# ECE 275A: Parameter Estimation I
# Classical Estimation

**Florian Meyer**

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

# Statistical Family

- $\boldsymbol{\theta} \in \mathbb{R}^p$: unknown parameter of interest to be estimated (deterministic)

- $\mathbf{y} \in \mathbb{R}^m$: data which depends on the unknown parameter $\boldsymbol{\theta}$ (random)

- To mathematically model the data, which is inherently random, we consider a **statistical family** of parameterized distributions, i.e.,

$$\mathcal{P} = \left\{ p(\mathbf{y}; \boldsymbol{\theta}) \big| \boldsymbol{\theta} \in \Theta \subset \mathbb{R}^p, \mathbf{y} \in \mathbb{Y} \subset \mathbb{R}^m \right\}$$

- Depending on the context, $p(\mathbf{y}; \boldsymbol{\theta})$ is either a probability density function (PDF) or a probability mass function (PMF)

# Estimators

- An estimator $\hat{\boldsymbol{\theta}}(\cdot) : \mathbb{Y} \to \Theta$ is a function that does not depend on the unknown parameter $\boldsymbol{\theta}$; it can be considered as a rule that assigns a value to $\boldsymbol{\theta}$ for *each realization* of **y**

- An estimate $\hat{\theta}(\boldsymbol{y})$ is the value of $\boldsymbol{\theta}$ for a *fixed realization* of **y**

- The primary goal of statistical parameter estimation is to find an estimator $\hat{\boldsymbol{\theta}}(\cdot)$ with the property of providing an estimate $\hat{\theta}(\boldsymbol{y})$ that is *accurate* (close to the true unknown parameter $\boldsymbol{\theta}$) for most parameter values $\boldsymbol{\theta}$ and data realizations **y**

- Another important property is that the estimator $\hat{\boldsymbol{\theta}}(\cdot)$ is *robust* to model mismatch (small changes in $p(\boldsymbol{y}; \boldsymbol{\theta})$ do not severely affect the performance of the estimator $\hat{\boldsymbol{\theta}}(\cdot)$)

# The Mean Square Error

- A natural optimality criterion for estimators is the mean square error ("little mse"),

$$\mathsf{mse}_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}}) = E\big(\|\tilde{\boldsymbol{\theta}}\|^2; \boldsymbol{\theta}\big)$$

where $\tilde{\boldsymbol{\theta}} = \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}$ is the estimation error

- For future reference we also consider the matrix mean square error ("big MSE"),

$$\mathsf{MSE}_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}}) = E\big(\tilde{\boldsymbol{\theta}}\tilde{\boldsymbol{\theta}}^{\mathsf{T}}; \boldsymbol{\theta}\big)$$

- Note that $\mathsf{mse}_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}}) = \mathsf{tr}\big(\mathsf{MSE}_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}})\big)$

- Unfortunately, this natural criterion leads to estimators that cannot be expressed only as a function of the data **y** and thus to unrealizable estimators

# Example: Unrealizability of Mean Square Error Estimation

- Consider the data model $y = A + n$, with scalar parameter $A$ to be estimated and additive Gaussian measurement noise $n \sim \mathcal{N}(0, \sigma^2)$ with known $\sigma^2$

- Consider the estimator $\hat{A} = ay$ based on an arbitrary constant $a \in \mathbb{R}$, i.e., $\hat{A} \sim \mathcal{N}(aA, a^2\sigma^2)$

- We aim to find the $a$ which results in the minimum mean square error

- First, we rewrite the mean square error for a general unknown parameter $\boldsymbol{\theta}$ as

$$
\begin{aligned}
\mathsf{mse}_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}}) &= E\big(\|\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}\|^2\big) \\
&= E\big(\|(\hat{\boldsymbol{\theta}} - E(\hat{\boldsymbol{\theta}})) + (E(\hat{\boldsymbol{\theta}}) - \boldsymbol{\theta})\|^2\big) \\
&= E\big(\|\hat{\boldsymbol{\theta}} - E(\hat{\boldsymbol{\theta}})\|^2\big) + \|E(\hat{\boldsymbol{\theta}}) - \boldsymbol{\theta}\|^2 \\
&= \mathsf{tr}\big(\underbrace{\mathsf{Cov}_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}})}_{variance}\big) + \|\underbrace{E(\hat{\boldsymbol{\theta}}) - \boldsymbol{\theta}}_{bias}\|^2
\end{aligned}
$$

