

UNIT 1

Estimation Theory



Johannes Brandstetter

Institute for Machine Learning

Copyright statement:

This material, no matter whether in printed or electronic form, may be used for personal and non-commercial educational use only. Any reproduction of this material, no matter whether as a whole or in parts, no matter whether in printed or in electronic form, requires explicit prior acceptance of the authors.

Basics of Estimation Theory

- Given: Data distribution $p(x, w)$ parameterized by $w \in \mathbb{R}^d$:
 - Remember $Z = \begin{pmatrix} X \\ y^T \end{pmatrix}$
 - For unsupervised learning: $Z = X$
- Goal: Estimate w by \hat{w} , given training data $X = (x_1, \dots, x_l)$ so that distribution of X is reflected accurately.
- Unbiased estimator: $E_X(\hat{w}) = w$ for all possible w
 - $E_X(\cdot)$ is the expectation wrt. data distribution $p(x, w)$
 - Bias: $b(\hat{w}) = E_X(\hat{w}) - w$
 - Variance: $\text{var}(\hat{w}) = E_X((\hat{w} - E_X(\hat{w}))^T(\hat{w} - E_X(\hat{w})))$
- For the rest of presentation: w : vector valued, w : scalar valued ($d = 1$).

■ We define the MSE as expectation over the training set:

$$\begin{aligned}\text{mse}(\hat{\mathbf{w}}) &= \mathbb{E}_{\mathbf{X}} \left((\hat{\mathbf{w}} - \mathbf{w})^T (\hat{\mathbf{w}} - \mathbf{w}) \right) \\ &= \mathbb{E}_{\mathbf{X}} \left(((\hat{\mathbf{w}} - \mathbb{E}_{\mathbf{X}}(\hat{\mathbf{w}})) + (\mathbb{E}_{\mathbf{X}}(\hat{\mathbf{w}}) - \mathbf{w}))^T \right. \\ &\quad \left. ((\hat{\mathbf{w}} - \mathbb{E}_{\mathbf{X}}(\hat{\mathbf{w}})) + (\mathbb{E}_{\mathbf{X}}(\hat{\mathbf{w}}) - \mathbf{w})) \right) \\ &= \mathbb{E}_{\mathbf{X}} \left((\hat{\mathbf{w}} - \mathbb{E}_{\mathbf{X}}(\hat{\mathbf{w}}))^T (\hat{\mathbf{w}} - \mathbb{E}_{\mathbf{X}}(\hat{\mathbf{w}})) \right. \\ &\quad \left. + 2 \underbrace{(\hat{\mathbf{w}} - \mathbb{E}_{\mathbf{X}}(\hat{\mathbf{w}}))^T}_{=0} (\mathbb{E}_{\mathbf{X}}(\hat{\mathbf{w}}) - \mathbf{w}) \right. \\ &\quad \left. + (\mathbb{E}_{\mathbf{X}}(\hat{\mathbf{w}}) - \mathbf{w})^T (\mathbb{E}_{\mathbf{X}}(\hat{\mathbf{w}}) - \mathbf{w}) \right) \\ &= \underbrace{\mathbb{E}_{\mathbf{X}} \left((\hat{\mathbf{w}} - \mathbb{E}_{\mathbf{X}}(\hat{\mathbf{w}}))^T (\hat{\mathbf{w}} - \mathbb{E}_{\mathbf{X}}(\hat{\mathbf{w}})) \right)}_{\text{var}(\hat{\mathbf{w}})} \\ &\quad + \underbrace{(\mathbb{E}_{\mathbf{X}}(\hat{\mathbf{w}}) - \mathbf{w})^T (\mathbb{E}_{\mathbf{X}}(\hat{\mathbf{w}}) - \mathbf{w})}_{b^2(\hat{\mathbf{w}})}\end{aligned}$$

Example MSE of estimator

- Example: $x_i = \varepsilon_i$, $\varepsilon_i \sim \mathcal{N}(\mu, \sigma^2)$ i.i.d. (x_i is scalar-valued for simplicity)
- We consider the estimator

