

Summary - Fundamentals of Mathematical Statistics

Estimator: Given observations $X = (X_1, \dots, X_n)$, realizations of a random variable X with distribution $P_\theta \in \{P_\theta | \theta \in \Theta\}$. Then $T(X)$ is an estimator. $T(\cdot)$ is measurable, and not allowed to depend on unknown parameters. (i.e. Θ) (aka. statistic or decision).

Method of moments: $X \in \mathbb{R}$, let $\theta \in \Theta \subseteq \mathbb{R}^p$ be the parameter of interest. The first p moments of X are given by

$$\mu_j(\theta) = E_\theta[X^j] = \int x^j dF_\theta(x) \quad j=1, \dots, p$$

Further assume that $m: \Theta \rightarrow \mathbb{R}^p$, $m(\theta) \mapsto [\mu_1(\theta), \dots, \mu_p(\theta)]$ has an inverse $m^{-1}(\mu_1, \dots, \mu_p) \mapsto M_1, \dots, M_p$.

By estimating the μ_j by their sample counterparts

$$\hat{\mu}_j = \frac{1}{n} \sum_{i=1}^n X_i^j \left(= \int x^j d\hat{F}_n(x) \right) \quad j=1, \dots, p$$

we obtain the method of moments estimator

$$\hat{\theta} := m^{-1}(\hat{\mu}_1, \dots, \hat{\mu}_p)$$

example: negative Binomial distribution $\theta \in (0, 1)$ (thus $p=1$)

$$P_\theta(X=x) = \binom{k+x-1}{x} \theta^k (1-\theta)^x \quad x \in \{0, 1, 2, \dots\}$$

$$\begin{aligned} \text{then } E_\theta[X] &= k \frac{1-\theta}{\theta} := m(\theta) \Rightarrow m^{-1}(\mu) = \frac{k}{\mu+k} \\ \Rightarrow \hat{\theta} &= \frac{k}{\hat{\mu}+k} = \frac{k}{\bar{x}+k} = \frac{nk}{\sum_{i=1}^n X_i + nk} \end{aligned}$$

example: Pareto distribution $P_\theta(x) = \theta (1+x)^{-\theta}$, $x > 0$

$$\begin{aligned} E_\theta[X] &= \frac{1}{\theta-1} := m(\theta) \quad (>0 \text{ by definition; Thus } \theta > 1) \\ \Rightarrow m^{-1}(\mu) &= \frac{1+\mu}{\mu} \Rightarrow \hat{\theta} = \frac{1+\bar{X}}{\bar{X}} \end{aligned}$$

Maximum Likelihood Estimator (MLE): For data $X = (X_1, \dots, X_n)$, the MLE is given by

$$\hat{\theta} := \arg \max_{\theta \in \Theta} L_X(\theta)$$

where $L_X(\theta): \Theta \rightarrow \mathbb{R}$ is the Likelihood function, given by

$$L_X(\theta) := \prod_{i=1}^n P_\theta(X_i) \quad \theta \in \Theta.$$

Alternatively, the MLE can be written as the maximizer of the log-likelihood:

$$\hat{\theta} = \arg \max_{\theta \in \Theta} \log(L_X(\theta)) = \arg \max_{\theta \in \Theta} \sum_{i=1}^n \log(P_\theta(X_i))$$

(this is easier to differentiate for e.g. finding a maximum of the fct.)

↓(MLE)

maxima • The likelihood function may have local maxima, is not always unique,
or may not exist (unbounded)

The MLE is a plug-in estimator of θ , because we can write

$$\hat{\theta} = \arg \max_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^n \log(P_\theta(X_i)) \Rightarrow \arg \max_{\theta \in \Theta} E_\theta[\log(P_\theta(X))] = \theta$$

this is an empirical expected value

example: Pareto distribution, $\theta \in (0, \infty)$, $P_\theta(x) = \theta(1+x)^{-1+\theta}$ $x > 0$
For $X = (X_1, \dots, X_n)$ we maximize the function

$$\log(L_X(\theta)) = \log\left(\prod_{i=1}^n P_\theta(X_i)\right) = \sum_{i=1}^n \log(P_\theta(X_i)) \quad \text{w.r.t. } \theta$$

we take the derivative $D = \sum_{i=1}^n \left(\frac{1}{\theta} - \log(1+X_i)\right)$

$$\Rightarrow \hat{\theta} = \frac{n}{\sum_{i=1}^n \log(1+X_i)}$$

example: standard normal distribution, $\theta = (\mu, \sigma^2) \in \mathbb{R} \times (0, \infty)$

$$P_\theta(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right) \quad x \in \mathbb{R}$$

the respective MLE's are given by

$$\hat{\mu} = \bar{X}, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$$

example: (MLE does not exist). Each observation is distributed either $\mathcal{N}(\mu, 1)$ or $\mathcal{N}(\mu, \sigma^2)$, each with probability $\frac{1}{2}$. $\theta = (\mu, \sigma^2)$

$$P_\theta(x) = \frac{1}{2} \phi(x-\mu) + \frac{1}{2\sigma} \phi\left(\frac{x-\mu}{\sigma}\right) \quad \phi: \text{density of normal distribution}$$

then $L_X(\tilde{\mu}, \tilde{\sigma}^2) = \prod_{i=1}^n \left(\frac{1}{2} \phi(x_i - \tilde{\mu}) + \frac{1}{2\tilde{\sigma}} \phi\left(\frac{x_i - \tilde{\mu}}{\tilde{\sigma}}\right) \right)$

We can show that L_X for certain $\tilde{\mu}, \tilde{\sigma}^2$ is unbounded (thus no maximum exists)
take $\tilde{\mu} = x_1$: $L_X(x_1, \tilde{\sigma}^2) = \underbrace{\frac{1}{\sqrt{2\pi}} \left(\frac{1}{2} + \frac{1}{2\tilde{\sigma}}\right)}_{\rightarrow +\infty} \underbrace{\prod_{i=2}^n \left(\frac{1}{2} \phi(x_i - x_1) + \frac{1}{2\tilde{\sigma}} \phi\left(\frac{x_i - x_1}{\tilde{\sigma}}\right)\right)}_{> 0 \ (*)}$

(*): further note that $\forall z \neq 0 \lim_{\tilde{\sigma} \rightarrow 0} \frac{1}{\tilde{\sigma}} \phi\left(\frac{z}{\tilde{\sigma}}\right) = 0$, thus

$$\lim_{\tilde{\sigma} \rightarrow 0} (*) = \prod_{i=2}^n \frac{1}{2} \phi(x_i - x_1) > 0$$

$$\Rightarrow \lim_{\tilde{\sigma} \rightarrow 0} L_X(x_1, \tilde{\sigma}^2) = \infty.$$

Conditional Distributions (background knowledge)

