixture Models



Prithwijit Guha
Dept. of EEE, IIT Guwahati

Unsupervised Learning



Parametric
Clustering
Algorithms

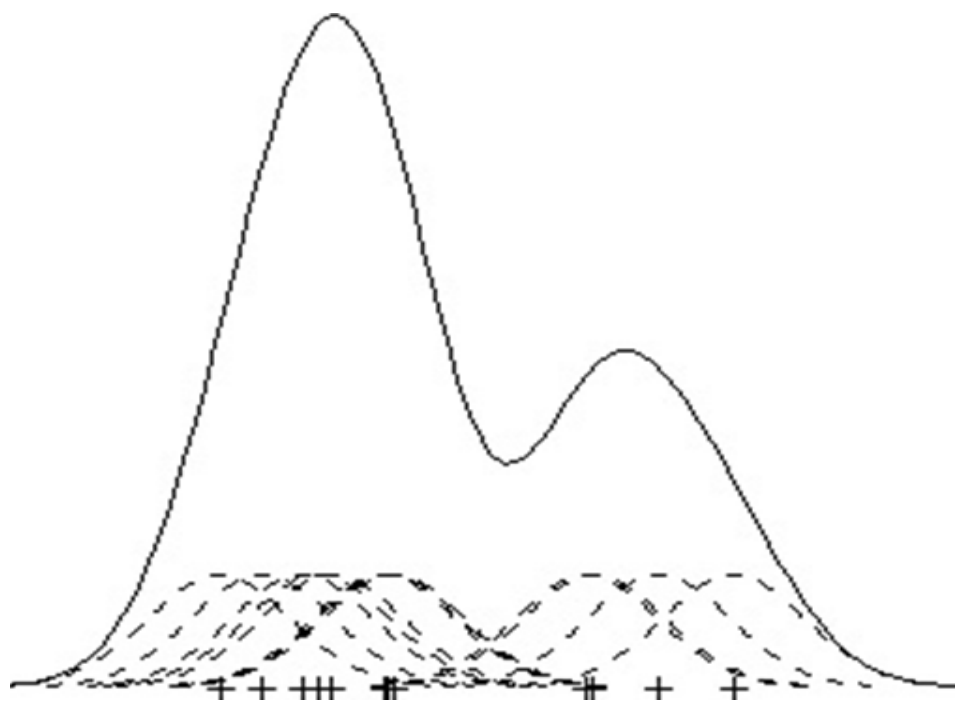
Generic
Clustering
Algorithms

**Estimation
Theory**

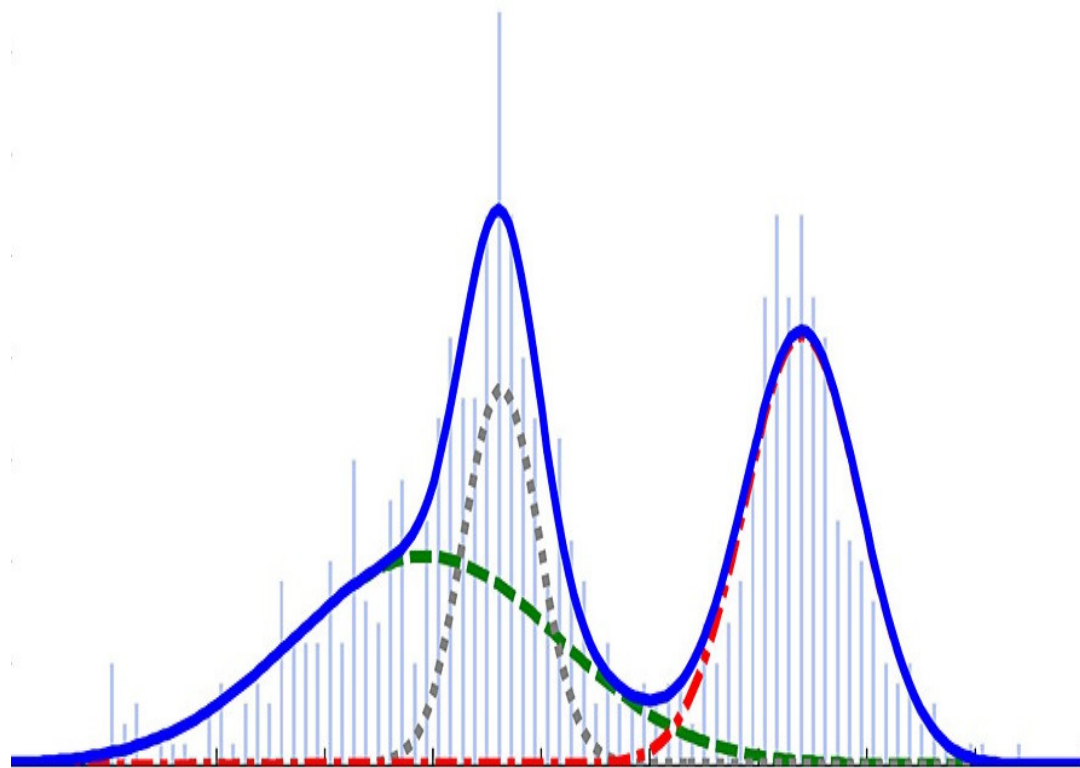
Generative
Models

Pattern
Mining

Estimation Theory



Non-Parametric Estimation



Parametric Estimation

Estimation Theory

- Introduction to Estimators ($\hat{\theta}$)
- Bias and Variance
- Analysis of Mean ($\hat{\mu}$) and Variance ($\hat{\sigma}^2$) Estimators
- Maximum Likelihood Estimation
- Learning Mixture Models
- Gaussian Mixture Models