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

- Bias $b^2(\hat{\mu})$ and variance $\text{var}(\hat{\mu})$ can be calculated as:

$$b^2(\hat{\mu}) = (\mathbf{E}_{\mathbf{X}}(\hat{\mu}) - \mu)^2 = \mu^2(\lambda - 1)^2$$

$$\text{var}(\hat{\mu}) = \mathbf{E}_{\mathbf{X}}((\hat{\mu} - \mathbf{E}_{\mathbf{X}}(\hat{\mu}))^T(\hat{\mu} - \mathbf{E}_{\mathbf{X}}(\hat{\mu})))$$

$$= \mathbf{E}_{\mathbf{X}} \left(\left(\frac{\lambda}{l} \sum_{i=1}^l x_i - \frac{\lambda}{l} \sum_{i=1}^l \mathbf{E}_{\mathbf{X}}(x_i) \right)^T \left(\frac{\lambda}{l} \sum_{i=1}^l x_i - \frac{\lambda}{l} \sum_{i=1}^l \mathbf{E}_{\mathbf{X}}(x_i) \right) \right)$$

$$= \mathbf{E}_{\mathbf{X}} \left(\frac{\lambda^2}{l^2} \left(\sum_{i=1}^l x_i - \sum_{i=1}^l \mathbf{E}_{\mathbf{X}}(x_i) \right)^T \left(\sum_{i=1}^l x_i - \sum_{i=1}^l \mathbf{E}_{\mathbf{X}}(x_i) \right) \right)$$

$$= \frac{\lambda^2}{l^2} \cdot l \sigma^2 = \frac{\lambda^2 \sigma^2}{l}$$

- Often: multiple estimates of the same parameters, i.e., $\{\hat{w}_1, \dots, \hat{w}_n\}$. Then they are combined to: $\hat{w} = \frac{1}{n} \sum_{j=1}^n \hat{w}_j$
- For unbiased and uncorrelated \hat{w}_j , that means:

$$E_{\mathbf{X}}(\hat{w}) = w$$

$$\text{Var}(\hat{w}) = \frac{1}{n^2} \sum_{j=1}^n \text{Var}(\hat{w}_j) = \frac{\text{Var}(\hat{w}_j)}{n}$$

- The more estimates are averaged, the more the variance decreases
- But if $E_{\mathbf{X}}(\hat{w}_j) = w + b^2(\hat{w}_j)$, which implies $E_{\mathbf{X}}(\hat{w}) = w + b^2(\hat{w}_j)$, no matter how many estimates are averaged, \hat{w} does not converge to the true w

Example unbiased estimator

- Example: $x_i = \varepsilon_i$, $\varepsilon_i \sim \mathcal{N}(\mu, \sigma^2)$ i.i.d. (x_i is scalar-valued for simplicity)
- Then, $w = \begin{pmatrix} \mu \\ \sigma \end{pmatrix}$
- We estimate $\hat{\mu}$ as:

$$\hat{\mu} = \frac{1}{l} \sum_{i=1}^l x_i, \quad E(\hat{\mu}) = \frac{1}{l} \sum_{i=1}^l E(x_i) = \mu$$

- For any value of μ , $E_X(\hat{\mu}) = \mu$
 $\Rightarrow \hat{\mu} = \frac{1}{l} \sum_{i=1}^l x_i$ is **(asymptotic) unbiased estimator** of μ

Example biased estimator

- Example: $x_i = \varepsilon_i$, $\varepsilon_i \sim \mathcal{N}(\mu, \sigma^2)$ i.i.d. (x_i is scalar-valued for simplicity)
- Then, $w = \begin{pmatrix} \mu \\ \sigma \end{pmatrix}$
- We estimate $\hat{\mu}$ as:

$$\hat{\mu} = \frac{1}{2l} \sum_{i=1}^l x_i, \quad E(\hat{\mu}) = \frac{1}{2l} \sum_{i=1}^l E(x_i) = \frac{1}{2}\mu$$

- For any other value than $\mu = 0$, $E_X(\hat{\mu}) \neq \mu$
 $\Rightarrow \hat{\mu} = \frac{1}{2l} \sum_{i=1}^l x_i$ is **biased estimator** of μ