The conditional probability of set A given set B is

$$P(A|B) = \frac{P(A \cap B)}{P(B)} \quad \text{where } P(B) \neq 0.$$

Law of total probability: For B_j a partition of Ω , then

$$P(A) = \sum_j P(A|B_j) \cdot P(B_j)$$

Marginal densities: $f(x, y)$, $(x, y) \in \mathbb{R}^n \times \mathbb{R}^m$ be a density

$$f_X(x) = \int_{y \in \mathbb{R}^m} f(x, y) dy \quad ; \quad f_Y(y) = \int_{x \in \mathbb{R}^n} f(x, y) dx, \quad y \in \mathbb{R}^m$$

conditional density: The cond. density of X given $Y=y$ is

$$f_{X|Y}(x|y) := \frac{f_{X,Y}(x,y)}{f_Y(y)}, \quad x \in \mathbb{R}^n$$

thus we get, similarly to the law of total probability

$$f_X(x) = \int f_{X|Y}(x|y) f_Y(y) dy$$

conditional expectation: Let $g: \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}$. The conditional expectation of $g(X, Y)$ given $Y=y$ is

$$E[g(X, Y) | Y=y] := \int f_{X|Y}(x|y) \cdot g(x, y) dx \quad (= h(y))$$

thus note $E[g_1(X) g_2(Y) | Y=y] = g_2(y) E[g_1(X) | Y=y]$

→ this is a function dependent on y .

→ this defines a random variable $h(y)$

Law of iterated expectation:

$$E[E[g(X, Y) | Y]] = E[g(X, Y)]$$

proof: $h(y) = E[g(X, y) | Y=y]$, $E[h(y)] = \int h(y) f_Y(y) dy = \int E[g(X, y) | Y=y] f_Y(y) dy$
 $= \int \int g(x, y) f_{X|Y}(x, y) dx dy = E[g(X, Y)]$

Multinomial distribution: $(N_1, \dots, N_k) \sim \text{Multinomial}(n, p_1, \dots, p_k)$

where $\sum p_i = 1$

$$P(N_1=n_1, \dots, N_k=n_k) = \underbrace{\binom{n}{n_1 \dots n_k}}_{:= \frac{n!}{n_1! \dots n_k!}} p_1^{n_1} \cdot \dots \cdot p_k^{n_k} \quad (n=n_1+\dots+n_k)$$

(multinomial coefficient)

Poisson distribution: $X \sim \text{Poisson}(\lambda)$, $\lambda > 0$, $x \in \{0, 1, 2, \dots\}$

$$P(X=x) = e^{-\lambda} \frac{\lambda^x}{x!}$$

• For $X \sim \text{Poi}(\lambda_1)$, $Y \sim \text{Poi}(\lambda_2)$ independent: $Z = X+Y \sim \text{Poi}(\lambda_1 + \lambda_2)$
proof: follows with binomial theorem $(x+y)^n = \sum_{k=0}^n \binom{n}{k} x^{n-k} y^k$

• For $X_i \sim \text{Poisson}(\lambda_i)$, X_1, \dots, X_n independent, $Z = \sum_{i=1}^n X_i$. Then
 $(X_1, \dots, X_n) | Z=z \sim \text{Multinomial}(z, p_1, \dots, p_k)$, $p_i = \frac{\lambda_i}{\sum_{j=1}^n \lambda_j}$
proof: compute explicitly $P(X_1, \dots, X_n | z) = \frac{P(X_1, \dots, X_n)}{P(Z)}$

$\max(X_1, X_2)$ -distribution: Let X_1, X_2 i.i.d., $Z = \max\{X_1, X_2\}$.

The density of Z is $f_Z(z) = 2 F_X(z) \cdot f_X(z)$

↑
Cumul. distr.

$$\text{proof: } F_Z(z) = P(Z \leq z) = P(\max\{X_1, X_2\} \leq z)$$

$$= (F_X(z))^2 \quad \text{then take derivative.}$$

Sufficiency

We write X for (X_1, \dots, X_n) copies of a random variable with distribution $P_\theta \{ P_\theta | \Theta \in \Theta \}$.

A statistic S is sufficient for $\Theta \in \Theta$, if for all Θ , and all possible S the cond. distribution $P_\theta(X=x | S=s)$ does not depend on Θ .

example: Bernoulli distribution, $\Theta \in \{0,1\}$, $P_\theta(X_i=1) = 1 - P(X_i=0) = \theta$. Take $S = \sum_{i=1}^n X_i$. Then S is sufficient for Θ , for all possible s .
 $P_\theta(X_1=x_1, \dots, X_n=x_n | S=s) = \frac{P(X_1=x_1, \dots, X_n=x_n)}{P(S=s)} = \frac{\theta^{x_1}(1-\theta)^{1-x_1} \dots \theta^{x_n}(1-\theta)^{1-x_n}}{\binom{n}{s} s^s (1-\theta)^{n-s}} = \frac{1}{\binom{n}{s}}$

example: Poisson distribution, $\Theta \in (0, \infty)$, $S = \sum_{i=1}^n X_i$ is sufficient
 $P_\theta(X_1=x_1, \dots, X_n=x_n | S=s) = \frac{P(X_1=x_1, \dots, X_n=x_n)}{P(S=s)} = \dots = \binom{s}{x_1 \dots x_n} \left(\frac{1}{\theta}\right)^s$ (multinomial)

example: exponential distribution, $X_1, X_2 \sim \exp(\theta)$ independent
then $f_{X_1, X_2}(x; \theta) = \theta e^{-\theta x}$, $x > 0$. $S = X_1 + X_2$ is sufficient.
Note $f_S(s; \theta) = s\theta^2 e^{-\theta s}$, $s > 0$ (\rightarrow Gamma(2, θ) distribution)
Thus $f_{X_1, X_2}(X_1, X_2 | S=s) = \frac{1}{s!}$, so independent of θ .

example: Order statistics, (X_1, \dots, X_n) i.i.d. from continuous distribution F .
 $S = (X_{(1)}, \dots, X_{(n)})$ is sufficient for F . Because

$$P_\theta((X_1, \dots, X_n) = (x_1, \dots, x_n) | (X_{(1)}, \dots, X_{(n)}) = s) = \frac{1}{n!}$$

Factorization Theorem of Neyman: A statistic S is sufficient iff.

We can write $P_\theta(x) = g_\theta(S(x)) \cdot h(x) \quad \forall x, \theta$ $\rightarrow g_\theta$ is allowed to depend on θ
density with $g_\theta(\cdot), h(\cdot) \geq 0$ are arbitrary functions.

proof (discrete case): Suppose X only takes discrete values $\{a_1, a_2, \dots\} \setminus \Theta$.
Let Q_θ be the distribution of S , thus $Q_\theta(s) := \sum_j P_\theta(X=a_j)$ where $S(a_j)=s$

Then the conditional distribution is

$$P_\theta(X=x | S=s) = \frac{P_\theta(X=x)}{Q_\theta(s)}, \quad S(x)=s$$

" \Rightarrow ": Assume S is sufficient. Then the above term depends only on x , say $h(x)$, by definition of sufficiency. Thus
 $P_\theta(X=x) = h(x) \cdot Q_\theta(s)$, where $g_\theta(s) = Q_\theta(s)$, $s = S(x)$.