Estimators

An Estimator is a Rule for Computing an Estimate of a certain Quantity from an Observed Data

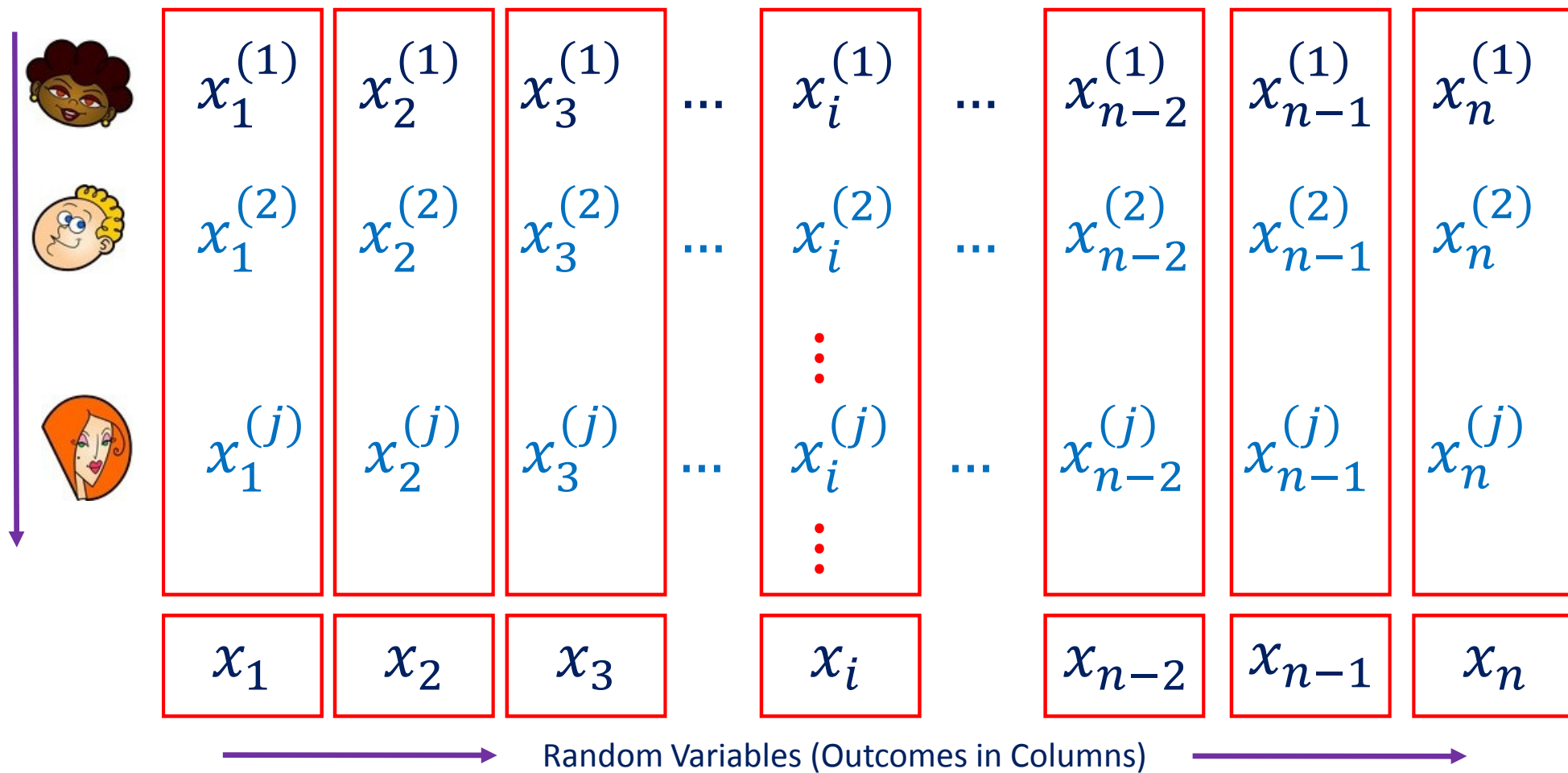
An Estimator is a Statistic that Estimates some Fact about the Population. Estimator can be seen as a Rule to Create the Estimate.

Sample Mean is an Estimate of Population Mean. The Quantity being Estimated is called the Estimand.

