

Midterm2

Sufficient Statistic

$p(x_1, \dots, x_n | T(x_1, \dots, x_n; \theta))$ does not depend on θ for any T

The factorization theorem (Neyman-Fisher)

$$p(x_1, \dots, x_n; \theta) = h(x_1, \dots, x_n) g(T(x_1, \dots, x_n; \theta))$$

Example: $x_1, \dots, x_n \sim N(\mu, \sigma^2)$

$$\begin{aligned} p(x_1, \dots, x_n; \mu, \sigma^2) &= \prod \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right) \\ &= \underbrace{\frac{1}{\sqrt{2\pi\sigma^2}^n}}_{h} \underbrace{\exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2\right)}_{g} \end{aligned}$$

$$T(x_1, \dots, x_n) = (\sum_{i=1}^n x_i, \sum_{i=1}^n x_i) \quad \text{if } \mu, \sigma^2 \text{ unknown}$$

$$T(x_1, \dots, x_n) = (\sum x_i). \quad \text{if } \mu \text{ unknown, } \sigma^2 \text{ known}$$

Proof of factorization theorem

$$L(\theta) = \underbrace{g(T(x_1, \dots, x_n; \theta))}_{\text{ignore constants do not depend on } \theta} h(x_1, \dots, x_n)$$

$$L(\theta) = g(T(x_1, \dots, x_n; \theta))$$

Minimal sufficiency,

$$\text{Define ratio: } R = \frac{p(y_1, \dots, y_n; \theta)}{p(x_1, \dots, x_n; \theta)}$$

T is MSS if R does not depend on θ iff $T(x_1, \dots, x_n) = T(y_1, \dots, y_n)$

Rao-Blackwell theorem:

$$R(\hat{\theta}, \theta) \leq R(\vec{\theta} | \theta). \quad \hat{\theta} = E(\vec{\theta} | T)$$

Example: $x_1, \dots, x_n \sim \text{Ber}(\theta)$.

$$\hat{\theta} = x_1$$

$$\hat{\theta} = E[x_1 \mid \sum^n x_i] = 1 \cdot p(x_1=1 \mid \sum^n x_i) = \frac{p \cdot C_{n-1}^{T-1} p^{T-1} (1-p)^{n-T}}{C_n^T p^T (1-p)^{n-T}}$$

$$= \frac{T}{n}$$

$$R(\hat{\theta}) = \theta(-\theta)$$

$$R(\tilde{\theta}) = \frac{\theta(-\theta)}{n} < \theta(-\theta)$$

Proof of Rao-Blackwell.

$$\begin{aligned} R(\tilde{\theta}, \theta) &= E[(E[\tilde{\theta}|T] - \theta)^2] \\ &= E[(E[\hat{\theta} - \theta|T])^2] \\ &\leq E[E[(\hat{\theta} - \theta)^2|T]] \quad \text{by Jensen.} \\ &= R(\hat{\theta}, \theta) \end{aligned}$$

Exponential Family

$$P(x; \theta) = \exp \left[\sum_i y_i(\theta) T_i(x) - A(\theta) \right] h(x)$$

$$A: \Theta \rightarrow \mathbb{R}$$

canonical parametrization:

$$P(x; \theta) = \exp \left[\sum_i \theta_i T_i(x) - A(\theta) \right] h(x)$$

θ : natural parameters.

properties of Exponential Families

- Random sampling.

The exponential structure is preserved for an iid sample $\{x_1, \dots, x_n\}_{n \in \mathbb{N}}$

$$P(x_1, \dots, x_n; \theta) = \prod h(x_i) \exp \left[\sum_i \theta_i \sum_j T_j(x_i) - n A(\theta) \right]$$

same natural parameters
 $T_1, \dots, \underset{n}{\overbrace{T \dots}}$

$$l(x_1, \dots, x_n) = \sum l_i(x_j)$$

- $A(\theta)$ log-normalization constant

log-partition function
cumulant function

$$A(\theta) = \log \int_x \exp \left[\sum \theta_i T_i(x) \right] h(x) dx$$

$$\frac{A(\theta)}{\theta_i} = \frac{\int_x \exp \left[\sum \theta_i T_i(x) \right] h(x) dx}{\int_x \exp \left[\sum \theta_i T_i(x) \right] h(x) dx} = E[T_i(x)]$$

$$\frac{\partial A(\theta)}{\partial \theta_i \partial \theta_j} = \text{cov}(T_i(x), T_j(x))$$