" \Leftarrow ": Insert $P_\theta(x) = g_\theta(S(x)) \cdot h(x)$ into the definition of $Q_\theta(s)$.

$$Q_\theta(s) = g_\theta(s) \sum_j h(a_j) \quad \text{where } S(a_j)=s$$

Then the conditional distribution is given by

$$P_\theta(X=x | S=s) = \frac{h(x)}{\sum_j h(a_j)} \cdot \frac{g_\theta(s)}{g_\theta(s)} \quad \text{which does not depend on } \theta. \quad \square$$

↳ sufficiency

example: uniform distribution with unknown endpoint, on $[0, \theta]$.

then $P_\theta(x_1, \dots, x_n) = \frac{1}{\theta^n} \mathbb{1}_{\{0 \leq \min\{x_1, \dots, x_n\} \leq \max\{x_1, \dots, x_n\} \leq \theta\}}$

$$= g_\theta(S(x_1, \dots, x_n)) \cdot h(x_1, \dots, x_n)$$

with $g_\theta(s) = \frac{1}{\theta^n} \mathbb{1}_{\{s \leq \theta\}}$ and $h(x_1, \dots, x_n) = \mathbb{1}_{\{0 \leq \min\{x_1, \dots, x_n\}\}}$
Thus $S = \max\{x_1, \dots, x_n\}$ is sufficient.

Exponential Families

$\{P_\theta | \theta \in \Theta\}$

A distribution P_θ belongs to the k -dimensional exponential family.
All P_θ can be written in the form

$$P_\theta(x) = \exp \left[\sum_{j=1}^k c_j(\theta) T_j(x) - d(\theta) \right] h(x)$$

for any functions T_j , c_j , d and h . Note that $\{T_j\}$, $\{c_j\}$ are not unique.

Note: The k -dim. Statistic $S(x) = (T_1(x), \dots, T_k(x))$ is sufficient for θ .

Note: If each x_1, \dots, x_n is from a k -dim. exp. Family, then so is $\tilde{X} = (x_1, \dots, x_n)$

$$\tilde{P}_\theta(x) = \prod_{i=1}^n P_\theta(x_i)$$

Canonical Form: $\{P_\theta | \theta \in \Theta\}$ is an exponential family in canonical form, if

$$P_\theta(x) = \exp \left[\sum_{j=1}^k \theta_j T_j(x) - d(\theta) \right] h(x)$$

Note: $d(\theta)$ is the normalizing constant, with

$$d(\theta) = \log \left(\int \exp \left(\sum_{j=1}^k \theta_j T_j(x) \right) h(x) d\nu(x) \right) \quad \rightarrow T(x) = \begin{pmatrix} T_1(x) \\ \vdots \\ T_k(x) \end{pmatrix}$$

Note: We write $\text{Cov}_\theta(T(X)) = E_\theta[T(X)T^T(X)] - E_\theta[T(X)]E_\theta[T^T(X)] \in \mathbb{R}^{K \times K}$

Lemma 4.5.1.: It holds that

$$E_\theta[T(X)] = \dot{d}(\theta) \quad \text{and} \quad \text{Cov}_\theta(T(X)) = \ddot{d}(\theta)$$

Proof:

$$\begin{aligned} \dot{d}(\theta) &= \frac{\partial}{\partial \theta} d(\theta) = \frac{\partial}{\partial \theta} \log \left(\int \exp(\theta^T T(x)) h(x) d\nu(x) \right) \\ &= \int \exp(\theta^T T(x)) T(x) h(x) d\nu(x) \cdot \frac{1}{\int \exp(\theta^T T(x)) h(x) d\nu(x)} = \exp(d(\theta)) \\ &= \int \exp(\theta^T T(x) - d(\theta)) T(x) h(x) d\nu(x) \\ &= \int P_\theta(x) T(x) d\nu(x) = E_\theta[T(X)] \\ \ddot{d}(\theta) &= \frac{\int \exp(\theta^T T(x)) T(x) T^T(x) h(x) d\nu(x)}{\int \exp(\theta^T T(x)) h(x) d\nu(x)} - \frac{(\int \exp(\theta^T T) T h d\nu(x)) \cdot (\int \exp(\theta^T T) T h d\nu)}{(\int \exp(\theta^T T) h d\nu)^2} \\ &= \int \exp(\theta^T T - d(\theta)) T T^T h d\nu - (\int \exp(\theta^T T - d(\theta)) T h d\nu) \cdot (\int \exp(\theta^T T - d(\theta)) T^T h d\nu) \\ &= \int T T^T P_\theta d\nu - (\int P_\theta T d\nu) \cdot (\int P_\theta T^T d\nu) \\ &= E_\theta[T T^T] - E_\theta[T] E_\theta[T^T] \\ &= \text{Cov}_\theta(T(X)) \end{aligned}$$

□

Intermezzo: Score function and Fisher information

The score function is given by $s_\theta(x) := \frac{d}{d\theta} \log(p_\theta(x))$

The Fisher information for estimating θ :

$$I(\theta) := \text{Var}_\theta(s_\theta(X))$$

Further, we have that $E_\theta[s_\theta(X)] = 0$ and $I(\theta) = -E_\theta[\dot{s}_\theta(X)]$

For a 1-dim. exponential family we have that

$$s_\theta(x) = c(\theta) T(x) - d(x) \quad \text{and} \quad E_\theta[T(x)] = \frac{\dot{c}(\theta)}{c(\theta)}$$

Combining the above equations, we can receive

$$I(\theta) = \ddot{d}(\theta) - \frac{\dot{d}(\theta)}{c(\theta)} \ddot{c}(\theta) \quad (\text{since } \text{Var}_\theta(s_\theta(X)) = -E_\theta[\dot{s}_\theta(X)])$$

Bias: $X = (X_1, \dots, X_n)$ observations from a distribution $p \in \{P_\theta | \theta \in \Theta\}$. ($\Theta \subseteq \mathbb{R}^p$)

Let $y := g(\theta) \in \mathbb{R}$ be the parameter of interest.

Let $T(X)$ be an estimator of y .