Example biased estimator

- Example: $x_i = \varepsilon_i$, $\varepsilon_i \sim \mathcal{N}(\mu, \sigma^2)$ i.i.d. (x_i is scalar-valued)
- Then, $\hat{\sigma}^2 = \frac{1}{l} \sum_{i=1}^l (x_i - \hat{\mu})^2$ with $\hat{\mu} = \frac{1}{l} \sum_{i=1}^l x_i$ is a **biased estimator**:

$$\begin{aligned}\text{E}_{\mathbf{X}}(\hat{\sigma}^2) &= \text{E}_{\mathbf{X}} \left(\frac{1}{l} \sum_{i=1}^l (x_i - \hat{\mu})^2 \right) \\ &= \frac{1}{l} \text{E}_{\mathbf{X}} \left(\sum_{i=1}^l x_i^2 - \sum_{i=1}^l 2x_i \hat{\mu} + \sum_{i=1}^l \hat{\mu}^2 \right) \\ &= \frac{1}{l} \text{E}_{\mathbf{X}} \left(\sum_{i=1}^l x_i^2 - 2\hat{\mu} \sum_{i=1}^l x_i + l\hat{\mu}^2 \right) \\ &= \frac{1}{l} \text{E}_{\mathbf{X}} \left(\sum_{i=1}^l x_i^2 - l\hat{\mu}^2 \right) \\ &= \frac{1}{l} \sum_{i=1}^l \text{E}_{\mathbf{X}}(x_i^2) - \frac{1}{l} \sum_{i=1}^l l\text{E}_{\mathbf{X}}(\hat{\mu}^2)\end{aligned}$$

Example biased estimator

- Using the following properties:

$$\text{var}(x_i) = E_{\mathbf{X}}(x_i^2) - [E_{\mathbf{X}}(x_i)]^2 \Rightarrow E_{\mathbf{X}}(x_i^2) = \sigma^2 + \mu^2$$

$$\text{var}(\hat{\mu}) = E_{\mathbf{X}}(\hat{\mu}^2) - [E_{\mathbf{X}}(\hat{\mu})]^2 \Rightarrow E_{\mathbf{X}}(\hat{\mu}^2) = \frac{\sigma^2}{l} + \mu^2$$

- We get:

$$\begin{aligned} E_{\mathbf{X}}(\hat{\sigma}^2) &= \frac{1}{l} \sum_{i=1}^l E_{\mathbf{X}}(x_i^2) - \frac{1}{l} \sum_{i=1}^l l E_{\mathbf{X}}(\hat{\mu}^2) \\ &= \frac{1}{l} (l\sigma^2 + l\mu^2 - \sigma^2 - l\mu^2) \\ &= \frac{(l-1)}{l} \sigma^2 \\ &\neq \sigma^2 \end{aligned}$$

- The unbiased estimator of the sample variance is

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

Quality criterion, MVU estimator

- Mean squared error as quality criterion:

$$\text{mse}(\hat{\boldsymbol{w}}) = \text{var}(\hat{\boldsymbol{w}}) + b^2(\hat{\boldsymbol{w}})$$

- Goal: find estimator with minimal mse
- However, if bias term – and thus optimal value – depend on \boldsymbol{w} , this goal is not realizable
- \Rightarrow Constrain bias term to be zero and find estimators that minimize variance \Rightarrow **MVU (minimal variance unbiased)** estimator
- Does MVU estimator exist? In general: no

How to find MVU?

- Even if MVU estimator exists, there is no known procedure which always produces it, but we can do the following:
- We place a lower bound on the variance of all unbiased estimators. If one estimator attains this bound for all values of the parameter \Rightarrow MVU estimator
- This bound is called **Cramér-Rao lower bound (CRLB)**

Maximum Likelihood Estimator (MLE)

- MVU estimator often unknown or does not exist \Rightarrow How to construct estimator that approximates MVU estimator?
- This estimator should be **asymptotically efficient** and **asymptotically unbiased**
- **Likelihood** of i.i.d. data \mathbf{X} :

$$\mathcal{L}(\mathbf{X}; \mathbf{w}) = \prod_{i=1}^l p(x_i; \mathbf{w})$$

- I.e., probability of the (data-distribution) model $p(\mathbf{x}; \mathbf{w})$ to (re-)produce the dataset
- We want to maximize the likelihood, or equivalently minimize the **negative log-likelihood**:

$$L = -\ln(\mathcal{L}) = -\sum_{i=1}^l \ln(p(x_i; \mathbf{w}))$$

MLE of normal distribution

- Example: $x_i = \varepsilon_i$, $\varepsilon_i \sim \mathcal{N}(\mu, \sigma^2)$ i.i.d. (x_i is scalar-valued for simplicity)
- Then, $\mathbf{w} = \begin{pmatrix} \mu \\ \sigma \end{pmatrix}$:

$$p(x_i; \mathbf{w}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2}\frac{(x_i - \mu)^2}{\sigma^2}\right)$$