- The log likelihood in an exponential family is concave

$$LL(\theta; x_1, \dots, x_n) \propto \left[\sum \theta_i \sum T_i(x_j) - nA(\theta) \right]$$

Hessian is $(-n) \times$ Hessian of A . A is convex

- minimal

no set of coefficient. such that $\sum a_i T_i(x) = \text{const}$

Over complete exponential families are not statistically identifiable.

- The exponential families arise naturally
 - 1. maximize the entropy of distribution
 - 2. constraint $\hat{m}_i = E_p[T_i(x)]$

- MLE coincide with MOM

$$LL(\theta; x_1, \dots, x_n) \propto \left[\sum \theta_i \sum T_i(x_j) - nA(\theta) \right] \Rightarrow \text{concave}.$$

$$\partial LL(\theta; x_1, \dots, x_n) = \sum^n T_i(x_j) - nA'(\theta)$$

$$= \underbrace{\sum_j T_i(x_j)}_{j} - n \underbrace{E_p[T_i(x)]}_{\bar{T}_i(x)}$$

Point Estimation

The Invariant of moments

Maximum Likelihood

The MLE is equivariant

Bayes Estimator

$$P(\theta | x_1, \dots, x_n) \propto L(\theta) P(\theta) \quad \text{Likelihood} \times \text{prior}$$

Example 1: $x_1, \dots, x_n \sim \text{Ber}(\theta)$. $\theta \sim \text{Beta}(2, \beta)$ prior

$$P(\theta) \propto \theta^{2-1} (1-\theta)^{\beta-1}$$

$$P(\theta | x) \propto L(\theta) P(\theta) = \theta^{\sum x_i} (1-\theta)^{n-\sum x_i} \theta^{2-1} (1-\theta)^{\beta-1}$$

$$P(\theta | x) = \text{Beta}(\alpha + \sum x_i, \beta + n - \sum x_i)$$

$$\hat{\theta}_{\text{Bayes}} = \frac{\alpha + \sum x_i}{n + \alpha + \beta} = (1 - \lambda) \hat{\theta}_{\text{MLE}} + \lambda \bar{x}$$

Example 2: $x_1, \dots, x_n \sim N(\mu, \sigma^2)$, $\mu \sim N(m, \tau^2)$

$$P(\mu | x) \propto \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^n \exp\left(-\frac{\sum (x_i - \mu)^2}{2\sigma^2}\right) \frac{1}{\sqrt{2\pi\tau^2}} \exp\left(-\frac{(m - \mu)^2}{2\tau^2}\right)$$

$$\frac{1}{\sqrt{2\pi\tau^2}} \exp\left(-\frac{\sum (x_i - \mu)^2 + (m - \mu)^2}{2\tau^2}\right)$$

$$\frac{1}{\sqrt{2\pi\tau^2}} \left(n\mu^2 - 2\sum x_i \mu + \sum x_i^2 \right) +$$

$$\sigma^2 (m^2 - 2m\mu + \mu^2)$$

$$(n\tau^2 + \sigma^2)\mu^2 - (2\sum x_i \tau^2 + 2m)\mu + \dots$$

$$(1 - \frac{\sum x_i \tau^2 + m^2}{n\tau^2 + \sigma^2})^2$$

$$\hat{\mu}_{\text{Bayes}} = \frac{\frac{1}{n} \sum x_i \tau^2 + \frac{1}{n} m \sigma^2}{\frac{\tau^2}{n} + \frac{\sigma^2}{n}}$$

$$\frac{2\tau^2 \nu^2}{n\tau^2 + \sigma^2}$$

$$\text{Var}(\mu | x) = \frac{\sigma^2 \nu^2}{n\tau^2 + \sigma^2}$$

Evaluating Point Estimators.