The bias of T is given by

$$\text{bias}_\theta(T) := E_\theta[T] - g(\theta)$$

We call the estimator T unbiased, if

$$\text{bias}_\theta(T) = 0 \quad \forall \theta \in \Theta.$$

Unbiasedness means that there is no systematic error $E_\theta[T] = g(\theta)$.

example: $X \sim \text{Binomial}(n, \theta)$, $\theta \in (0, 1)$. Then

$$E_\theta[T(X)] = \sum_{k=0}^n \binom{n}{k} \theta^k (1-\theta)^{n-k} T(k) =: q(\theta)$$

Since $q(\theta)$ is a polynomial of at most degree n , only parameters $g(\theta)$ of such form can be estimated unbiasedly. (NOT: $\sqrt{\theta}$, $\frac{1}{\theta}$, ...)

example: $X \sim \text{Poisson}(\theta)$. $E_\theta[T(X)] = \sum_{k=0}^{\infty} e^{-\theta} \frac{\theta^k}{k!} T(k) =: \bar{e}^\theta p(\theta)$

$p(\theta)$ is a power series in θ . Unbiased estimators have to be of such form.

example: X_1, \dots, X_n i.i.d. $\mathcal{N}(\mu, \sigma^2)$, $\theta = (\mu, \sigma^2) \in \mathbb{R} \times (0, \infty)$. Then

$$S^2 := \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \text{ is an unbiased estimator of } \sigma^2.$$

Mean Square Error: The MSE of an estimator T is

$$\text{MSE}_\theta(T) := E_\theta[(T(X) - g(\theta))^2]$$

We can do the following decomposition: "bias-variance decomposition"

$$\text{MSE}_\theta(T) = \text{bias}_\theta^2(T) + \text{Var}_\theta(T)$$

proof:

$$\begin{aligned} \text{MSE}_\theta(T) &= E_\theta[(T - g(\theta))^2] = E_\theta[(T - g(\theta) + \underbrace{E_\theta[T] - E_\theta[T]}_{\text{const}})^2] \\ &= \underbrace{E_\theta[(T - E_\theta[T])^2]}_{\text{Var}_\theta(T)} + \underbrace{(E_\theta[T] - g(\theta))^2}_{\text{bias}_\theta^2(T)} + 2(E_\theta[T] - g(\theta)) \underbrace{E_\theta[T - E_\theta[T]]}_{=0} \end{aligned}$$

Uniform Minimum Variance Unbiased (UMVU):

We call an unbiased estimator T^* UMVU, if for any other unbiased estimator T it holds

$$\text{Var}_\theta(T^*) \leq \text{Var}_\theta(T) \quad \forall \theta$$

not super relevant

(For any T unbiased and S sufficient, we have that
 $T^* := E[T|S]$ is an unbiased estimator ($E_\theta[T^*] = g(\theta)$))

Complete: We call a sufficient statistic S complete if the following implication holds:

$$E_\theta[h(S)] = 0 \quad \forall \theta \Rightarrow h(S) = 0 \quad P_\theta\text{-a.s.} \quad \forall \theta$$

with h any function, not depending on θ .

(Lemma 5.4.1) Completeness for exponential families: remember that for $\theta \in \Theta \subseteq \mathbb{R}^k$

$$P_\theta(x) = \exp\left[\sum_{j=1}^k c_j(\theta) T_j(x) - d(\theta)\right] h(x)$$

consider the set $\mathcal{C} := \{(c_1(\theta), \dots, c_k(\theta)) \mid \theta \in \Theta\} \subseteq \mathbb{R}^k$.

If \mathcal{C} contains an open ball in \mathbb{R}^k (i.e. is truly k -dim.), then the statistic $S := (T_1(x), \dots, T_k(x))$ is complete.

Lehmann-Scheffé Lemma: Let T be an unbiased estimator of $g(\theta)$, with finite variance for all θ . If further another statistic S is sufficient and complete, then

$$T^* := E(T|S) \quad \text{is UMVU.}$$

proof:

$$(1) T^* = T^*(S) \text{ is unbiased: } E_\theta[T^*] = E_\theta[E(T|S)] = E_\theta[T] = g(\theta)$$

(2) Show $\text{Var}_\theta(T^*) \leq \text{Var}_\theta(T) \quad \forall \theta$: equivalently $\text{Var}_\theta(E[T|S]) \leq \text{Var}_\theta(T) \quad \forall \theta$.

Note that this follows from the following equality for any y, z random variables: $\text{Var}(y) = \text{var}(E[y|z]) + E[\text{var}(y|z)]$

proof of equality:

$$\text{var}(E[y|z]) = E[E[y|z]^2] - (E[E[y|z]])^2$$

$$= \underbrace{E[E[y|z]^2]}_{(*)} - E[y]^2$$

and

$$E[\text{var}(y|z)] = E[E[y^2|z] - E[y|z]^2]$$

$$= E[y^2] - \underbrace{E[E[y|z]^2]}_{(*)}$$

Note that adding both these terms, $(*)$ cancels out, and we are left exactly with the variance

$$\text{Var}(y) = E[y^2] - E[y]^2$$

Substituting $T=y$ and $S=z$ proves the inequality, together with the fact that $E[\text{var}(y|z)] \geq 0$.

(3) UMVU: Let $T'(S)$ be any other unbiased estimator of $g(\theta)$.

$$\text{Then } E_\theta[T^*(S) - T'(S)] = E_\theta[T^*(S)] - E_\theta[T'(S)] = g(\theta) - g(\theta) = 0 \quad \forall \theta.$$

But since S is complete, we have that

$$h(S) = T^*(S) - T'(S) = 0 \quad P_\theta\text{-a.s.} \quad \forall \theta \Rightarrow T^*(S) = T'(S) \quad P_\theta\text{-a.s.} \quad \forall \theta.$$

(θ>0)

example: X_1, \dots, X_n i.i.d. Poisson(θ), $g(\theta) = e^{-\theta}$ (early failure probability)
 $T = \sum_{\{X_i=0\}} 1$ is an unbiased estimator, and $S = \sum_{i=1}^n X_i$ is a sufficient statistic.

We now check for completeness of S . Thus let h be arbitrary

$$E_\theta[h(S)] = \sum_{k=0}^{\infty} e^{-\theta} \frac{(\theta)^k}{k!} h(k) = 0 \quad \forall \theta.$$

It follows that $\sum_{k=0}^{\infty} \frac{(\theta)^k}{k!} h(k) = 0 \quad \forall \theta$. By Taylor expansion at zero: $f(k) = \sum_{k=0}^{\infty} \frac{x^k}{k!} f^k(0)$

For this sum to be zero for all x (here of the form $n\theta + \theta$) all $f^k(0) = h(k)$ need to be equal to zero $\forall k$. Since S can only take values in \mathbb{N} , it follows that $h=0$ P_θ -q.s. and S is complete.

Now, by Lehmann-Scheffé, $T^* = E[T|S]$ is UMVU. Compute T^*

$$T^* = E[T|S] = 1 \cdot P(X_1=0|S) + 0 \cdot P(X_1=1|S)$$

$$\stackrel{\text{Bayes}}{=} \frac{P(S=s|X_1=0) \cdot P(X_1=0)}{P(S=s)} \quad (\text{note: } S \text{ is poisson distributed})$$

$$= \left(\frac{n-1}{n}\right)^s$$

example: X_1, \dots, X_n i.i.d. Uniform([0, θ]), $g(\theta) = \theta$. We know (previous example) that $S = \max\{X_1, \dots, X_n\}$ is sufficient. We now show that it is complete.