- The mean μ and the variance σ^2 are the two parameters that need to be estimated.
- The likelihood function is:

$$\mathcal{L}(\mu, \sigma^2, x_1, \dots, x_l) = (2\pi\sigma^2)^{-\frac{l}{2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^l (x_i - \mu)^2\right)$$

MLE of normal distribution

- The negative log-likelihood function is:

$$L = -\ln(\mathcal{L})(\mu, \sigma^2, x_1, \dots, x_l) = \frac{l}{2} \ln(2\pi) + \frac{l}{2} \ln(\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^l (x_i - \mu)^2$$

- Maximum likelihood estimate of μ :

$$\begin{aligned}\frac{\partial L}{\partial \mu} &= \frac{\partial}{\partial \mu} \left(\frac{l}{2} \ln(2\pi) + \frac{l}{2} \ln(\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^l (x_i - \mu)^2 \right) \\ &= -\frac{1}{\sigma^2} \sum_{i=1}^l (x_i - \mu) \\ \Rightarrow \sum_{i=1}^l x_i - l\mu &= 0 \Rightarrow \hat{\mu} = \frac{1}{l} \sum_{i=1}^l x_i\end{aligned}$$

- Estimator $\hat{\mu} = \frac{1}{l} \sum_{i=1}^l x_i$ is **unbiased**

MLE of normal distribution

- The negative log-likelihood function is:

$$L = -\ln(\mathcal{L})(\mu, \sigma^2, x_1, \dots, x_l) = \frac{l}{2} \ln(2\pi) + \frac{l}{2} \ln(\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^l (x_i - \mu)^2$$

- Maximum likelihood estimate of σ^2 :

$$\begin{aligned}\frac{\partial L}{\partial \sigma^2} &= \frac{\partial}{\partial \sigma^2} \left(\frac{l}{2} \ln(2\pi) + \frac{l}{2} \ln(\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^l (x_i - \mu)^2 \right) \\ &= \frac{l}{2\sigma^2} - \frac{1}{2(\sigma^2)^2} \sum_{i=1}^l (x_i - \mu)^2 = \frac{1}{2\sigma^2} \left[l - \left(\frac{1}{\sigma^2} \sum_{i=1}^l (x_i - \mu)^2 \right) \right] \\ \Rightarrow l - \left(\frac{1}{\sigma^2} \sum_{i=1}^l (x_i - \mu)^2 \right) &= 0 \Rightarrow \hat{\sigma}^2 = \frac{1}{l} \sum_{i=1}^l (x_i - \mu)^2\end{aligned}$$

- Estimator $\hat{\sigma}^2 = \frac{1}{l} \sum_{i=1}^l (x_i - \mu)^2$ is **asymptotically unbiased**

Fisher information matrix

- The **Fisher information matrix** $\mathbf{I}_F(\mathbf{w})$ for a parametrized model:

$$\mathbf{I}_F(\mathbf{w}) : [\mathbf{I}_F(\mathbf{w})]_{ij} = \mathbb{E}_{p(\mathbf{x}; \mathbf{w})} \left(\frac{\partial \ln p(\mathbf{x}; \mathbf{w})}{\partial w_i} \frac{\partial \ln p(\mathbf{x}; \mathbf{w})}{\partial w_j} \right),$$

where $[\mathbf{A}]_{ij}$ selects the ij th component of a matrix \mathbf{A} .

- If the density function $p(\mathbf{x}; \mathbf{w})$ satisfies:

$$\forall \mathbf{w} : \mathbb{E}_{p(\mathbf{x}; \mathbf{w})} \underbrace{\left(\frac{\partial \ln p(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w}} \right)}_{\text{score}} = \mathbf{0},$$

then

$$\mathbf{I}_F(\mathbf{w}) = -\mathbb{E}_{p(\mathbf{x}; \mathbf{w})} \left(\frac{\partial^2 \ln p(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}} \right).$$