# Example: Unrealizability of Mean Square Error Estimation

- For the scalar $\hat{A}$, the mean square error is obtained as

$$\mathsf{mse}_A\big(\hat{A}\big) = \mathsf{var}(\hat{A}) + \big(E(\hat{A}) - A\big)^2$$
$$= a^2\sigma^2 + (a-1)^2 A^2$$

- Next, we differentiate the mean square error with respect to $a$ and set it to zero, i.e.,

$$\frac{\partial \mathsf{mse}_A\big(\hat{A}\big)}{\partial a} = 2a\sigma^2 + 2(a-1)A^2 \triangleq 0$$

- The optimal value for $a$ is now obtained as

$$a_* = \frac{A^2}{\sigma^2 + A^2}$$

- $a_*$ depends on the unknown parameter $A$ and is thus unrealizable!

# Uniformly Unbiased Estimators (UUBE)

- An alternative approach to obtain realizable estimators is to constrain the bias $E(\hat{\boldsymbol{\theta}}; \boldsymbol{\theta}) - \boldsymbol{\theta}$ to be zero:

## Uniformly Unbiased Estimators (UUBE)

A UUBE $\hat{\boldsymbol{\theta}}$ is an estimator that satisfies

$$E(\hat{\boldsymbol{\theta}}; \boldsymbol{\theta}) = \boldsymbol{\theta}, \forall \boldsymbol{\theta} \in \Theta$$

- Note that if $\hat{\boldsymbol{\theta}}$ is an UUBE, then $E(\tilde{\boldsymbol{\theta}}; \boldsymbol{\theta}) = \mathbf{0}, \forall \boldsymbol{\theta} \in \Theta$

# Uniformly Minimum Variance Unbiased Estimators (UMVUE)

- For all $\hat{\boldsymbol{\theta}}$ that are UUBEs, we have

$$\text{MSE}_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}}) = \text{Cov}_{\boldsymbol{\theta}}(\tilde{\boldsymbol{\theta}}) = \text{Cov}_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}})$$

- Thus, of particular interest, is the UUBE that minimizes the variance

## Uniformly Minimum Variance Unbiased Estimator (UMVUE)

A UMVUE $\hat{\boldsymbol{\theta}}_*$ is an estimator that is defined as follows

$\hat{\boldsymbol{\theta}}_*$ is an UUBE and $\text{Cov}_{\boldsymbol{\theta}}(\tilde{\boldsymbol{\theta}}_*) \preccurlyeq \text{Cov}_{\boldsymbol{\theta}}(\tilde{\boldsymbol{\theta}}), \; \forall \boldsymbol{\theta} \in \Theta, \forall \hat{\boldsymbol{\theta}}$ that are UUBEs

## How to construct a UMVUE?

- General Question: **How to construct a UMVUE?**

- One approach is to work with statistical families for which there exists a uniform lower bound of the error covariance matrix $\boldsymbol{B_\theta}$, i.e.,

$$\boldsymbol{B_\theta} \preccurlyeq \mathrm{Cov}_{\boldsymbol{\theta}}(\tilde{\boldsymbol{\theta}}), \quad \forall \boldsymbol{\theta} \in \Theta, \hat{\boldsymbol{\theta}} \text{ that are UUBEs}$$

- If such a (matrix) lower bound exists and an UUBE $\hat{\boldsymbol{\theta}}'$ can be found such that

$$\mathrm{Cov}_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}}') = \boldsymbol{B_\theta}, \forall \boldsymbol{\theta} \in \Theta$$

then $\hat{\boldsymbol{\theta}}'$ is the UMVUE, i.e., $\hat{\boldsymbol{\theta}}' = \hat{\boldsymbol{\theta}}_*$

- For so-called **Regular Statistical Families (RSF)**, such a uniform lower bound exists and is referred to as the **Cramér-Rao Lower Bound (CRLB)**

# Regular Statistical Families (RSF)

- Recall the definition to a statistical family

$$\mathcal{P} = \left\{ p(\boldsymbol{y}; \boldsymbol{\theta}) \middle| \boldsymbol{\theta} \in \Theta \subseteq \mathbb{R}^p, \boldsymbol{y} \in \mathbb{Y} \subseteq \mathbb{R}^m \right\}$$

---

**Regular Statistical Families (RSF)**

An RSF is a statistical family $\mathcal{P}$ that satisfies the three conditions:

**R1** The support of $p(\boldsymbol{y}; \boldsymbol{\theta})$ independent of the parameter vector $\boldsymbol{\theta}$

**R2** $p(\boldsymbol{y}; \boldsymbol{\theta})$ is differentiable (i.e. $\nabla_{\boldsymbol{\theta}} p(\boldsymbol{y}; \boldsymbol{\theta})$ exists)

**R3** $p(\boldsymbol{y}; \boldsymbol{\theta})$ is doubly-differentiable (i.e. $\nabla_{\boldsymbol{\theta}}^2 p(\boldsymbol{y}; \boldsymbol{\theta})$ exists)

---

- Let us define the score of $\mathcal{P}$ as $\boldsymbol{s}_{\boldsymbol{\theta}}(\mathbf{y}) = \nabla_{\boldsymbol{\theta}} \ln p(\boldsymbol{y}; \boldsymbol{\theta})$
- If $\mathcal{P}$ is an RSF, then $E\left( \boldsymbol{s}_{\boldsymbol{\theta}}(\mathbf{y}) \right) = \boldsymbol{0}$

# Cramér-Rao Lower Bound (CRLB)

- If $\mathcal{P}$ is an RSF and $\hat{\boldsymbol{\theta}}$ is an UUBE, we also have

$$E\left(\boldsymbol{s_\theta}(\mathbf{y})\tilde{\boldsymbol{\theta}}^\mathsf{T}\right) = \boldsymbol{I}$$

- The **Fisher Information Matrix (FIM) $\boldsymbol{J_\theta}$** of the RSF $\mathcal{P}$ is defined as the covariance matrix of the score

$$\boldsymbol{J_\theta} = \mathrm{Cov}_{\boldsymbol{\theta}}\big(\boldsymbol{s_\theta}(\mathbf{y})\big) = E\left(\boldsymbol{s_\theta}(\mathbf{y})\boldsymbol{s_\theta^\mathsf{T}}(\mathbf{y})\right) = -E\left(\nabla_{\boldsymbol{\theta}}^2 \ln p(\mathbf{y};\boldsymbol{\theta})\right)$$

- If the FIM is positive definite, $\boldsymbol{J_\theta}^{-1}$ exists and is equal to the CRLB, i.e.,

$$\mathrm{MSE}_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}}) = \mathrm{Cov}_{\boldsymbol{\theta}}(\hat{\boldsymbol{\theta}}) \succcurlyeq \mathrm{CRLB} = \boldsymbol{J_\theta}^{-1}, \quad \forall\boldsymbol{\theta}\in\Theta, \forall\hat{\boldsymbol{\theta}} \text{ that are UUBEs}$$

with equality iff $\tilde{\boldsymbol{\theta}} = \boldsymbol{J_\theta}^{-1}\boldsymbol{s_\theta}(\mathbf{y})$

## Example: Linear-Gaussian Model

- Let $\mathbf{y} = \boldsymbol{A}\boldsymbol{\theta} + \mathbf{n}$ where the additive noise $\mathbf{n}$ is Gaussian distributed, i.e., $\mathbf{n} \sim \mathcal{N}(\mathbf{0}, \boldsymbol{C})$, and both $\boldsymbol{A} \in \mathbb{R}^{m \times p}$ and $\boldsymbol{C} \in \mathbb{R}^{m \times m}$ are known

- Note that $\mathbf{y} \sim \mathcal{N}(\boldsymbol{A}\boldsymbol{\theta}, \boldsymbol{C})$

- **Goal:** Estimate the unknown parameter $\boldsymbol{\theta}$ from observed data $\mathbf{y}$

- A common choice is the Maximum Likelihood (ML) estimator
  $\hat{\boldsymbol{\theta}}_{\mathrm{ML}} = \arg\max_{\boldsymbol{\theta}} p(\mathbf{y}; \boldsymbol{\theta})$

- Assuming that $\boldsymbol{A}$ is a tall matrix and both $\boldsymbol{C}$ and $\boldsymbol{A}$ are full rank, the ML estimator is given by

$$\hat{\boldsymbol{\theta}}_{\mathrm{ML}} = \arg\max_{\boldsymbol{\theta}} \ln p(\mathbf{y}; \boldsymbol{\theta}) = (\boldsymbol{A}^\mathsf{T} \boldsymbol{C}^{-1} \boldsymbol{A})^{-1} \boldsymbol{A}^\mathsf{T} \boldsymbol{C}^{-1} \mathbf{y}$$