First find the probability density of S , $F_S(s) = P_\theta(\max\{X_1, \dots, X_n\} \leq s) = \left(\frac{s}{\theta}\right)^n$
 take the derivative to receive $f_S(s) = ns^{n-1}/\theta^n \quad 0 \leq s \leq \theta$

Thus $E_\theta[h(S)] = \int_0^\theta h(s) \frac{ns^{n-1}}{\theta^n} ds = 0 \quad \forall \theta$

so $\int_0^\theta h(s) s^{n-1} ds = 0 \quad \forall \theta$ differentiating wrt. to θ gives

$$h(\theta) \theta^{n-1} = 0 \quad \forall \theta \Rightarrow h=0. \text{ Thus } S \text{ is complec.}$$

Now use the fact from Lehmann-Scheffé that any unbiased estimator $T(S)$ is equal to the unique UMVU. Thus it suffices to find an unbiased estimator.

For this we check if S is unbiased:

$$E_\theta[S] = \int_0^\theta s \frac{ns^{n-1}}{\theta^n} ds = \frac{n}{n+1} \theta$$

S is biased, but we see, that we can easily remove bias with a factor, to receive $T(S) = \frac{n+1}{n} S$.

Score function: If $p_\theta(x)$ is differentiable for all x , we have

$$S_\theta(x) = \frac{d}{d\theta} \log(p_\theta(x)) = \frac{\dot{p}_\theta(x)}{p_\theta(x)} \quad \text{where } \dot{p}_\theta(x) = \frac{d}{d\theta} p_\theta(x)$$

Fisher Information / score function: X_1, \dots, X_n i.i.d. with density p_θ and $S_\theta = \dot{p}_\theta/p_\theta$.

The joint density is $p_\theta^{\text{joint}}(x) = \prod_{i=1}^n p_\theta(x_i)$

The score function of n observations is $S_\theta^{\text{joint}}(x) = \sum_{i=1}^n S_\theta(x_i)$

The Fisher information of n observations is

$$I(\theta) = \text{Var}_\theta(S_\theta^{\text{joint}}(x)) = \sum_{i=1}^n \text{Var}_\theta(S_\theta(x_i)) = n I(\theta)$$

$$\text{where } I(\theta) = \text{Var}_\theta(S_\theta(x)) = E_\theta[S_\theta(x)^2] - (E_\theta[S_\theta(x)])^2$$

Cramér-Rao Lower bound (CRLB): Let T be an unbiased estimator of $g(\theta)$ with finite variance. Then $g(\theta)$ is differentiable wrt. θ with $\dot{g}(\theta) = \frac{d}{d\theta} g(\theta) = \text{cov}(T, S_\theta(X))$

Also $\text{Var}_\theta(T) \geq \frac{(\dot{g}(\theta))^2}{I(\theta)} \quad \forall \theta.$

We call this the Cramér-Rao lower bound.

Proof: (1) first show differentiability of $g(\theta)$. Since T is unbiased

$$\begin{aligned}\frac{g(\theta+h) - g(\theta)}{h} &= \frac{E_{\theta+h}[T(X)] - E_\theta[T(X)]}{h} \\ &= \frac{1}{h} \int T \cdot (P_{\theta+h} - P_\theta) d\gamma = \int T \cdot \frac{P_{\theta+h} - P_\theta}{h} \frac{P_\theta}{P_\theta} d\gamma \\ &= E_\theta[T(X) \cdot \frac{P_{\theta+h}(X) - P_\theta(X)}{h P_\theta(X)} + T(X) S_\theta(X) - T(X) S_\theta(X)] \\ &= E_\theta[T(X) \cdot \left(\frac{P_{\theta+h}(X) - P_\theta(X)}{h P_\theta(X)} - S_\theta(X) \right)] + E_\theta[T(X) S_\theta(X)] \\ &\stackrel{?}{=} E_\theta[(T(X) - g(\theta)) \left(\frac{P_{\theta+h}(X) - P_\theta(X)}{h P_\theta(X)} - S_\theta(X) \right)] + E_\theta[T(X) S_\theta(X)] \\ \text{as } h \rightarrow 0 &\rightarrow 0 + E_\theta[T(X) S_\theta(X)] \quad (*)\end{aligned}$$

because by Cauchy-Schwarz we have $|E[XY]|^2 \leq E[X^2] \cdot E[Y^2]$
thus

$$\begin{aligned}(*) &\leq \underbrace{\text{Var}_\theta(T(X))}_{\substack{\text{const.} \\ \xrightarrow{h \rightarrow 0} 0}} \cdot E\left[\left(\frac{P_{\theta+h}(X) - P_\theta(X)}{h P_\theta(X)} - S_\theta(X)\right)^2\right] \\ &\xrightarrow{h \rightarrow 0} 0 \quad \hookrightarrow S_\theta(X).\end{aligned}$$

(2) This gives us also the derivative of $g(\theta)$

$$\dot{g}(\theta) = E_\theta[T(X) S_\theta(X)] = \text{cov}(T(X), S_\theta(X)) \quad (\text{because } E[S_\theta(X)] = 0)$$

(3) The CRLB follows from Cauchy-Schwarz

$$\begin{aligned}\dot{g}(\theta)^2 &= (E_\theta[T(X) S_\theta(X)])^2 \stackrel{?}{\leq} \text{Var}_\theta(T) \text{Var}_\theta(S_\theta(X)) \\ &= \text{Var}_\theta(T) \cdot I(\theta) \quad \square\end{aligned}$$

generally $\text{cov}(X, Y)^2 \leq \text{Var}(X) \cdot \text{Var}(Y)$

example: X_1, \dots, X_n i.i.d Exponential(θ), $\theta > 0$, $P_\theta(X) = \theta e^{-\theta X}$, $g(\theta) = \frac{1}{\theta}$

$T = \bar{X}$ is unbiased with $\text{Var}_\theta(T) = 1/(n\theta)^2$.

