

ECE 275A: Parameter Estimation I

Bayesian Estimation – Part I

Florian Meyer

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

- Loss functions $\ell(\hat{\theta}(\mathbf{y}), \theta)$ have the following properties

$$\ell(\hat{\theta}(\mathbf{y}), \theta) \geq 0 \quad \text{with} \quad \ell(\hat{\theta}(\mathbf{y}), \theta) = 0 \Leftrightarrow \hat{\theta}(\mathbf{y}) = \theta$$

- Possible choices of $\ell(\hat{\theta}(\mathbf{y}), \theta)$ include

① Quadratic error: $\ell(\hat{\theta}(\mathbf{y}), \theta) = \|\hat{\theta}(\mathbf{y}) - \theta\|^2$

② Hit-or-miss error: $\ell(\hat{\theta}(\mathbf{y}), \theta) = \lim_{\delta \rightarrow 0} \ell_{\delta}(\hat{\theta}(\mathbf{y}), \theta)$ with

$$\ell_{\delta}(\hat{\theta}(\mathbf{y}), \theta) = \begin{cases} 0 & \|\hat{\theta}(\mathbf{y}) - \theta\| \leq \delta \\ 1 & \text{others} \end{cases}$$

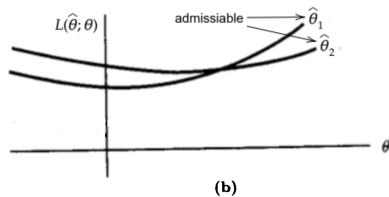
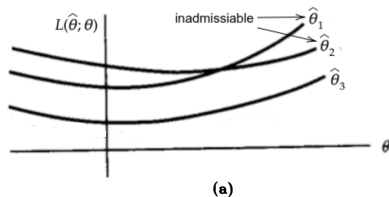
③ Absolute error: $\ell(\hat{\theta}(\mathbf{y}), \theta) = \|\hat{\theta}(\mathbf{y}) - \theta\|_1$

Average Loss and Admissibility of Estimators

- The average loss is given by $L(\hat{\theta}; \theta) = E_{\mathbf{y}}(\ell(\hat{\theta}(\mathbf{y}), \theta)) \geq 0$
- If the quadratic error is used, the average loss is the mean square error $L(\hat{\theta}; \theta) = \text{mse}_{\theta}(\hat{\theta}(\mathbf{y}))$
- An estimator $\hat{\theta}(\mathbf{y})$ is inadmissible if another estimator $\hat{\theta}'(\mathbf{y})$ never has greater average loss than $\hat{\theta}(\mathbf{y})$ but sometimes has strictly lower average loss
- In other words, $\hat{\theta}(\mathbf{y})$ is inadmissible if there exists an $\hat{\theta}'(\mathbf{y})$ such that $L(\hat{\theta}'; \theta) \leq L(\hat{\theta}; \theta)$ for all θ and $L(\hat{\theta}'; \theta') < L(\hat{\theta}; \theta')$ for some θ'

Admissible Estimators – Classical Frequentist Statistics

- In classical frequentist statistics θ is deterministic
- To find the optimal estimator we need to find all admissible estimators
- Recall that an estimator $\hat{\theta}$ that is optimal $\forall \theta$ is typically not realizable



- In (a), estimator $\hat{\theta}_3$ is admissible and optimal for all θ
- In (b), both $\hat{\theta}_1$ and $\hat{\theta}_2$ are admissible, and there is no estimator that is optimal for all θ

Bayesian Statistics

- In Bayesian statistics, both \mathbf{y} and $\boldsymbol{\theta}$ are random
- The joint statistics $p(\mathbf{y}, \boldsymbol{\theta}) = p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta})$ are assumed known
- The Bayes loss $L(\hat{\boldsymbol{\theta}})$ is defined as

$$\begin{aligned} L(\hat{\boldsymbol{\theta}}) &= E_{\mathbf{y}, \boldsymbol{\theta}}(\ell(\hat{\boldsymbol{\theta}}(\mathbf{y}), \boldsymbol{\theta})) \\ &= E_{\boldsymbol{\theta}}\left(E_{\mathbf{y}|\boldsymbol{\theta}}(\ell(\hat{\boldsymbol{\theta}}(\mathbf{y}), \boldsymbol{\theta}) | \boldsymbol{\theta} = \boldsymbol{\theta})\right) \\ &= E_{\boldsymbol{\theta}}(L(\hat{\boldsymbol{\theta}}|\boldsymbol{\theta})) \\ &= \int L(\hat{\boldsymbol{\theta}}|\boldsymbol{\theta})p(\boldsymbol{\theta})d\boldsymbol{\theta} \end{aligned}$$

where $L(\hat{\boldsymbol{\theta}}|\boldsymbol{\theta}) = E_{\mathbf{y}|\boldsymbol{\theta}}(\ell(\hat{\boldsymbol{\theta}}(\mathbf{y}), \boldsymbol{\theta}) | \boldsymbol{\theta} = \boldsymbol{\theta})$