Mean Squared Error

$$\begin{aligned} E_{\theta} (\hat{\theta} - \theta)^2 &= \left(E_{\theta} (\hat{\theta} - \theta) \right)^2 + \text{Var}_{\theta}(\hat{\theta}) \\ &= B^2 + V \end{aligned}$$

Log-likelihood: $LL(\theta) = \sum_i^n \log p(x_i; \theta)$

Score: $S(\theta) = \frac{\partial LL(\theta)}{\partial \theta}$

Fish Information: $I(\theta) = E [S(\theta) S(\theta)^T]$

Score function has mean 0

Proof:

$$\begin{aligned} E_{\theta}(S(\theta)) &= \sum_x^n \nabla_{\theta} \log p(x_i; \theta) p(x_i; \theta) dx_i \\ &= n \underbrace{\int \nabla_{\theta} \log p(x; \theta) p(x; \theta) dx}_{\text{dominated convergence}} \\ &= \int \frac{\nabla_{\theta} p(x; \theta)}{p(x; \theta)} p(x; \theta) dx \\ &= \nabla_{\theta} \int p(x; \theta) dx = \nabla_{\theta} 1 = 0 \quad \text{theorem} \end{aligned}$$

$I_{11}(\theta) = E [-\nabla_{\theta}^2 \log p(x; \theta)]$

Proof:

$$\begin{aligned} \nabla_{\theta}^2 \log p(x; \theta) &= \nabla_{\theta} \frac{\nabla_{\theta} p(x; \theta)}{p(x; \theta)} \\ &= \frac{\nabla_{\theta} p(x; \theta)}{p(x; \theta)} - \frac{\nabla_{\theta} p(x; \theta) \nabla_{\theta} p(x; \theta)^T}{p(x; \theta)^2} \\ &= \frac{\nabla_{\theta} p(x; \theta)}{p(x; \theta)} - S(\theta) S(\theta)^T \end{aligned}$$

$\rightarrow \rightarrow \rightarrow \rightarrow \rightarrow \rightarrow$

$$\begin{aligned}
 E \nabla_{\theta} \log p(x; \theta) &= E \frac{\nabla_{\theta} \log p(x; \theta)}{p(x; \theta)} - [E s(\theta)] s(\theta)^T \\
 &= -[E s(\theta)] s(\theta)^T \\
 &= -I(\theta)
 \end{aligned}$$

Cramér-Rao Bound

$\hat{\theta}$ unbiased estimator

$$\text{Var}(\hat{\theta}) \geq \frac{1}{n I_{11}(\theta)}$$

Multivariate Generalization

$$\text{Var}(\hat{\theta}) \geq I(\theta)^{-1}$$

KL loss:

$$KL(p(x; \theta), p(x; \alpha)) = E_{x \sim p(\theta)} \log \left(\frac{p(x; \theta)}{p(x; \alpha)} \right)$$

$$\text{Risk: } R(\theta, \hat{\theta}(x)) = E_{x \sim \theta} L(\theta, \hat{\theta}(x))$$

$$\text{Bayes Risk: } B_n(\hat{\theta}) = \int R(\theta, \hat{\theta}) \pi(\theta) d\theta.$$

$\alpha = \beta = \sqrt{n}$, $R(\hat{\theta}_B, \theta) = \frac{n}{4(n + \sqrt{n})}$