We now compute the CRLB $\frac{\dot{g}^2(\theta)}{I(\theta)^{\text{exact}}(\theta)} = \frac{\dot{g}^2(\theta)}{n I(\theta)}$

where $\dot{g}(\theta) = -\frac{1}{\theta^2}$, and

single observation

$$\begin{aligned}I(\theta) &= \text{Var}_\theta(S_\theta(X)) = \text{var}\left(\frac{d}{d\theta} \log(P_\theta(X))\right) = \text{var}\left(\frac{d}{d\theta}(\log(\theta) - \theta X)\right) \\ &= \text{var}\left(\frac{1}{\theta} - X\right) = \text{Var}(X) = \frac{1}{\theta^2}\end{aligned}$$

Thus $\text{CRLB} = \frac{1}{n\theta^2}$

↓ CRLB

example: X_1, \dots, X_n i.i.d. Poisson(θ), $\theta > 0$, then $g(\theta) = e^\theta$
 $\log(p_\theta(x)) = -\theta + x \log(\theta) - \log(x!)$ and $S_\theta(x) = -1 + \frac{x}{\theta} (= \frac{d}{d\theta} \log(p_\theta(x)))$
and $I(\theta) = \text{Var}_\theta\left(\frac{X}{\theta}\right) = \frac{\text{Var}(X)}{\theta^2} = \frac{1}{\theta}$
 $\rightarrow \text{CRLB} = \frac{\theta}{n} e^{-2\theta}$

Note: The UMVU estimator of $g(\theta)$ is $T := (1 - \frac{1}{n}) \sum_{i=1}^n X_i$.
It does not reach the CRLB, but the gap for large n is small.

CRLB and exponential families: CRLB can only be reached within exponential families.

Let T be unbiased for $g(\theta)$ and reaching the CRLB. Then $\{P_\theta | \theta \in \Theta\}$ forms a 1-dim. exp. family and $\exists c, d, h$ s.t. $\forall \theta$

$$P_\theta(x) = \exp(c(\theta)T(x) - d(\theta)) \cdot h(x) \quad \forall x$$

Also, c, d are diff.able with $g(\theta) = \frac{d(\theta)}{c(\theta)} \quad \forall \theta$.

NOT SURE (proof: relevant for exam?)

Tests and Confidence Intervals

Quantile functions: Let F be a cumulative distr. function on \mathbb{R} .
Then F is càdlàg (right-continuous, and limit from the left exists)
We define the quantile functions

$$q_{\text{sup}}^F(u) := \sup\{x \mid F(x) \leq u\}$$

$$q_{\text{inf}}^F(u) := \inf\{x \mid F(x) \geq u\} := F^{-1}(u) \quad \begin{matrix} \text{might not necessarily} \\ \text{exist} \end{matrix}$$

It further holds that

$$F(q_{\text{inf}}^F(u)) \geq u \quad \text{and} \quad F(q_{\text{sup}}^F(u) - h) \leq u \quad \forall h > 0.$$

(non-randomized) tests: Consider the model class $\{P_\theta \mid \theta \in \Theta\}$ and $g: \Theta \rightarrow \Gamma$, $\theta \mapsto g(\theta)$ the parameter of interest.
Let $\gamma_0 \in \Gamma$ and $\alpha \in [0, 1]$ be given a non-randomized test at level α for the hypothesis $H_0: \gamma = \gamma_0$
is a statistic $\Phi(X, \gamma_0) \in \{0, 1\}$

such that $P_\theta(\Phi(X, \gamma_0) = 1) \leq \alpha \quad \forall \theta \in \{\gamma \mid g(\gamma) = \gamma_0\}$
(i.e. prob that hypothesis rejects, $\Phi=1$, but γ_0 is actually true, is $\leq \alpha$)
 \uparrow iff.

Note: A test is typically based on a statistic T , of the form
 $\Phi(X) = \begin{cases} 0 & \text{if } T(X) > c \\ 1 & \text{else} \end{cases}$

c : "critical value".

Pivot: To test $H_{\gamma_0}: \gamma = \gamma_0$, we look for a Pivot. This is a function $Z(X, \gamma_1)$, such that $\forall \theta \in \Theta \quad P_\theta(Z(X, g(\theta)) \leq \cdot) =: G(\cdot)$
does not depend on θ .
A such pivot does not always exist.

\downarrow Pivot

However if a such Pivot $Z(X, \gamma_0)$, with cumulative distribution G , exists, then the test \leftarrow test statistic

$$\Phi(X, \gamma_0) := \begin{cases} 1 & \text{if } Z(X, \gamma_0) \notin [q_L, q_R] \\ 0 & \text{else} \end{cases}$$

$$\text{where } q_L := q_{\sup}^G\left(\frac{\alpha}{2}\right) \text{ and } q_R := q_{\inf}^G\left(1 - \frac{\alpha}{2}\right) \quad [\text{Quantiles}]$$

has level α for testing H_{γ_0} with $\gamma_0 = g(\theta_0)$, because

$$\begin{aligned} P_{\theta_0}(\Phi(X, g(\theta_0)) = 1) &= P_{\theta_0}(Z(X, g(\theta_0)) > q_R) + P_{\theta_0}(Z(X, g(\theta_0)) < q_L) \\ &= 1 - G(q_R) + G(q_L^-) \leq 1 - (1 - \frac{\alpha}{2}) + \frac{\alpha}{2} = \alpha. \end{aligned}$$

example: (location model) X_1, \dots, X_n i.i.d. copies of $X = \mu + \varepsilon$, $\varepsilon \sim F_0$.

$H_0: \mu = \mu_0$. Let $\hat{\mu}$ be an equivariant estimator (i.e. $\hat{\mu} - \mu$ distr. does not depend on μ)

Case 1: F_0 is known.

Take $Z(X, \gamma_0) = \hat{\mu} - \mu$ as pivot. By def. of equivariance, its distribution does not depend on θ .

Case 2: F_0 is of the form $\Phi(\frac{\cdot - \mu}{\sigma})$ for any $\sigma > 0$.

$$Z(X, \mu) := \frac{\sqrt{n}(\bar{X} - \mu)}{S_n} \quad \text{with } S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$$

Then G is the Student distribution with $n-1$ degrees of freedom.

If σ is known, then G is normally distributed

Case 3: F_0 is symmetric and continuous at $x=0$.

$$Z(X, \mu) := \sum_{i=1}^n \mathbb{1}_{\{X_i \geq \mu\}}$$

Then G is the Binomial(n, p) distribution with $p = \frac{1}{2}$.

Confidence set/interval: We call $I = I(X) \subseteq \Pi$ a confidence set for μ at level $1-\alpha$ if

$$P_{\theta}(\mu \in I) \geq 1-\alpha \quad \forall \theta \in \Theta$$

We call I a confidence interval, if it is of the form $I := [\underline{\mu}, \bar{\mu}]$, where the boundaries depend only on X .

Let for any $\gamma_0 \in \mathbb{R}$, $\Phi(X, \gamma_0) \in \{0, 1\}$ be a test at level α for the hypothesis $H_{\gamma_0}: \gamma_1 = \gamma_0$.

This means we reject H_{γ_0} iff $\Phi(X, \gamma_0) = 1$. It further holds that $P_{\theta: \gamma=\gamma_0}(\Phi(X, \gamma_0) = 1) \leq \alpha$.

Then

$$I(X) := \{\mu \mid \Phi(X, \mu) = 0\}$$

is a $(1-\alpha)$ -confidence set for μ .