Experiments & Samples

	$x_1^{(1)}$	$x_2^{(1)}$	$x_3^{(1)}$...	$x_i^{(1)}$...	$x_{n-2}^{(1)}$	$x_{n-1}^{(1)}$	$x_n^{(1)}$
	$x_1^{(2)}$	$x_2^{(2)}$	$x_3^{(2)}$...	$x_i^{(2)}$...	$x_{n-2}^{(2)}$	$x_{n-1}^{(2)}$	$x_n^{(2)}$
	$x_1^{(3)}$	$x_2^{(3)}$	$x_3^{(3)}$...	$x_i^{(3)}$...	$x_{n-2}^{(3)}$	$x_{n-1}^{(3)}$	$x_n^{(3)}$
					\vdots				
	$x_1^{(j)}$	$x_2^{(j)}$	$x_3^{(j)}$...	$x_i^{(j)}$...	$x_{n-2}^{(j)}$	$x_{n-1}^{(j)}$	$x_n^{(j)}$
					\vdots				

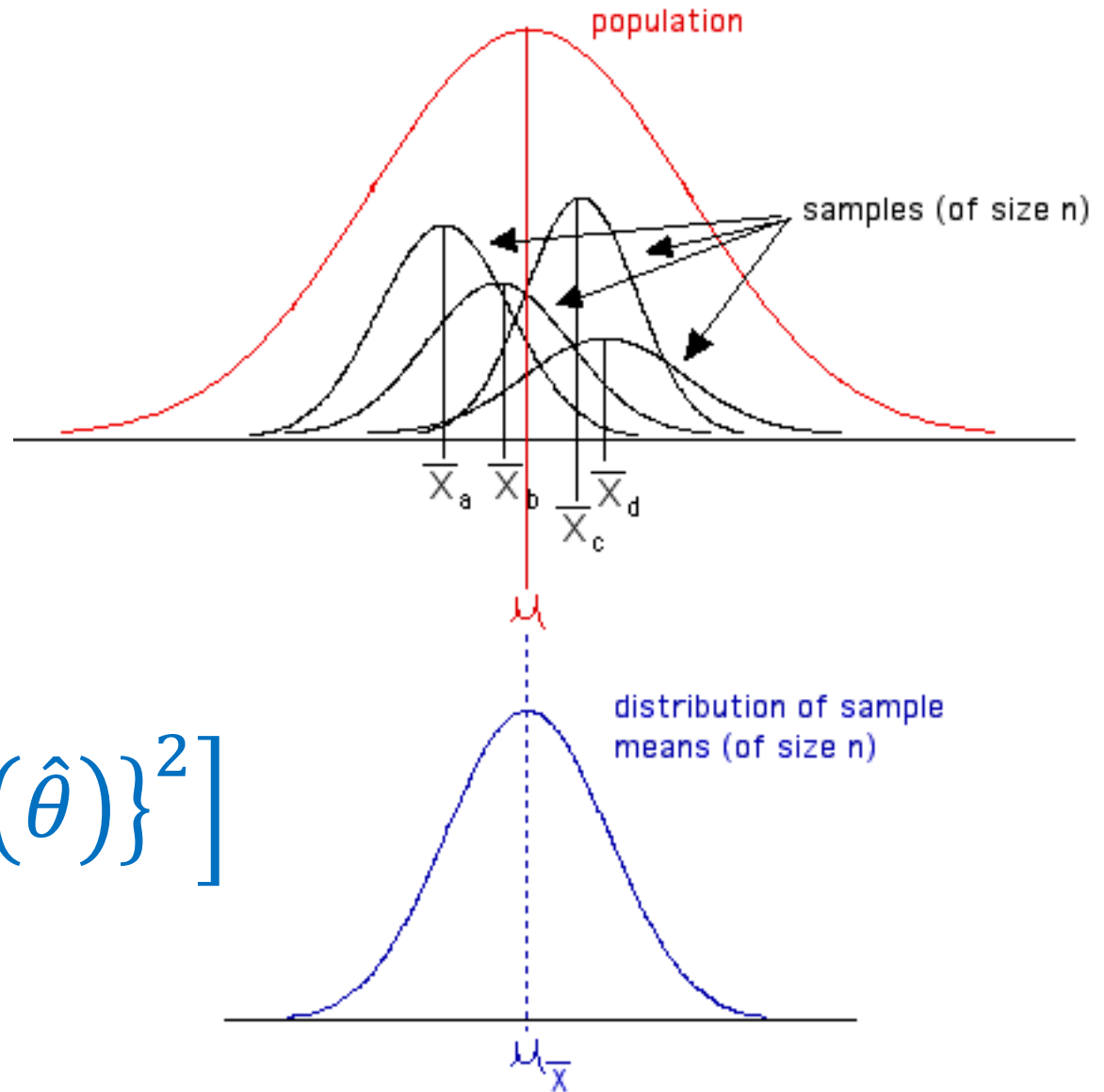
Experiments & Samples



Bias & Variance

$$\text{Bias}(\hat{\theta}) = \theta - E(\hat{\theta})$$

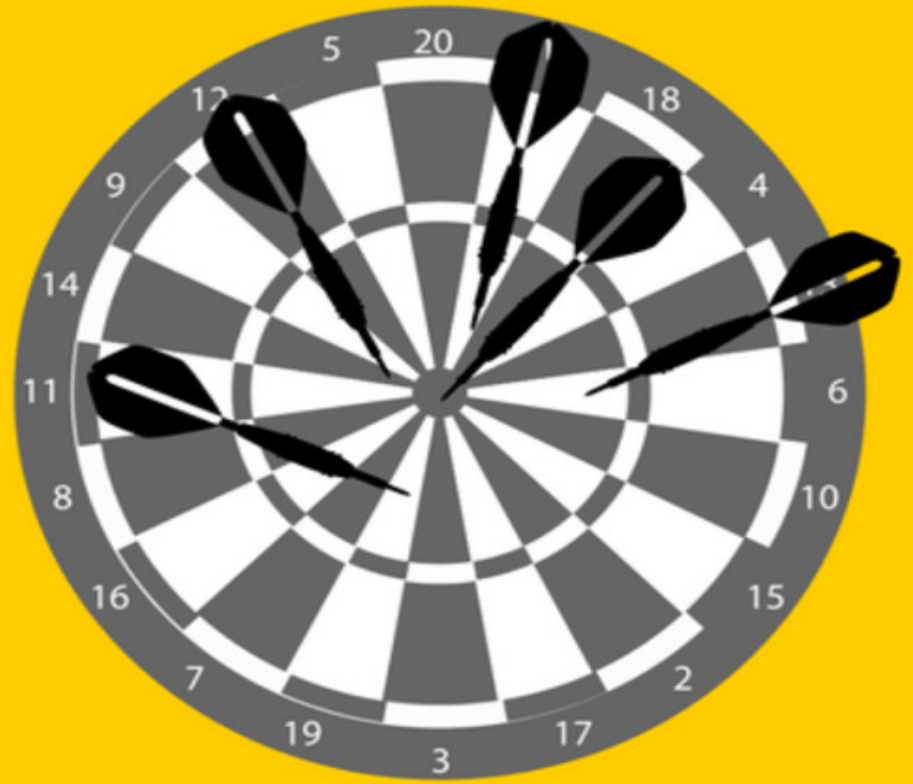
$$\text{Var}(\hat{\theta}) = E \left[\{\hat{\theta} - E(\hat{\theta})\}^2 \right]$$



High Bias
Low Variance



High Variance
Low Bias



Population Mean & Variance

$$\{x_1, x_2, \dots x_i, \dots x_{n-1}, x_n\}$$

Samples Drawn in Any Experiment

$$\forall i \ E(x_i) = E(x) = \mu$$

$$\forall i \ E(x_i^2) = E(x^2)$$

$$\sigma^2 = E[(x - \mu)^2] = E(x^2) - \mu^2$$

Sample Mean & Variance Estimator

$$\{x_1, x_2, \dots x_i, \dots x_{n-1}, x_n\}$$

Samples Drawn in Any Experiment

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\hat{v} = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \hat{\mu}^2$$

Bias: Mean Estimator