- The Bayes optimal estimator minimizes the Bayes loss, i.e., $\hat{\boldsymbol{\theta}}_* = \arg \min_{\hat{\boldsymbol{\theta}}} L(\hat{\boldsymbol{\theta}})$.

Bayesian Statistics

- The Bayes optimal estimation has the following two properties:
 - ① $\forall p(\theta), \exists$ an admissible $\hat{\theta}$ that is Bayes optimal
 - ② \forall admissible $\hat{\theta}, \exists p(\theta)$ for which $\hat{\theta}$ is Bayes optimal
- Note that $\hat{\theta}_* = \arg \min_{\hat{\theta}} L(\hat{\theta})$ is an optimization in function space
- How to do this tractably?
- $L(\hat{\theta})$ can also be written as

$$\begin{aligned} L(\hat{\theta}) &= E_{\mathbf{y}, \theta}(\ell(\hat{\theta}(\mathbf{y}), \theta)) \\ &= E_{\mathbf{y}}\left(E_{\theta|\mathbf{y}}(\ell(\hat{\theta}(\mathbf{y}), \theta) | \mathbf{y} = \mathbf{y})\right) \\ &= E_{\mathbf{y}}(L(\hat{\theta}|\mathbf{y})) \\ &= \int L(\hat{\theta}|\mathbf{y}) p(\mathbf{y}) d\mathbf{y} \end{aligned}$$

where $L(\hat{\theta}|\mathbf{y}) = E_{\theta|\mathbf{y}}(\ell(\hat{\theta}(\mathbf{y}), \theta) | \mathbf{y} = \mathbf{y})$

$$L(\hat{\theta}) = E_{\theta}(L(\hat{\theta}|\theta)) = E_{\mathbf{y}}(L(\hat{\theta}|\mathbf{y}))$$

$$L(\hat{\theta}) = \int \underbrace{L(\hat{\theta}|\theta)}_{\geq 0} \underbrace{p(\theta)}_{\geq 0} d\theta = \int \underbrace{L(\hat{\theta}|\mathbf{y})}_{\geq 0} \underbrace{p(\mathbf{y})}_{\geq 0} d\mathbf{y}$$

- Either minimizing $L(\hat{\theta}|\theta)$ or $L(\hat{\theta}|\mathbf{y})$ minimizes $L(\hat{\theta})$
- Recall that $\hat{\theta}_* = \arg \min_{\hat{\theta}} L(\hat{\theta})$ requires optimization in function space, but if we know the value of $\hat{\theta}_*(\mathbf{y})$ for each \mathbf{y} , then we know the function $\hat{\theta}(\cdot)$
- Let $\hat{\theta}_*(\mathbf{y}) = \arg \min_{\hat{\theta}} L(\hat{\theta}|\mathbf{y})$, which is a regular vector optimization for each $\mathbf{y} \in \mathbb{Y}$ and $\hat{\theta}_*(\mathbf{y}) \in \mathbb{R}^p$, whereas $\hat{\theta}_* \in \{f : \mathbb{Y} \subset \mathbb{R}^m \rightarrow \mathbb{R}^p\}$

- If the loss function is quadratic error (i.e., $\ell(\hat{\boldsymbol{\theta}}(\mathbf{y}), \boldsymbol{\theta}) = \|\hat{\boldsymbol{\theta}}(\mathbf{y}) - \boldsymbol{\theta}\|^2$), the Bayes loss $L(\hat{\boldsymbol{\theta}})$ becomes the Bayesian mean square error, i.e.,

$$\begin{aligned}\text{bmse}(\hat{\boldsymbol{\theta}}) &= E_{\mathbf{y}, \boldsymbol{\theta}} (\|\hat{\boldsymbol{\theta}}(\mathbf{y}) - \boldsymbol{\theta}\|^2) \\ &= E_{\mathbf{y}, \boldsymbol{\theta}} (\|\tilde{\boldsymbol{\theta}}(\mathbf{y})\|^2) \\ &= E_{\mathbf{y}} \left(\underbrace{E_{\boldsymbol{\theta}|\mathbf{y}} (\|\tilde{\boldsymbol{\theta}}(\mathbf{y})\|^2 | \mathbf{y} = \mathbf{y})}_{L(\hat{\boldsymbol{\theta}}|\mathbf{y})} \right)\end{aligned}$$