- It is easy to verify that $E(\hat{\boldsymbol{\theta}}_{\mathrm{ML}}; \boldsymbol{\theta}) = \boldsymbol{\theta} \Rightarrow \hat{\boldsymbol{\theta}}_{\mathrm{ML}}$ is an UUBE
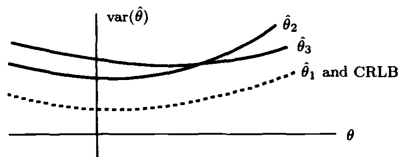
## Example: Linear-Gaussian Model

- Next, we calculate the score and FIM as follows

$$s_\theta(\mathbf{y}) = \mathbf{A}^\mathsf{T} \mathbf{C}^{-1}(\mathbf{y} - \mathbf{A}\theta)$$
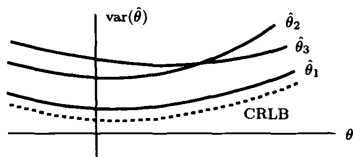$$\mathbf{J}_\theta = \mathbf{A}^\mathsf{T} \mathbf{C}^{-1} \mathbf{A}$$

- Since the ML estimation error can be expressed as $\tilde{\theta}_{\mathrm{ML}} = \mathbf{J}_\theta^{-1} s_\theta(\mathbf{y})$, the ML estimator $\hat{\theta}_{\mathrm{ML}}$ can attain the CRLB

- **The MLE for linear-Gaussian models is also the UMVUE!**

# Efficient Estimators

- An estimator which is unbiased and attains the CRLB is said to be *efficient* in that it efficiently uses the data

- An UMVUE may or may not be efficient [*Kay, 1993*] :



(a) $\hat{\theta}_1$ efficient and MVU      (b) $\hat{\theta}_1$ MVU but not efficient

- In (a), the UMVU $\hat{\theta}_1$ is efficient in that it attains the CRLB

- In (b), $\hat{\theta}_1$ is the UMVU but does not attain the CRLB, and hence it is not efficient

# Best Linear Unbiased Estimator (BLUE)

- Again we consider the linear model $\mathbf{y} = \boldsymbol{A\theta} + \mathbf{n}$, where the additive noise $\mathbf{n}$ has zero-mean and covariance $\boldsymbol{C}$ but does not have to be Gaussian

- Furthermore, the statistical family $\mathcal{P}$ induced by $p(\mathbf{y}; \boldsymbol{\theta})$ does not have to be an RSF

- Note that $E(\mathbf{y}; \boldsymbol{\theta}) = \boldsymbol{A\theta}$ and $\text{Cov}_{\boldsymbol{\theta}}(\mathbf{y}) = \boldsymbol{C}$

- **Linear Unbiased Estimator (LUE)**: We aim to develop an UUBE $\hat{\boldsymbol{\theta}}$ that is a linear function of the data (i.e., $\hat{\boldsymbol{\theta}} = \boldsymbol{K}\mathbf{y}$)

# Best Linear Unbiased Estimator (BLUE)

- For $\hat{\boldsymbol{\theta}}$ unbiased and linear

$$\Rightarrow E(\hat{\boldsymbol{\theta}}) = E(\boldsymbol{K}\mathbf{y}) = E(\boldsymbol{K}(\boldsymbol{A}\boldsymbol{\theta} + \mathbf{n})) = \boldsymbol{K}\boldsymbol{A}\boldsymbol{\theta} = \boldsymbol{\theta}$$
$$\Rightarrow \boldsymbol{K}\boldsymbol{A} = \boldsymbol{I}$$
$$\Rightarrow \tilde{\boldsymbol{\theta}} = \boldsymbol{K}\mathbf{n}$$

- Thus, the MSE can be obtained as

$$\mathrm{MSE}_\theta(\hat{\boldsymbol{\theta}}) = \mathrm{Cov}_\theta(\hat{\boldsymbol{\theta}}) = \mathrm{Cov}_\theta(\tilde{\boldsymbol{\theta}}) = E(\boldsymbol{K}\mathbf{n}\mathbf{n}^\mathsf{T}\boldsymbol{K}^\mathsf{T}) = \boldsymbol{K}\boldsymbol{C}\boldsymbol{K}^\mathsf{T}$$