$\hat{\theta}_{\text{Bayes}} \text{ minimize } r(\hat{\theta} x^n)$	posterior risk $r(\hat{\theta} x^n) = \int L(\theta, \hat{\theta}(x^n)) \pi(\theta x^n) d\theta.$
---	---

$$\text{Proof: } B_n(\hat{\theta}) = \int R(\theta, \hat{\theta}) \pi(\theta) d\theta = \int \left(\int L(\theta, \hat{\theta}) p(x | \theta) dx^n \right) \pi(\theta) d\theta$$

$$= \iint L(\theta, \hat{\theta}(x^n)) \underline{p(x, \theta)} dx^n d\theta$$

$$= \iint L(\theta, \hat{\theta}(x^n)) \pi(\theta | x^n) m(x^n) dx^n d\theta$$

$$= \iint L(\theta, \hat{\theta}(x^n)) \pi(\theta | x^n) dx^n d\theta$$

$$= \iint r(\hat{\theta} | x^n) m(x^n) dx^n$$

$\hat{\theta}_{\text{Bayes}} \text{ minimize } r(\hat{\theta} | x^n)$

If $L(\theta, \hat{\theta}) = |\theta - \hat{\theta}|$, then $\hat{\theta}_{\text{Bayes}} = \overbrace{\mathbb{E}(\hat{\theta} | x^n)}$

$$\int |\theta - \hat{\theta}| p(\theta | x) d\theta$$

$$= \int (\theta - \hat{\theta}) p(\theta | x) d\theta \Rightarrow \hat{\theta} = \int \theta p(\theta | x) d\theta = E_{\theta}(\theta | x)$$

The risk: $R(\theta, \hat{\theta}) = \mathbb{E}_{x \sim p(x)} L(\theta, \hat{\theta}) = \int_X L(\theta, \hat{\theta}(x)) p(x; \theta) dx_n$

when $L(\theta, \hat{\theta})$ is squared loss. MSE is the risk,

Example: Comparing risk functions

$$X_1, \dots, X_n \sim \text{Ber}(p)$$

$$\hat{p}_1 = \bar{x}, \quad \hat{p}_2 = \frac{n\bar{x} + \lambda}{n + \lambda + \beta}$$

$$r(p, \hat{p}_1) = \text{Var}(\bar{x}) = \frac{p(1-p)}{n}$$

$$r(p, \hat{p}_2) = \left(\frac{\lambda(1-p) - \beta p}{n + \lambda + \beta} \right)^2 + \frac{n p(1-p)}{(n + \lambda + \beta)^2}$$

$$\text{Let } \lambda = \beta = \sqrt{\frac{n}{4}}, \quad R(p, \hat{p}_2) = \frac{n}{4(n + \sqrt{n})^2}$$

$$\text{maximum risk: } \bar{R}(\hat{\theta}) = \max_{\theta} R(\theta, \hat{\theta})$$

$$\text{Bayes risk: } B_n(\hat{\theta}) = \int R(\theta, \hat{\theta}) \pi(\theta) d\theta$$

$$P(\theta | x) = \frac{p(x | \theta) \pi(\theta)}{p(x)}, \quad \pi(\theta | x^n) = \frac{p(x^n | \theta) \pi(\theta)}{m(x^n)} \rightarrow \text{marginal distribution of }$$

Posterior risk: $r(\hat{\theta} | x^n) = \int_{\theta} L(\theta, \hat{\theta}(x^n)) \pi(\theta | x^n) d\theta$

* The difference between $R(\theta, \hat{\theta})$ and $r(\hat{\theta} | x^n)$

$$\text{Bayes risk: } B_n(\hat{\theta}) = \int_{\theta} R(\theta, \hat{\theta}) \pi(\theta) d\theta$$

$$\begin{aligned} &= \int_{\theta} \int_x L(\theta, \hat{\theta}) p(x | \theta) dx \pi(\theta) d\theta \\ &= \int_{\theta} \int_x L(\theta, \hat{\theta}) p(x) dx d\theta \end{aligned}$$

$$\begin{aligned}
&= \int_{\Theta} \int_X L(\theta, \hat{\theta}) p(\theta|x) m(x) dx d\theta \\
&= \int_X \int_{\Theta} L(\theta, \hat{\theta}) p(\theta|x) d\theta m(x) dx \\
&= \int_X r(\theta|x^n) m(x^n) dx \quad \star
\end{aligned}$$