Fisher information matrix

$$\begin{aligned}\mathbf{0} &= \mathbb{E}_{p(\mathbf{x}; \mathbf{w})} \left(\frac{\partial \ln p(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w}} \right) \\&= \int_{\mathbf{X}} \frac{\partial \ln p(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w}} p(\mathbf{x}; \mathbf{w}) d\mathbf{x} \quad \text{assuming } \mathbf{X} \text{ is distributed as } p(\mathbf{x}; \mathbf{w}) \\&\Rightarrow \frac{\partial}{\partial \mathbf{w}} \int_{\mathbf{X}} \frac{\partial \ln p(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w}} p(\mathbf{x}; \mathbf{w}) d\mathbf{x} = \mathbf{0} \\&\Rightarrow \int_{\mathbf{X}} \left(\frac{\partial^2 \ln p(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}} p(\mathbf{x}; \mathbf{w}) + \frac{\partial \ln p(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w}} \frac{\partial p(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w}} \right) d\mathbf{x} \\&\Rightarrow \int_{\mathbf{X}} \left(\frac{\partial^2 \ln p(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}} p(\mathbf{x}; \mathbf{w}) + \frac{\partial \ln p(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w}} \frac{\partial \ln p(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w}} p(\mathbf{x}; \mathbf{w}) \right) d\mathbf{x} \\&\Rightarrow -\mathbb{E}_{p(\mathbf{x}; \mathbf{w})} \left(\frac{\partial^2 \ln p(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w} \partial \mathbf{w}} \right) = \mathbb{E}_{p(\mathbf{x}; \mathbf{w})} \left(\frac{\partial \ln p(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w}} \frac{\partial \ln p(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w}} \right)\end{aligned}$$

Fisher information matrix

- Fisher information is the average curvature of the negative log-likelihood
- Main rationale: the larger this quantity, the smaller the variance of $p(x, w)$
- The Fisher information gives the amount of information that an observable random variable x carries about an unobservable parameter w upon which the parameterized density function $p(x; w)$ of x depends

Cramér-Rao lower bound (CRLB)

- Let \hat{w} be an unbiased estimator of w
- Let $Cov(\hat{w})$ be the covariance matrix and

$$[\mathbf{I}_F(w)]_{ij} = -\mathbb{E}_{p(x;w)} \left(\frac{\partial^2 \ln p(x;w)}{\partial w_i \partial w_j} \right)$$

- CRLB:**

$$Cov(\hat{w}) \geq \frac{1}{\mathbf{I}_F(w)}$$

- An estimator is said to be **efficient** if it reaches the CRLB. It is efficient in that it efficiently makes use of the data and extracts information to estimate the parameter
- A MVU estimator may or may not be efficient. This means it could have minimum variance but without reaching the CRLB

CRLB of normal distribution

- First we calculate $\text{var}(\hat{\mu})$ and $\text{var}(\hat{\sigma^2})$:

$$\begin{aligned}\text{var}(\hat{\mu}) &= \mathbb{E}_{\mathbf{X}} \left((\hat{\mu} - \mathbb{E}_{\mathbf{X}}(\hat{\mu}))^T (\hat{\mu} - \mathbb{E}_{\mathbf{X}}(\hat{\mu})) \right) \\ &= \mathbb{E}_{\mathbf{X}} \left(\left(\frac{1}{l} \sum_{i=1}^l x_i - \frac{1}{l} \sum_{i=1}^l \mathbb{E}_{\mathbf{X}}(x_i) \right)^T \left(\frac{1}{l} \sum_{i=1}^l x_i - \frac{1}{l} \sum_{i=1}^l \mathbb{E}_{\mathbf{X}}(x_i) \right) \right) \\ &= \mathbb{E}_{\mathbf{X}} \left(\frac{1}{l^2} \left(\sum_{i=1}^l x_i - \sum_{i=1}^l \mathbb{E}_{\mathbf{X}}(x_i) \right)^T \left(\sum_{i=1}^l x_i - \sum_{i=1}^l \mathbb{E}_{\mathbf{X}}(x_i) \right) \right) \\ &= \frac{1}{l^2} \cdot l \sigma^2 = \frac{\sigma^2}{l} \\ \text{var}(\hat{\sigma^2}) &= \frac{2\sigma^4}{l-1}\end{aligned}$$

- Derivation of $\text{var}(\hat{\sigma^2})$ via $\frac{(l-1)\hat{\sigma^2}}{\sigma^2} \sim \chi_{l-1}^2$