$$\text{Bias}(\hat{\mu}) = \mu - E(\hat{\mu})$$

$$E(\hat{\mu}) = E\left(\frac{1}{n} \sum_{i=1}^n x_i\right) = \frac{1}{n} \sum_{i=1}^n E(x_i) = \frac{1}{n} \times (n\mu) = \mu$$

$$\text{Bias}(\hat{\mu}) = 0$$

Sample Mean ($\hat{\mu}$) is an Unbiased Estimator

Variance: Mean Estimator

$$\text{Var}(\hat{\mu}) = E[(\hat{\mu} - E(\hat{\mu}))^2] = E[(\hat{\mu} - \mu)^2] = E(\hat{\mu}^2) - \mu^2$$

$$E(\hat{\mu}^2) = E\left[\left\{\frac{1}{n}\sum_{i=1}^n x_i\right\}\left\{\frac{1}{n}\sum_{j=1}^n x_j\right\}\right] = \frac{1}{n^2}E\left[\left\{\sum_{i=1}^n x_i^2\right\} + \left\{\sum_{\substack{i,j \\ i \neq j}} x_i x_j\right\}\right]$$

$$E(\hat{\mu}^2) = \frac{1}{n^2}\left[\sum_{i=1}^n E(x_i^2) + \sum_{\substack{i,j \\ i \neq j}} E(x_i)E(x_j)\right]$$

Variance: Mean Estimator

$$E(\hat{\mu}^2) = \frac{1}{n^2} \left[\sum_{i=1}^n E(x_i^2) + \sum_{\substack{i,j \\ i \neq j}} E(x_i)E(x_j) \right] = \frac{nE(x^2) + (n^2 - n)\mu^2}{n^2}$$

$$E(\hat{\mu}^2) = \frac{E(x^2) - \mu^2}{n} + \mu^2$$

$$E(\hat{\mu}^2) = \frac{\sigma^2}{n} + \mu^2$$

Variance: Mean Estimator

$$E(\hat{\mu}^2) = \frac{\sigma^2}{n} + \mu^2$$

$$Var(\hat{\mu}) = E(\hat{\mu}^2) - \mu^2 = \frac{\sigma^2}{n}$$

$$\lim_{n \rightarrow \infty} Var(\hat{\mu}) = \lim_{n \rightarrow \infty} \frac{\sigma^2}{n} = 0$$