- **Gauss-Markov Theorem (GMT):** The **Best Linear Unbiased Estimator (BLUE)** is given by

$$\hat{\boldsymbol{\theta}}_o = \boldsymbol{K}_o\mathbf{y} \qquad \boldsymbol{K}_o = (\boldsymbol{A}^T\boldsymbol{C}^{-1}\boldsymbol{A})^{-1}\boldsymbol{A}^T\boldsymbol{C}^{-1} \qquad (1)$$

## Proof of the Gauss-Markov Theorem

**Proof:**

- Based on Eq. (1) the $\text{MSE}_{\theta}(\hat{\boldsymbol{\theta}}_o)$ can be obtained as

$$\begin{aligned}
\text{MSE}_{\theta}(\hat{\boldsymbol{\theta}}_o) &= \boldsymbol{K}_o \boldsymbol{C} \boldsymbol{K}_o^{\mathsf{T}} \\
&= (\boldsymbol{A}^T \boldsymbol{C}^{-1} \boldsymbol{A})^{-1} \boldsymbol{A}^T \boldsymbol{C}^{-1} \boldsymbol{C} \boldsymbol{K}_o^{\mathsf{T}} \\
&= (\boldsymbol{A}^{\mathsf{T}} \boldsymbol{C}^{-1} \boldsymbol{A})^{-1}
\end{aligned}$$

where we used $\boldsymbol{K}_o \boldsymbol{A} = \boldsymbol{I}$

- Furthermore, note that

$$\begin{aligned}
E(\tilde{\boldsymbol{\theta}}_o \tilde{\boldsymbol{\theta}}^{\mathsf{T}}) &= E(\boldsymbol{K}_o \mathbf{n} \mathbf{n}^{\mathsf{T}} \boldsymbol{K}^{\mathsf{T}}) \\
&= \boldsymbol{K}_o \boldsymbol{C} \boldsymbol{K}^{\mathsf{T}} \\
&= (\boldsymbol{A}^{\mathsf{T}} \boldsymbol{C}^{-1} \boldsymbol{A})^{-1} \boldsymbol{A}^{\mathsf{T}} \boldsymbol{C}^{-1} \boldsymbol{C} \boldsymbol{K}^{\mathsf{T}} \\
&= (\boldsymbol{A}^{\mathsf{T}} \boldsymbol{C}^{-1} \boldsymbol{A})^{-1}
\end{aligned}$$

where we used $\boldsymbol{K} \boldsymbol{A} = \boldsymbol{I}$

## Proof of the Gauss-Markov Theorem

- We can now develop the following positive-semidefinite expression

$$
\begin{aligned}
E((\tilde{\boldsymbol{\theta}}_o - \tilde{\boldsymbol{\theta}})(\tilde{\boldsymbol{\theta}}_o - \tilde{\boldsymbol{\theta}})^{\mathsf{T}}) &= E(\tilde{\boldsymbol{\theta}}_o \tilde{\boldsymbol{\theta}}_o^{\mathsf{T}} - \tilde{\boldsymbol{\theta}} \tilde{\boldsymbol{\theta}}_o^{\mathsf{T}} - \tilde{\boldsymbol{\theta}}_o \tilde{\boldsymbol{\theta}}^{\mathsf{T}} + \tilde{\boldsymbol{\theta}} \tilde{\boldsymbol{\theta}}^{\mathsf{T}}) \\
&= \underbrace{E(\tilde{\boldsymbol{\theta}}_o \tilde{\boldsymbol{\theta}}_o^{\mathsf{T}})}_{\mathsf{MSE}_{\theta}(\hat{\boldsymbol{\theta}}_o)} - E(\tilde{\boldsymbol{\theta}} \tilde{\boldsymbol{\theta}}_o^{\mathsf{T}}) - E(\tilde{\boldsymbol{\theta}}_o \tilde{\boldsymbol{\theta}}^{\mathsf{T}}) + \underbrace{E(\tilde{\boldsymbol{\theta}} \tilde{\boldsymbol{\theta}}^{\mathsf{T}})}_{\mathsf{MSE}_{\theta}(\hat{\boldsymbol{\theta}})} \\
&= (\boldsymbol{A}^{\mathsf{T}} \boldsymbol{C}^{-1} \boldsymbol{A})^{-1} - 2(\boldsymbol{A}^{\mathsf{T}} \boldsymbol{C}^{-1} \boldsymbol{A})^{-1} + \boldsymbol{K} \boldsymbol{C} \boldsymbol{K}^{\mathsf{T}} \\
&= -(\boldsymbol{A}^{\mathsf{T}} \boldsymbol{C}^{-1} \boldsymbol{A})^{-1} + \boldsymbol{K} \boldsymbol{C} \boldsymbol{K}^{\mathsf{T}} \\
&= -\mathsf{MSE}_{\theta}(\hat{\boldsymbol{\theta}}_o) + \mathsf{MSE}_{\theta}(\hat{\boldsymbol{\theta}}) \\
&\succcurlyeq \boldsymbol{0}
\end{aligned}
$$