CRLB of normal distribution

- The negative log-likelihood function is:

$$L = -\ln(\mathcal{L})(\mu, \sigma^2, x_1, \dots, x_l) = \frac{l}{2} \ln(2\pi) + \frac{l}{2} \ln(\sigma^2) + \frac{1}{2\sigma^2} \sum_{i=1}^l (x_i - \mu)^2$$

$$[\mathbf{I}_F(\mathbf{w})]_{ij} = -\mathbb{E}_{p(\mathbf{x}; \mathbf{w})} \left(\frac{\partial^2 \ln p(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w}_i \partial \mathbf{w}_j} \right):$$

$$\Rightarrow \nabla_{\mu, \sigma^2}^2 L = \begin{pmatrix} \frac{\partial^2}{\partial \mu^2} & \frac{\partial^2}{\partial \mu \partial \sigma^2} \\ \frac{\partial^2}{\partial \sigma^2 \partial \mu} & \frac{\partial^2}{\partial (\sigma^2)^2} \end{pmatrix}$$

$$= \begin{pmatrix} \frac{l}{\sigma^2} & \frac{1}{\sigma^4} \sum_{i=1}^l (x_i - \mu) \\ \frac{1}{\sigma^4} \sum_{i=1}^l (x_i - \mu) & \frac{1}{(\sigma^2)^3} \sum_{i=1}^l (x_i - \mu)^2 - \frac{l}{2(\sigma^2)^2} \end{pmatrix}$$

$$\Rightarrow -\mathbb{E}_{p(\mathbf{x}; \mathbf{w})} \left(\frac{\partial^2 \ln p(\mathbf{x}; \mathbf{w})}{\partial \mathbf{w}_i \partial \mathbf{w}_j} \right) = \mathbb{E} (\nabla_{\mu, \sigma^2}^2 L) = \begin{pmatrix} \frac{l}{\sigma^2} & 0 \\ 0 & \frac{l}{2(\sigma^2)^2} \end{pmatrix}$$

■ CRLB:

$$\begin{aligned} Cov(\hat{\boldsymbol{w}}) &\geq \frac{1}{\mathbf{I}_F(\boldsymbol{w})} \\ \begin{pmatrix} \frac{\sigma^2}{l} & 0 \\ 0 & \frac{2\sigma^4}{l-1} \end{pmatrix} &\geq \begin{pmatrix} \frac{l}{\sigma^2} & 0 \\ 0 & \frac{l}{2(\sigma^2)^2} \end{pmatrix}^{-1} \\ &\geq \begin{pmatrix} \frac{\sigma^2}{l} & 0 \\ 0 & \frac{2(\sigma^2)^2}{l} \end{pmatrix} \end{aligned}$$

- MLE estimate $\hat{\mu} = \frac{1}{l} \sum_{i=1}^l x_i$ is **efficient**
- MLE estimate $\hat{\sigma^2} = \frac{1}{l} \sum_{i=1}^l (x_i - \mu)^2$ is **asymptotically efficient**

Properties of MLE

- Summary: MLE is a popular estimator on which almost all practical estimation tasks are based.
- The popularity stems from the fact that it can be applied to a broad range of problems and it approximates the MVU estimator for large data sets.
- MLE is asymptotically efficient and unbiased.

Properties of MLE

- **Invariance under parameter change:** Let g be a function changing the parameter \mathbf{w} into parameter $\mathbf{u} : \mathbf{u} = g(\mathbf{w})$, then

$$\hat{\mathbf{u}} = g(\hat{\mathbf{w}}) ,$$

where \mathbf{u} and \mathbf{w} are ML estimators. If g changes \mathbf{w} into different \mathbf{w} then

$\hat{\mathbf{u}} = g(\hat{\mathbf{w}})$ maximizes the likelihood function: $\max_{\mathbf{w}: \mathbf{u}=g(\mathbf{w})} p(\{\mathbf{x}\}; \mathbf{w})$

- **Asymptotically unbiased:**

$$\mathrm{E}_{p(\mathbf{x}; \mathbf{w})}(\hat{\mathbf{w}}) \xrightarrow{l \rightarrow \infty} \mathbf{w}$$

- **Asymptotically efficient:**

$$\mathrm{Cov}(\hat{\mathbf{w}}) \xrightarrow{l \rightarrow \infty} \text{CRLB}$$