Variance of Sample Mean
Estimator Becomes Smaller
for Larger Sample Sizes



Bias: Variance Estimator

$$\text{Bias}(\hat{v}) = \sigma^2 - E(\hat{v})$$

$$E(\hat{v}) = E\left[\frac{1}{n} \sum_{i=1}^n x_i^2 - \hat{\mu}^2\right] = \frac{1}{n} \sum_{i=1}^n E(x_i^2) - E(\hat{\mu}^2)$$

$$E(\hat{v}) = \frac{1}{n} \times \{nE(x^2)\} - \left(\frac{\sigma^2}{n} + \mu^2\right) = \{E(x^2) - \mu^2\} - \frac{\sigma^2}{n}$$

$$E(\hat{v}) = \sigma^2 - \frac{\sigma^2}{n} = \left(\frac{n-1}{n}\right) \sigma^2$$

Bias: Variance Estimator

$$E(\hat{v}) = \left(\frac{n-1}{n} \right) \sigma^2$$

$$\text{Bias}(\hat{v}) = \sigma^2 - E(\hat{v}) = \frac{\sigma^2}{n}$$

$$\lim_{n \rightarrow \infty} \text{Bias}(\hat{v}) = \lim_{n \rightarrow \infty} \frac{\sigma^2}{n} = 0$$

Sample Variance (\hat{v}) is an Asymptotically Unbiased Estimator

Bias: Variance Estimator

$$E(\hat{v}) = \left(\frac{n-1}{n}\right) \sigma^2 \Rightarrow \left(\frac{n}{n-1}\right) E(\hat{v}) = E\left(\frac{n}{n-1} \hat{v}\right) = \sigma^2$$

$$\hat{v}_{new} = \frac{n}{n-1} \hat{v} = \frac{n}{n-1} \left[\frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu})^2 \right] = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu})^2$$

$$Bias(\hat{v}_{new}) = \sigma^2 - E(\hat{v}_{new}) = 0$$

Sample Variance (\hat{v}_{new}) is an Unbiased Estimator

Likelihood & Log-Likelihood

$$\mathbf{D} = \{x_1, x_2, \dots x_i, \dots x_{n-1}, x_n\}$$

i.i.d.



$$P(\mathbf{D} \mid \boldsymbol{\theta}) = \prod_{i=1}^n P(x_i \mid \boldsymbol{\theta})$$

$$\mathbf{L}(\boldsymbol{\theta}) = \ln P(\mathbf{D} \mid \boldsymbol{\theta}) = \sum_{i=1}^n \ln P(x_i \mid \boldsymbol{\theta})$$

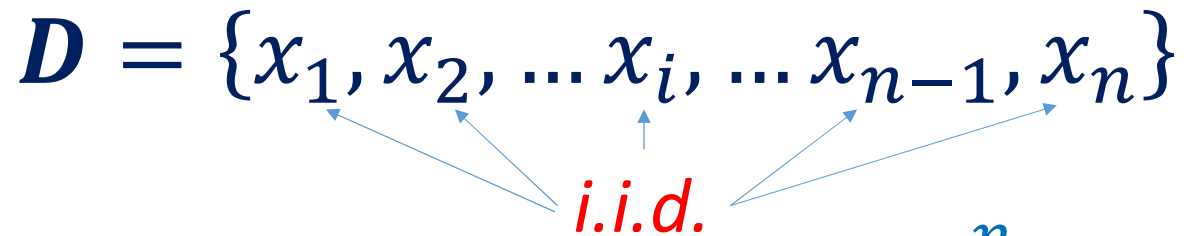
Parameter Estimation: Gaussian Distribution

$$\mathbf{D} = \{x_1, x_2, \dots x_i, \dots x_{n-1}, x_n\}$$

The dataset \mathbf{D} is drawn from a Gaussian Distribution with mean μ and variance $v = \sigma^2$. The most likely distribution parameters $\hat{\mu}_{MLE}$ and \hat{v}_{MLE} need to be estimated from available data.

$$P(x \mid \mu, v) = \frac{1}{\sqrt{2\pi v}} e^{-\frac{(x-\mu)^2}{2v}}$$

Parameter Estimation: Gaussian Distribution

$$\mathbf{D} = \{x_1, x_2, \dots x_i, \dots x_{n-1}, x_n\}$$


A diagram illustrating the independent and identically distributed (i.i.d.) property. The text *i.i.d.* is written in red. Five blue arrows point from this text to the individual data points $x_1, x_2, x_i, x_{n-1},$ and x_n in the set \mathbf{D} , indicating that each data point is drawn from the same distribution independently of the others.

