

ECE 275A: Parameter Estimation I

Method of Moments

Florian Meyer

*Electrical and Computer Engineering Department
University of California San Diego*

Convergence in Probability

- A sequence of random vectors $\{\mathbf{X}_k\}$ with $\mathbf{X}_k \in \mathbb{R}^n$, converges to random vector $\mathbf{X} \in \mathbb{R}^n$ “in probability” if

$$\lim_{k \rightarrow \infty} P(\|\mathbf{X}_k - \mathbf{X}\| \geq \epsilon) = 0, \forall \epsilon > 0$$

- This is typically denoted either as $\mathbf{X}_k \xrightarrow{\text{prob.}} \mathbf{X}$ or as $\text{p-lim}_{k \rightarrow \infty} \mathbf{X}_k = \mathbf{X}$

Carry-Over Property of Convergence in Probability

- Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be a continuous vector-valued function
- For any two vectors \mathbf{x}_k and \mathbf{x} and $\forall \delta > 0$, there exists $\epsilon > 0$ such that we obtain

$$\|\mathbf{x}_k - \mathbf{x}\| < \epsilon \Rightarrow \|f(\mathbf{x}_k) - f(\mathbf{x})\| < \delta$$

- Note that the contrapositive statement is

$$\|f(\mathbf{x}_k) - f(\mathbf{x})\| \geq \delta \Rightarrow \|\mathbf{x}_k - \mathbf{x}\| \geq \epsilon \quad (1)$$

Carry-Over Property of Convergence in Probability

- Let $\{\mathbf{X}_k\}$ be a sequence that converges in probability, i.e.,
 $\text{p-lim}_{k \rightarrow \infty} \mathbf{X}_k = \mathbf{X}$

- Using equation (1), we have

$$\forall \delta > 0, \exists \epsilon > 0, \text{ such that } P(\|f(\mathbf{X}_k) - f(\mathbf{X})\| \geq \delta) \leq P(\|\mathbf{X}_k - \mathbf{X}\| \geq \epsilon)$$

- Hence, from $\text{p-lim}_{k \rightarrow \infty} \mathbf{X}_k = \mathbf{X}$ it follows that

$$\lim_{k \rightarrow \infty} P(\|f(\mathbf{X}_k) - f(\mathbf{X})\| \geq \delta) = 0, \forall \delta > 0,$$

$$\text{i.e., } \text{p-lim}_{k \rightarrow \infty} f(\mathbf{X}_k) = f(\mathbf{X})$$

- This carry-over property of convergence in probability can also be denoted as

$$\mathbf{X}_k \xrightarrow{\text{prob.}} \mathbf{X} \Rightarrow f(\mathbf{X}_k) \xrightarrow{\text{prob.}} f(\mathbf{X})$$

Convergence in Mean Square

- A sequence $\{\mathbf{X}_k\}$ converges to \mathbf{X} “in mean square” if

$$\lim_{k \rightarrow \infty} E(\|\mathbf{X}_k - \mathbf{X}\|^2) = 0$$

- This is typically denoted as $\mathbf{X}_k \xrightarrow{\text{m.s.}} \mathbf{X}$
- By using the Markov inequality, we obtain

$$P(\|\mathbf{X}_k - \mathbf{X}\| \geq \epsilon) = P(\|\mathbf{X}_k - \mathbf{X}\|^2 \geq \epsilon^2) \leq \frac{E(\|\mathbf{X}_k - \mathbf{X}\|^2)}{\epsilon^2}$$

- We have $\lim_{k \rightarrow \infty} P(\|\mathbf{X}_k - \mathbf{X}\| \geq \epsilon) \leq \lim_{k \rightarrow \infty} \frac{E(\|\mathbf{X}_k - \mathbf{X}\|^2)}{\epsilon^2} = 0, \forall \epsilon > 0$ and thus shown that convergence in mean square implies convergence in probability, i.e.,

$$\mathbf{X}_k \xrightarrow{\text{m.s.}} \mathbf{X} \Rightarrow \mathbf{X}_k \xrightarrow{\text{prob.}} \mathbf{X}$$

Method of Moments

- Let $\theta \in \mathbb{R}^p$ be a parameter vector that parameterizes a (statistical family) of model distributions $p_\theta(x) = p(x; \theta)$
- Assume the statistical family is **identifiable**

$$p(x; \theta) = p(x; \theta'), \forall x \Leftrightarrow \theta = \theta'$$

and **well-specified**

$$\exists \theta \in \Theta \text{ such that } p_{\text{true}}(x) = p(x; \theta)$$

where p_{true} is the true distribution

- If the family is identifiable and well-specified, then the truth distribution is uniquely represented by a “true parameter vector” θ_{true} :

$$\exists \theta_{\text{true}} \text{ such that } p_{\text{true}}(x) = p(x; \theta_{\text{true}})$$