- Finally, we get $\mathsf{MSE}_{\theta}(\hat{\boldsymbol{\theta}}_o) \preccurlyeq \mathsf{MSE}_{\theta}(\hat{\boldsymbol{\theta}})$ with equality iff $\hat{\boldsymbol{\theta}}_o = \hat{\boldsymbol{\theta}}$.

# Asymptotic Properties of ML Estimation

- Recall that $\mathbf{y} \in \mathbb{R}^m$; for a finite number of data records $m$, in general, the ML estimator is not the UMVUE

- However, for $p(\mathbf{y}, \theta)$ being from a RSF and $m \to \infty$, it can be shown that the ML estimate is *asymptotically unbiased, asymptotically Gaussian, and asymptotically efficient/UMVUE*, i.e.,

$$\hat{\boldsymbol{\theta}}_{\mathrm{ML}}(\boldsymbol{y}) \sim \mathcal{N}(\boldsymbol{\theta}, \boldsymbol{J}_{\boldsymbol{\theta}}^{-1}) \quad \text{for } m \to \infty$$

where $\boldsymbol{J}_{\boldsymbol{\theta}}$ is the Fisher information matrix

- Since (apart from pathological situations) $\boldsymbol{J}_{\boldsymbol{\theta}}^{-1} \to 0$ for $m \to \infty$, it follows that $\hat{\boldsymbol{\theta}}_{\mathrm{ML}}(\boldsymbol{y})$ is consistent, i.e.,

$$\hat{\boldsymbol{\theta}}_{\mathrm{ML}} \to \boldsymbol{\theta} \quad \text{for } m \to \infty$$

- Proof: See Kay, Appendix 7B

# Parameter Transformation

- Let $\boldsymbol{\alpha} = \boldsymbol{g}(\boldsymbol{\theta})$, where $\boldsymbol{\alpha}$ and $\boldsymbol{\theta}$ may have different dimensions

- We recall that $\hat{\boldsymbol{\theta}}_{\mathrm{ML}}(\boldsymbol{y}) \triangleq \arg\max_{\boldsymbol{\theta}} p(\boldsymbol{y}; \boldsymbol{\theta})$

- The ML estimate of $\boldsymbol{\alpha}$ can be defined as

$$\hat{\boldsymbol{\alpha}}_{\mathrm{ML}}(\boldsymbol{y}) \triangleq \arg\max_{\boldsymbol{\alpha}} \tilde{p}(\boldsymbol{y}; \boldsymbol{\alpha})$$

where $\tilde{p}(\boldsymbol{y}; \boldsymbol{\alpha})$ is given as follows

1. if $\boldsymbol{g}(\boldsymbol{\theta})$ is invertible, then $\tilde{p}(\boldsymbol{y}; \boldsymbol{\alpha}) = p(\boldsymbol{y}; \boldsymbol{g}^{-1}(\boldsymbol{\alpha}))$
2. if $\boldsymbol{g}(\boldsymbol{\theta})$ is not invertible, i.e., to a given $\boldsymbol{\alpha}$ belonging to the range of $\boldsymbol{g}(\cdot)$, there exist several $\boldsymbol{\theta}_i$ such that $\boldsymbol{\alpha} = \boldsymbol{g}(\boldsymbol{\theta}_i)$, then $\tilde{p}(\boldsymbol{y}; \boldsymbol{\alpha}) = \max_{i:\boldsymbol{g}(\boldsymbol{\theta}_i)=\boldsymbol{\alpha}} p(\boldsymbol{y}; \boldsymbol{\theta}_i)$

- In either case, we have $\hat{\boldsymbol{\alpha}}_{\mathrm{ML}}(\boldsymbol{y}) = \boldsymbol{g}(\hat{\boldsymbol{\theta}}_{\mathrm{ML}}(\boldsymbol{y}))$