$$L(\boldsymbol{\theta}) = \ln P(\mathbf{D} \mid \mu, \nu) = \sum_{i=1}^n \ln P(x_i \mid \mu, \nu)$$

$$P(x \mid \mu, \nu) = \frac{1}{\sqrt{2\pi\nu}} e^{-\frac{(x-\mu)^2}{2\nu}}$$

Parameter Estimation: Gaussian Distribution

$$L(\mu, v) = \sum_{i=1}^n \ln \frac{1}{\sqrt{2\pi v}} e^{-\frac{(x_i - \mu)^2}{2v}} = -\frac{n}{2} \ln 2\pi - \frac{n}{2} \ln v - \frac{1}{2v} \sum_{i=1}^n (x_i - \mu)^2$$

$$\frac{\partial L(\mu, v)}{\partial \mu} = -\frac{1}{2v} \sum_{i=1}^n 2(x_i - \mu)(-1) = \frac{1}{v} \sum_{i=1}^n (x_i - \mu)$$

$$\frac{\partial L(\mu, v)}{\partial \mu} = 0$$



$$\hat{\mu}_{MLE} = \frac{1}{n} \sum_{i=1}^n x_i$$

Parameter Estimation : Gaussian Distribution

$$L(\mu, v) = \sum_{i=1}^n \ln \frac{1}{\sqrt{2\pi v}} e^{-\frac{(x_i - \mu)^2}{2v}} = -\frac{n}{2} \ln 2\pi - \frac{n}{2} \ln v - \frac{1}{2v} \sum_{i=1}^n (x_i - \mu)^2$$

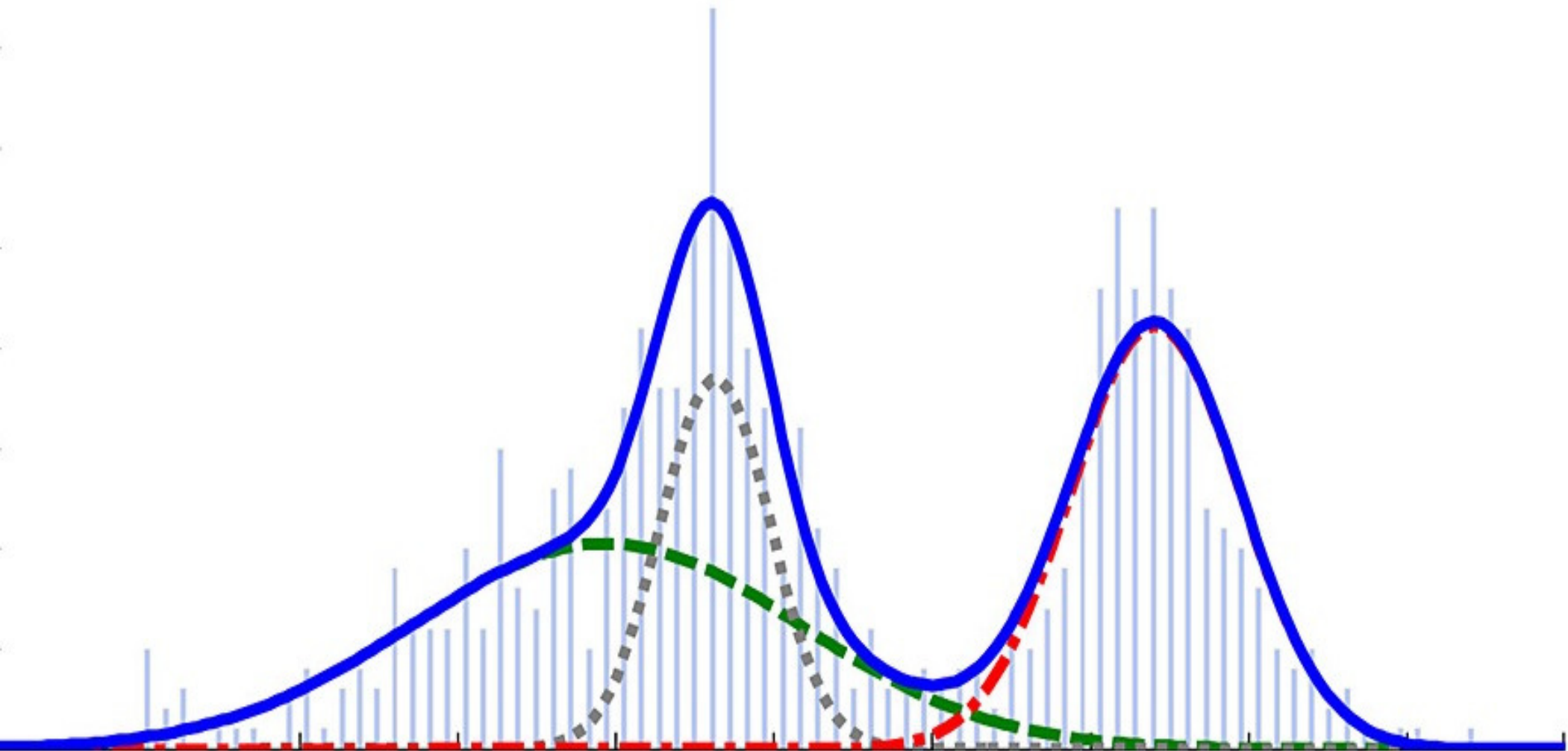
$$\frac{\partial L(\mu, v)}{\partial v} = -\frac{n}{2v} + \frac{1}{2v^2} \sum_{i=1}^n (x_i - \mu)^2$$

$$\frac{\partial L(\mu, v)}{\partial \mu} = 0$$



$$\hat{v}_{MLE} = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}_{MLE})^2$$

Learning Mixture Models



The Mixture Model

$$\mathbf{D} = \{x_1, x_2, \dots x_i, \dots x_{n-1}, x_n\}$$

i.i.d.