- In this case, learning the true distribution is equivalent to learning the true parameter θ_{true}

Method of Moments

- Let $\mu_k(\boldsymbol{\theta}) = E(X^k; \boldsymbol{\theta})$ be the k -th non-central moment ($k = 1, 2, \dots$) of the model distribution $p(x; \boldsymbol{\theta})$
- Define $\mathbf{m}(\cdot): \mathbb{R}^p \rightarrow \mathbb{R}^p$ as a vector-valued function, i.e.,

$$\mathbf{m}(\boldsymbol{\theta}) = \begin{bmatrix} \mu_{k_1}(\boldsymbol{\theta}) \\ \mu_{k_2}(\boldsymbol{\theta}) \\ \vdots \\ \mu_{k_p}(\boldsymbol{\theta}) \end{bmatrix}$$

where the p elements are non-central moments $\mu_{k_i}(\boldsymbol{\theta})$ selected such that the $p \times p$ Jacobian matrix $\frac{\partial}{\partial \boldsymbol{\theta}} \mathbf{m}(\boldsymbol{\theta})$ is nonsingular for all $\boldsymbol{\theta}$

Method of Moments

- Let \mathbf{m}_{true} be the vector of true moments corresponding to the vector of model moments $\mathbf{m}(\boldsymbol{\theta})$
- Under the nonsingular Jacobian matrix assumption, the vector equation

$$\mathbf{m}(\boldsymbol{\theta}) = \mathbf{m}_{\text{true}}$$

is a system of p independent scalar equations that can be solved for $\boldsymbol{\theta}_{\text{true}}$, i.e.,

$$\boldsymbol{\theta}_{\text{true}} = \mathbf{m}^{-1}(\mathbf{m}_{\text{true}})$$

where $\mathbf{m}^{-1}(\cdot)$ is the inverse function of $\mathbf{m}(\cdot)$ which is also continuously differentiable

- However, \mathbf{m}_{true} is unknown

Method of Moments

- Assume that we observe N i.i.d. samples $\{x_1, \dots, x_N\}$ drawn from the true distribution $p_{\text{true}}(x)$
- Let $\hat{\mu}_k^{(N)} = \frac{1}{N} \sum_{i=1}^N x_i^k$ be the k -th non-central sample moment
- Define $\hat{\mathbf{m}}^{(N)}$ be the vector whose elements are the sample moments of the elements (which are moments) of $\mathbf{m}(\cdot)$

$$\hat{\mathbf{m}}^{(N)} = \begin{bmatrix} \hat{\mu}_{k_1}^{(N)} \\ \hat{\mu}_{k_2}^{(N)} \\ \vdots \\ \hat{\mu}_{k_p}^{(N)} \end{bmatrix}$$

- Note that $E(\hat{\mu}_k^{(N)}) = \mu_k(\theta_{\text{true}}) = \mu_{k*}$

Method of Moments

- Assume $\text{Var}(X^k) < \infty, \forall k \in \{k_1, k_2, \dots, k_p\}$

$$\begin{aligned} E((\hat{\mu}_k^{(N)} - \mu_{k*})^2) &= \text{Var}(\hat{\mu}_k^{(N)}) \\ &= \text{Var}\left(\frac{1}{N} \sum_{i=1}^N X_i^k\right) \\ &= \frac{1}{N^2} \sum_{i=1}^N \text{Var}(X_i^k) \\ &= \frac{1}{N} \text{Var}(X^k) \end{aligned}$$

- Thus, we obtain

$$\lim_{N \rightarrow \infty} E((\hat{\mu}_k^{(N)} - \mu_{k*})^2) = 0 \Rightarrow \hat{\mu}_k^{(N)} \xrightarrow{\text{m.s.}} \mu_{k*} \Rightarrow \hat{\mu}_k^{(N)} \xrightarrow{\text{prob.}} \mu_{k*}$$

$$\text{and } \text{p-lim}_{N \rightarrow \infty} \hat{\mathbf{m}}^{(N)} = \mathbf{m}_{\text{true}}$$

Method of Moments

- Now define the **Method of Moments (MOM)** estimate

$$\hat{\theta}^{(N)} = \mathbf{m}^{-1}(\hat{\mathbf{m}}^{(N)})$$

- The *Carry-Over Property of Convergence in Probability* yields consistency, i.e.,

$$\text{p-lim}_{N \rightarrow \infty} \hat{\theta}^{(N)} = \text{p-lim}_{N \rightarrow \infty} \mathbf{m}^{-1}(\hat{\mathbf{m}}^{(N)}) = \mathbf{m}^{-1}(\mathbf{m}_{\text{true}}) = \theta_{\text{true}},$$

- Thus, we have

$$\hat{\theta}^{(N)} \xrightarrow{\text{prob.}} \theta_{\text{true}}$$