Conversely, if $I(X)$ is a $(1-\alpha)$ -confidence set for μ , then for any γ_0 , the test

$$\Phi(X, \gamma_0) = \begin{cases} 1 & \text{if } \gamma_0 \notin I(X) \\ 0 & \text{else} \end{cases}$$

is a test at level α of H_{γ_0} .

Level and power of a test:

	H_0 is true	H_0 is false
Test rejects H_0	α (level, Type I error)	$1-\beta$ (power)
Test doesn't reject H_0	$1-\alpha$	β (Type II error)

• p-value: (assume H_0 is true) the probability of obtaining a result equal, or more extreme than what was actually observed.
(e.g. if p is large, then it was very likely to make this observation, given H_0 is true).

Two-sample problem: We assume our data consists of two independent samples $X = (X_1, \dots, X_n)$ and $Y = (Y_1, \dots, Y_m)$ (each i.i.d.) sampled from (possibly unknown) distributions F_X and F_Y .

The testing problem is: $H_0: F_X = F_Y$ against a one-/two-sided alternative. (sometimes only one side of the distribution is relevant to us)

Two-sample student test: We assume that both F_X and F_Y are normally distributed with the same variance (but possibly different mean). $F_X = N(\mu, \sigma^2)$, $F_Y = N(\mu + \gamma, \sigma^2)$ where $\gamma \in \mathbb{R}$.
two-sided: $\Pi = \mathbb{R}$, one-sided: $\Pi = (-\infty, 0]$

The testing problem is then $H_0: \gamma = 0$.

For constructing a pivot $Z(X, Y, \gamma)$ we consider the sample means \bar{X} (expectation μ , var. σ^2/n) and \bar{Y} (expect. $\mu + \gamma$, var. σ^2/m). Then the quantity $\bar{X} - \bar{Y}$ is $N(\gamma, \sigma^2(n+m)/nm)$ distributed. But this distribution still depends on unknown parameters.

$\sqrt{\frac{nm}{n+m}} \left(\frac{\bar{Y} - \bar{X} - \gamma}{\sigma} \right)$ is $N(0, 1)$ distributed.

To arrive at our pivot, we plug in our estimate for σ , given by

$$S^2 = \frac{1}{m+n-2} \left(\sum_{i=1}^n (X_i - \bar{X})^2 + \sum_{j=1}^m (Y_j - \bar{Y})^2 \right)$$

we receive

$$Z(X, Y, \gamma) := \sqrt{\frac{nm}{n+m}} \left(\frac{\bar{Y} - \bar{X} - \gamma}{S} \right)$$

The distribution G is Student($n+m-2$).

The test statistic for $H_0: \gamma = 0$ is thus given by

$$T = Z(X, Y, 0)$$

and we receive the one-sided test at level α (against $H_1: \gamma > 0$)

$$\Phi(X, Y) = \begin{cases} 1 & \text{if } T < -t_{n+m-2} (1-\alpha) \\ 0 & \text{if } T \geq -t_{n+m-2} (1-\alpha) \end{cases} \rightarrow (1-\alpha)\text{-quantile of Student}(n+m-2)$$

Two-Sample Wilcoxon's test: Assume F_X and F_Y are continuous.

Let $Z = (X_1, \dots, X_n, Y_1, \dots, Y_m)$, $N = n+m$. Define $R_i = \text{rank}(z_i)$ $\xrightarrow{i=1, \dots, N}$ index; sorted accordingly
The Wilcoxon test statistic is

$$Z(X, Y) = T = \sum_{i=1}^n R_i \quad (\text{i.e. sum ranks of all } X_1, \dots, X_n \text{ among } Z) \\ = \#\{Y_j < X_i\} + \frac{n(n+1)}{2}$$

Large values of T indicate that X is typically larger than Y .

Under $H_0: F_X = F_Y$ the vector of ranks (R_1, \dots, R_n) has the same distribution as n random draws without replacement from $\{1, \dots, N\}$.

$$E_{H_0}[T] = \frac{n(N+1)}{2}$$

Neyman Pearson and UMP: Let $\{P_\theta | \Theta \in \mathbb{H}\}$ be a family of probability measures. Let $\mathbb{H}_0, \mathbb{H}_1$ be a partition of \mathbb{H} . We consider the general testing problem

$H_0: \theta \in \mathbb{H}_0$
against
 $H_1: \theta \in \mathbb{H}_1$.

Neyman Pearson test: Consider the case $\mathbb{H} = \{\theta_0, \theta_1\}$, and thus the test $H_0: \theta = \theta_0$ against $H_1: \theta = \theta_1$.

First, define the Risk of a test ϕ

$$R(\theta, \phi) := \begin{cases} E_\theta[\phi(X)], & \theta = \theta_0 \\ 1 - E_\theta[\phi(X)], & \theta = \theta_1 \end{cases}$$

→ probability of rejecting when H_0 is true
(Error of first kind)

→ probability of accepting when H_0 is not true
(Error of second kind)

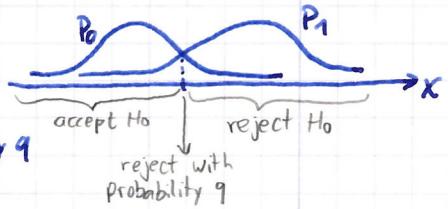
Generally we want the risk to be small.

$$\text{For a non-randomized test: } R(\theta, \phi) = \begin{cases} P_{\theta_0}(\phi(X)=1) & \theta = \theta_0 \\ 1 - P_{\theta_0}(\phi(X)=1) = P_{\theta_1}(\phi(X)=0) & \theta = \theta_1 \end{cases}$$

A Neyman Pearson Test is given by

$$\phi_{NP}(X) := \begin{cases} 1 & \frac{P_1(X)}{P_0(X)} > c \\ q & \frac{P_1(X)}{P_0(X)} = c \\ 0 & \frac{P_1(X)}{P_0(X)} < c \end{cases}$$

→ in this case we reject with probability q



For some constants $q \in [0, 1]$, $c \in [0, \infty)$.

Neyman Pearson Lemma: Take same hypotheses as above. For any test ϕ it holds true that

$$R(\theta_1, \phi) - R(\theta_1, \phi_{NP}) \geq c \cdot [R(\theta_0, \phi_{NP}) - R(\theta_0, \phi)]$$

= level of the test

Note that if ϕ_{NP} and ϕ have the same level, the right hand side is zero. Then the Risk (in case of θ_1) of ϕ_{NP} is strictly smaller than the Risk of ϕ .