Bayes estimator : minimize Bayes risk (definition)
minimize $r(\theta|x^n)$ posterior risk

$$\begin{aligned}
\hat{\theta}_{\text{Bayes}} &\quad L(\theta, \hat{\theta}) = (\theta - \hat{\theta})^2, \quad r(\hat{\theta}|x^n) = \int (\theta - \hat{\theta})^2 \pi(\theta|x^n) d\theta \\
&\quad \nabla r(\hat{\theta}|x^n) = \int -2(\theta - \hat{\theta}) \pi(\theta|x^n) d\theta = 0 \\
&\quad \hat{\theta} = \mathbb{E}(\theta|x^n) \\
L(\theta, \hat{\theta}) &= |\theta - \hat{\theta}| \quad r(\hat{\theta}|x^n) = \int_{-\infty}^{\hat{\theta}} (\hat{\theta} - \theta) \pi(\theta|x^n) d\theta \\
&\quad + \int_{\hat{\theta}}^{\infty} (\theta - \hat{\theta}) \pi(\theta|x^n) d\theta \\
\nabla r(\hat{\theta}|x^n) &= 0 \Rightarrow \int_{-\infty}^{\hat{\theta}} \pi(\theta|x^n) d\theta = \\
&\quad \int_{-\infty}^{\hat{\theta}} \pi(\theta|x^n) d\theta = \frac{1}{2} \\
\hat{\theta} &= \text{median}(\theta|x^n) \\
L(\theta, \hat{\theta}) &= \mathbb{1}(\theta \neq \hat{\theta}) \quad r(\hat{\theta}|x^n) = \int \mathbb{1}(\theta \neq \hat{\theta}) \pi(\theta|x^n) d\theta \\
&= 1 - \pi(\hat{\theta}|x^n) d\theta \\
\hat{\theta} &= \text{mode}(\theta|x^n)
\end{aligned}$$

Minimax Estimator through Bayes Estimator

→ Bounding the Minimax Risk.

$$B_{\pi}(\hat{\theta}_{\text{Bayes}}) \leq B_{\pi}(\hat{\theta}_{\text{minimax}}) \leq \inf_{\tilde{\theta}} \sup_{\theta \in \Theta} R(\theta, \tilde{\theta}) \leq \sup_{\theta \in \Theta} R(\theta, \hat{\theta}_{\text{Bayes}})$$

Example : $X_1, \dots, X_n \sim N(\theta, \text{Id})$

$\hat{\theta} = \frac{1}{n} \sum_{i=1}^n X_i$ show that $\hat{\theta}$ is minimax estimator.

$$\hat{\theta} \sim N(\theta, \frac{I_d}{n})$$

$$R(\theta, \hat{\theta}) = E \left[\sum_{i=1}^d (\hat{\theta}_i - \theta_i)^2 \right] = E \left[\sum_{i=1}^d z_i^2 \right].$$

$$z_i \sim N(0, \frac{1}{n})$$

$$\inf_{\hat{\theta}} \sup_{\theta \in \Theta} R(\theta, \hat{\theta}) \leq R(\theta, \hat{\theta}) = \frac{d}{n}$$

$$\text{take prior } \pi = N(0, c^2 I_d) \quad ?$$

→ Least Favorable Prior

$$R(\theta, \hat{\theta}) \leq B_n(\hat{\theta}) \text{ for all } \theta$$

π is least favorable prior, $\hat{\theta}$ is minimax

MLE Asymptotics

Consistency: $\hat{\theta}_{MLE} \xrightarrow{P} \theta$

Asymptotic distribution $\sqrt{n}(\hat{\theta}_{MLE} - \theta) \xrightarrow{d} N(0, \frac{1}{I(\theta)})$

MLE. maximize $\log L(\theta)$

equals minimize $\frac{1}{n} \sum \frac{\log p(x_i; \theta)}{\log p(x_i; \hat{\theta})} \rightarrow \text{empirical risk}$

population risk $\in E_{\theta} \underbrace{\frac{\log p(x; \theta)}{\log p(x; \hat{\theta})}}_{KL \text{ divergence.}}$

Conditions for consistency:

1. Identifiability: if $\theta_1 \neq \theta_2$, then $p(x; \theta_1) \neq p(x; \theta_2)$

2. Strong identifiability: $\forall \epsilon > 0, \inf_{\hat{\theta}: |\hat{\theta} - \theta| \geq \epsilon} KL(p(x; \theta) \| p(x; \hat{\theta})) > \epsilon$

3. Uniform LLN: $\sup_{\hat{\theta}} |R_n(\theta, \hat{\theta}) - R(\theta, \hat{\theta})| \xrightarrow{P} 0$

Asymptotic: $\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, \frac{1}{I(\theta)})$

Proof:

$$\nabla = LL'(\hat{\theta}) = LL'(\theta) + (\hat{\theta} - \theta)LL''(\hat{\theta})$$

$$\hat{\theta} - \theta = -\frac{LL'(\theta)}{LL''(\hat{\theta})}$$

$$\sqrt{n}(\hat{\theta} - \theta) = -\frac{\frac{LL'(\theta)}{n}}{\frac{LL''(\hat{\theta})}{n}}$$

$$\text{numerator: } \frac{LL'(\theta)}{\sqrt{n}} = \sqrt{n} \times \frac{1}{n} \sum \nabla_{\theta} \log p(x_i; \theta) \\ = \sqrt{n} \times \left(\frac{1}{n} \sum \nabla_{\theta} \log p(x_i; \theta) - E[\nabla_{\theta} \log p(x_i; \theta)] \right)$$

$$\text{by CLT, and } E(S(\theta)) \Big|_{\theta} = 0 \xrightarrow{d} N(0, \text{Var}(S(\theta))) \xrightarrow{d} N(0, I(\theta))$$

$$\text{denominator: } -\frac{LL''(\hat{\theta})}{n} = -\frac{1}{n} \sum \nabla_{\theta}^2 \log p(x_i; \hat{\theta}) \xrightarrow{d} I(\theta)$$

$$\hat{\theta} \xrightarrow{P} \theta$$

by Slutsky's

$$\sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} \frac{1}{I(\theta)} N(0, I(\theta)) \xrightarrow{d} N(0, \frac{1}{I(\theta)})$$

Example: $x_1, \dots, x_n \sim \text{Exp}(\theta)$

$$p(x) = \theta e^{-\theta x}$$

$$LL(\theta) = n \log \theta - \theta \sum x_i$$

$$S(\theta) = \frac{n}{\theta} - \sum x_i$$

$$I(\theta) = E\left[-\frac{n}{\theta^2}\right] = \frac{n}{\theta^2}$$

$$\hat{\theta}_{MLE} = \frac{\sum x_i}{n}, \sqrt{n}(\hat{\theta} - \theta) \xrightarrow{d} N(0, \frac{1}{I(\theta)})$$

Counterexample..

Uniform distribution.

Definitions: minimal representation: no set of $a \in \mathbb{R}^s$

$$\sum_{i=1}^s a_i T_i(x) = \text{const. for all } x \in X$$

Non-minimal exponential families are over-complete.

not identifiable.

full-rank exponential family: space of θ_i is s -dim.

\Downarrow sufficient statistics are minimal

Minimal Sufficiency:

$T(x_1, \dots, x_n)$ is sufficient, and for any other sufficient statistic

$S(x_1, \dots, x_n)$. we can write $T(x_1, \dots, x_n) = g(S(x_1, \dots, x_n))$

Condition:

$$R(x_1, \dots, x_n; y_1, \dots, y_n; \theta) = \frac{P(x_1, \dots, x_n; \theta)}{P(y_1, \dots, y_n; \theta)}$$

does not depend on θ

$$\text{iff } T(x_1, \dots, x_n) = T(y_1, \dots, y_n)$$

Inconsistency of the MLE:

- not identifiable

- the parameter space is too large, fail of Uniform

MLE under misspecification

KL

when q does not belong to our model \mathcal{H}

$$KL(q \parallel P_{\theta \text{ MLE}}) \leq KL(q \parallel P_\theta) \text{ for all } \theta \in \mathcal{H}$$

MLE is estimating the KL projection of q onto our model.