- We want $\hat{\boldsymbol{\theta}}_*(\mathbf{y}) = \arg \min_{\hat{\boldsymbol{\theta}}} L(\hat{\boldsymbol{\theta}}|\mathbf{y}) = \arg \min_{\hat{\boldsymbol{\theta}}} E_{\boldsymbol{\theta}|\mathbf{y}} (\|\tilde{\boldsymbol{\theta}}(\mathbf{y})\|^2 | \mathbf{y} = \mathbf{y})$

- To minimize $L(\hat{\boldsymbol{\theta}}|\mathbf{y})$, we take the derivative with respect to $\hat{\boldsymbol{\theta}}$ and set it to 0,

$$\begin{aligned}\nabla_{\hat{\boldsymbol{\theta}}} L(\hat{\boldsymbol{\theta}}|\mathbf{y}) &= E_{\boldsymbol{\theta}|\mathbf{y}}(\nabla_{\hat{\boldsymbol{\theta}}} \|(\tilde{\boldsymbol{\theta}}(\mathbf{y}))\|^2 | \mathbf{y} = \mathbf{y}) \\ &= 2E_{\boldsymbol{\theta}|\mathbf{y}}(\hat{\boldsymbol{\theta}}(\mathbf{y}) - \boldsymbol{\theta} | \mathbf{y} = \mathbf{y}) \stackrel{\text{set}}{=} 0\end{aligned}$$

We get $\hat{\boldsymbol{\theta}}_*(\mathbf{y}) = E_{\boldsymbol{\theta}|\mathbf{y}}(\boldsymbol{\theta}|\mathbf{y})$

- The Hessian $\nabla_{\hat{\boldsymbol{\theta}}}^2 L(\hat{\boldsymbol{\theta}}|\mathbf{y}) = 2\mathbf{I} \Rightarrow \hat{\boldsymbol{\theta}}_*(\mathbf{y})$ is optimal
- Besides $E_{\mathbf{y}}(\hat{\boldsymbol{\theta}}_*(\mathbf{y})) = E_{\mathbf{y}}(E_{\boldsymbol{\theta}|\mathbf{y}}(\boldsymbol{\theta}|\mathbf{y})) = E_{\boldsymbol{\theta}}(\boldsymbol{\theta})$, hence it is unbiased in a Bayesian sense

Example: Quadratic Error

- Consider a scalar observation $y = x + n$, where $x \sim \mathcal{N}(0, \sigma_x^2)$ is the unknown parameter and $n \sim \mathcal{N}(0, \sigma_n^2)$
- It is assumed that x and n are independent $\Rightarrow p(y|x) = \mathcal{N}(x, \sigma_n^2)$
- To find the minimum mean square error (MMSE) estimator $\hat{x}_{\text{MMSE}} = E(x|y)$, we first calculate the posterior distribution $p(x|y)$

$$\begin{aligned} p(x|y) &\propto p(y|x)p(x) \\ &\propto \exp\left(-\frac{(y-x)^2}{2\sigma_n^2} - \frac{x^2}{2\sigma_x^2}\right) \\ &= \exp\left(-\frac{\sigma_x^2 + \sigma_n^2}{2\sigma_x^2\sigma_n^2}x^2 + \frac{1}{\sigma_n^2}yx - \frac{1}{2\sigma_n^2}y^2\right) \\ &\propto \mathcal{N}\left(\frac{\sigma_x^2}{\sigma_x^2 + \sigma_n^2}y, \frac{\sigma_x^2\sigma_n^2}{\sigma_x^2 + \sigma_n^2}\right) \end{aligned}$$

- Thus $\hat{x}_{\text{MMSE}} = E(x|y) = \frac{\sigma_x^2}{\sigma_x^2 + \sigma_n^2}y$, which is different from the ML estimator $\hat{x}_{\text{ML}} = \arg \max_x p(y|x) = y$.