The diagram illustrates the relationship between the data set \mathbf{D} and the independent and identically distributed (i.i.d.) assumption. The data set is represented as a set of samples $\{x_1, x_2, \dots, x_i, \dots, x_{n-1}, x_n\}$. The label *i.i.d.* is positioned below the ellipsis, with blue arrows pointing from it to each of the individual samples $x_1, x_2, x_i, x_{n-1},$ and x_n , indicating that each sample is drawn independently from the same distribution.

$$P(x \mid \boldsymbol{\theta}) = \sum_{j=1}^m P(x \mid \boldsymbol{\omega}_j, \boldsymbol{\theta}_j) P(\boldsymbol{\omega}_j)$$

Evaluating Gradient of $L(\boldsymbol{\theta})$

$$L(\boldsymbol{\theta}) = \sum_{i=1}^n \ln P(x_i | \boldsymbol{\theta}) = \sum_{i=1}^n \ln \left\{ \sum_{j=1}^m P(x_i | \boldsymbol{\omega}_j, \boldsymbol{\theta}_j) P(\boldsymbol{\omega}_j) \right\}$$

$$\nabla_{\boldsymbol{\theta}_r} L(\boldsymbol{\theta}) = \sum_{i=1}^n \frac{P(\boldsymbol{\omega}_r)}{P(x_i | \boldsymbol{\theta})} \nabla_{\boldsymbol{\theta}_r} P(x_i | \boldsymbol{\omega}_r, \boldsymbol{\theta}_r)$$

Assumption: $\boldsymbol{\theta}_j$ and $\boldsymbol{\theta}_k$ are Functionally Independent

Introducing the Posterior

$$P(\boldsymbol{\omega}_j \mid x_k, \boldsymbol{\theta}) = \frac{P(x_k \mid \boldsymbol{\omega}_j, \boldsymbol{\theta}_j) P(\boldsymbol{\omega}_j)}{P(x_k \mid \boldsymbol{\theta})}$$

$$\nabla_{\boldsymbol{\theta}_r} L(\boldsymbol{\theta}) = \sum_{i=1}^n \frac{P(\boldsymbol{\omega}_r)}{P(x_i \mid \boldsymbol{\theta})} \nabla_{\boldsymbol{\theta}_r} P(x_i \mid \boldsymbol{\omega}_r, \boldsymbol{\theta}_r)$$

$$\nabla_{\boldsymbol{\theta}_r} L(\boldsymbol{\theta}) = \sum_{i=1}^n \frac{P(\boldsymbol{\omega}_r \mid x_i, \boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}_r} P(x_i \mid \boldsymbol{\omega}_r, \boldsymbol{\theta}_r)}{P(x_i \mid \boldsymbol{\omega}_r, \boldsymbol{\theta}_r)}$$

Gradient of Log-Likelihood

$$\nabla_{\boldsymbol{\theta}_r} L(\boldsymbol{\theta}) = \sum_{i=1}^n \frac{P(\boldsymbol{\omega}_r | x_i, \boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}_r} P(x_i | \boldsymbol{\omega}_r, \boldsymbol{\theta}_r)}{P(x_i | \boldsymbol{\omega}_r, \boldsymbol{\theta}_r)}$$

$$\nabla_{\boldsymbol{\theta}_r} L(\boldsymbol{\theta}) = \sum_{i=1}^n P(\boldsymbol{\omega}_r | x_k, \boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}_r} \ln P(x_i | \boldsymbol{\omega}_r, \boldsymbol{\theta}_r)$$

Gaussian Mixture Models (GMM)

$$P(x \mid \boldsymbol{\theta}) = \sum_{j=1}^m P(x \mid \boldsymbol{\omega}_j, \boldsymbol{\theta}_j) P(\boldsymbol{\omega}_j)$$



$$P(x \mid \boldsymbol{\omega}_j, \boldsymbol{\theta}_j) = \frac{1}{\sqrt{2\pi v_j}} e^{-\frac{(x - \mu_j)^2}{2v_j}}$$

GMM: Component Mean

$$\ln P(x_i \mid \omega_j, \theta_j) = -\frac{1}{2} \ln 2\pi - \frac{1}{2} \ln v_j - \frac{(x_i - \mu_j)^2}{2v_j}$$

$$\frac{\partial L(\theta)}{\partial \mu_r} = \sum_{i=1}^n P(\omega_r \mid x_i, \theta) \frac{\partial}{\partial \mu_r} \{\ln P(x_i \mid \omega_r, \theta_r)\}$$

$$\frac{\partial L(\theta)}{\partial \mu_r} = \sum_{i=1}^n P(\omega_r \mid x_i, \theta) \left(\frac{x_i - \mu_r}{v_r} \right)$$

GMM: Component Mean

$$\frac{\partial L(\boldsymbol{\theta})}{\partial \mu_r} = \sum_{i=1}^n P(\boldsymbol{\omega}_r \mid x_i, \boldsymbol{\theta}) \left(\frac{x_i - \mu_r}{v_r} \right)$$

$$\frac{\partial L(\boldsymbol{\theta})}{\partial \mu_r} = 0 \quad \longrightarrow \quad \hat{\mu}_r = \frac{\sum_{i=1}^n x_i P(\boldsymbol{\omega}_r \mid x_i, \boldsymbol{\theta})}{\sum_{i=1}^n P(\boldsymbol{\omega}_r \mid x_i, \boldsymbol{\theta})}$$

GMM: Component Variance

$$\ln P(x_i | \omega_j, \theta_j) = -\frac{1}{2} \ln 2\pi - \frac{1}{2} \ln v_j - \frac{(x_i - \mu_j)^2}{2v_j}$$

$$\frac{\partial L(\theta)}{\partial v_r} = \sum_{i=1}^n P(\omega_r | x_i, \theta) \frac{\partial}{\partial v_r} \{\ln P(x_i | \omega_r, \theta_r)\}$$

$$\frac{\partial L(\theta)}{\partial v_r} = \sum_{i=1}^n P(\omega_r | x_i, \theta) \left\{ -\frac{1}{2v_r} + \frac{(x_i - \mu_r)^2}{2v_r^2} \right\}$$

GMM: Component Variance

$$\frac{\partial L(\boldsymbol{\theta})}{\partial v_r} = \sum_{i=1}^n P(\boldsymbol{\omega}_r \mid x_i, \boldsymbol{\theta}) \left\{ -\frac{1}{2v_r} + \frac{(x_i - \mu_r)^2}{2v_r^2} \right\}$$

$$\frac{\partial L(\boldsymbol{\theta})}{\partial v_r} = 0 \quad \longrightarrow \quad \hat{v}_r = \frac{\sum_{i=1}^n (x_i - \mu_r)^2 P(\boldsymbol{\omega}_r \mid x_i, \boldsymbol{\theta})}{\sum_{i=1}^n P(\boldsymbol{\omega}_r \mid x_i, \boldsymbol{\theta})}$$

GMM: Mean & Variance Update

$$\hat{\mu}_r^{(t+1)} = \frac{\sum_{i=1}^n x_i P(\boldsymbol{\omega}_r \mid x_i, \boldsymbol{\theta}^{(t)})}{\sum_{i=1}^n P(\boldsymbol{\omega}_r \mid x_i, \boldsymbol{\theta}^{(t)})}$$

$$\hat{v}_r^{(t+1)} = \frac{\sum_{i=1}^n \left\{ x_i - \mu_r^{(t)} \right\}^2 P(\boldsymbol{\omega}_r \mid x_i, \boldsymbol{\theta}^{(t)})}{\sum_{i=1}^n P(\boldsymbol{\omega}_r \mid x_i, \boldsymbol{\theta}^{(t)})}$$

$$x \in \mathbb{R}^1$$

GMM: Mean & Covariance Update

$$\hat{\boldsymbol{\mu}}_r^{(t+1)} = \frac{\sum_{i=1}^n \mathbf{x}_i P(\boldsymbol{\omega}_r \mid \mathbf{x}_i, \boldsymbol{\theta}^{(t)})}{\sum_{i=1}^n P(\boldsymbol{\omega}_r \mid \mathbf{x}_i, \boldsymbol{\theta}^{(t)})}$$

$$\hat{\mathbf{C}}_r^{(t+1)} = \frac{\sum_{i=1}^n \left\{ \mathbf{x}_i - \boldsymbol{\mu}_r^{(t)} \right\} \left\{ \mathbf{x}_i - \boldsymbol{\mu}_r^{(t)} \right\}^T P(\boldsymbol{\omega}_r \mid \mathbf{x}_i, \boldsymbol{\theta}^{(t)})}{\sum_{i=1}^n P(\boldsymbol{\omega}_r \mid \mathbf{x}_i, \boldsymbol{\theta}^{(t)})}$$

$$\mathbf{x} \in \mathbb{R}^d$$

GMM: Applications

- Parametric Estimation of Distribution from Data
- Data Analytics using Estimated Distribution
- GMM Likelihoods used for Classification
- Posterioigrams as Features Constructed from Data
- Numerous Applications in Different Domains

Summary

- Introduction to Estimators ($\hat{\theta}$)
- Bias and Variance
- Analysis of Mean ($\hat{\mu}$) and Variance ($\hat{\sigma}^2$) Estimators
- Maximum Likelihood Estimation
- Learning Mixture Models
- Gaussian Mixture Models



Thank You