Proof:

$$\begin{aligned} R(\theta_1, \phi) - R(\theta_1, \phi_{NP}) &= E_{\theta_1}[\phi_{NP}(X)] - E_{\theta_1}[\phi(X)] \\ &= \int (\phi_{NP} - \phi) P_1 \\ &= \int_{P_1/P_0 > c} (\phi_{NP} - \phi) P_1 + \int_{P_1/P_0 = c} (\phi_{NP} - \phi) P_1 + \int_{P_1/P_0 < c} (\phi_{NP} - \phi) P_1 \\ &\geq \int_{P_1/P_0 > c} (\phi_{NP} - \phi) \cdot c P_0 + \int_{P_1/P_0 = c} (\phi_{NP} - \phi) \cdot c P_0 + \int_{P_1/P_0 < c} (\phi_{NP} - \phi) \cdot c P_0 \\ &= c \cdot \int (\phi_{NP} - \phi) P_0 = c \cdot [E_{\theta_0}[\phi_{NP}(X)] - E_{\theta_0}[\phi(X)]] \\ &= c \cdot [R(\theta_0, \phi_{NP}) - R(\theta_0, \phi)] \quad \square \end{aligned}$$

Level of a (randomized) test: Let $\alpha \in [0, 1]$, $\phi: \mathcal{X} \rightarrow [0, 1]$. We say ϕ is a test at level α if

$$\sup_{\theta \in \mathbb{H}_0} E_\theta[\phi(X)] \leq \alpha$$

→ probability of rejecting, when $\theta \in \mathbb{H}_0$
(corresponds to this)

special case: non-randomized

$$E_\theta[\phi(X)] = P_\theta(\phi(X)=1)$$

Uniformly most powerful: We call a test ϕ UMP, if

- it has level α , and
- for all tests ϕ' with level α we have $E_\theta[\phi'(X)] \leq E_\theta[\phi(X)] \quad \forall \theta \in \Theta$ (i.e. ϕ rejects H_0 (correctly) more often than ϕ')

One-sided UMP test and exponential families: Let $\Theta \subseteq \mathbb{R}$ be an interval.

Regard the testing problem $H_0: \theta \leq \theta_0$ against $H_1: \theta > \theta_0$.
Let $\{P_\theta | \theta \in \Theta\}$ be a 1-dim. exponential family

$$P_\theta(x) = \exp[c(\theta) T(x) - d(\theta)] h(x)$$

further assume that $c(\theta)$ is a strictly increasing function of θ .
then

$$\phi(T(X)) = \begin{cases} 1 & T(X) > t_0 \\ q & T(X) = t_0 \\ 0 & T(X) < t_0 \end{cases}$$

is UMP.

q and t_0 are chosen s.t. $E_{\theta_0}[\phi(T)] = \alpha$.

$$(\theta_0 < \theta_1)$$

proof: Recall the NP-test for $H_0: \theta = \theta_0$ against $H_1: \theta = \theta_1$, where

$$\phi_{NP}(X) = \begin{cases} 1 & P_1(x)/P_0(x) > c_0 \\ q_0 & P_1(x)/P_0(x) = c_0 \\ 0 & P_1(x)/P_0(x) < c_0 \end{cases}$$

with q_0 and c_0 chosen s.t. $E_{\theta_0}[\phi_{NP}(X)] = \alpha$.

Note that then

$$\frac{P_{\theta_1}(x)}{P_{\theta_0}(x)} = \exp \left[\underbrace{(c(\theta_1) - c(\theta_0))}_{>0} T(x) - (d(\theta_1) - d(\theta_0)) \right]$$

> 0 because $\theta_1 > \theta_0$, and c increasing.

is increasing in $T(X)$. Thus $\left(\frac{P_1(x)}{P_0(x)} \right) \stackrel{>}{=} c \Leftrightarrow T(X) \stackrel{>}{=} t$ for some t and c

Thus ϕ is equivalent to ϕ_{NP} and so most powerful for $H_0: \theta = \theta_0$ against $H_1: \theta = \theta_1$.

Since above we chose $\theta_1 (> \theta_0)$ arbitrarily, the above statement holds true for all $\theta_1 > \theta_0$. Thus ϕ is UMP for $H_0: \theta = \theta_0$ against $H_1: \theta > \theta_0$.

Hence for the testing problem $H_0: \theta \leq \theta_0$ against $H_1: \theta > \theta_0$ ϕ fulfills the second condition for UMP.

It is left to show that ϕ has level α , i.e. $\sup_{\theta \in \Theta} E_\theta[\phi(X)] \leq \alpha$.

For this we define $\beta(\theta) = E_\theta[\phi(T)]$ and use the fact that β is increasing in θ together with $\beta(\theta_0) = \alpha$ (because $\phi = \phi_{NP}$). \rightarrow see above
Then $\sup_{\theta \in \Theta} \beta(\theta) = \beta(\theta_0) = E_{\theta_0}[\phi(X)] = \alpha$

and it follows that ϕ also has level α (for the relevant testing problem).
Thus it is UMP.

Left to prove: $\beta(\theta)$ is increasing. Let $\theta, \vartheta \in \Theta, \vartheta > \theta$.

We know that $\frac{\bar{P}_\vartheta(t)}{\bar{P}_\theta(t)} = \exp \left[\underbrace{(c(\vartheta) - c(\theta)) \cdot t}_{>0} - (d(\vartheta) - d(\theta)) \right]$ is increasing in t . \rightarrow new probability measure wrt. T thus different \bar{P}_θ .

Since further $1 = \int \bar{P}_{\vartheta}(t) \bar{V}(t) = \int \bar{P}_\theta(t) \bar{V}(t)$, there must be an So, where $\bar{P}_{\vartheta}(S_0) = \bar{P}_\theta(S_0)$ the densities cross. (and $\bar{P}_\vartheta(t) \leq \bar{P}_\theta(t)$ for $t \leq S_0$) \rightarrow and the other way around

$$\beta(\vartheta) - \beta(\theta) = \int \phi(t) (\bar{P}_\vartheta(t) - \bar{P}_\theta(t)) d\bar{V}(t)$$

$$= \int_{t \leq S_0} \phi(t) (\bar{P}_\vartheta(t) - \bar{P}_\theta(t)) d\bar{V}(t) + \int_{t > S_0} \phi(t) (\bar{P}_\vartheta(t) - \bar{P}_\theta(t)) d\bar{V}(t)$$

(note: ϕ is increasing in t) $\geq \phi(S_0) \int \bar{P}_\vartheta(t) - \bar{P}_\theta(t) d\bar{V}(t) = 0$

□

UMP

example: X_1, \dots, X_n i.i.d. $\text{Exp}(\theta)$, $\theta > 0$, $H_0: \theta \leq \theta_0$, $H_1: \theta > \theta_0$.

$$p_\theta(x) = \theta e^{-\theta x} = \exp[-\theta x - (-\log(\theta))] \cdot 1_{\{x > 0\}}(x)$$

$$\text{For } n \text{ observations: } \prod_{i=1}^n p_\theta(x_i) = \exp\left[-\theta \underbrace{\sum_{i=1}^n x_i}_{c(\theta)} - (n \log(\theta))\right] \cdot 1_{\{x_1, \dots, x_n > 0\}}$$

$$\Rightarrow \phi(T) = \begin{cases} 1 & \text{if } T \leq t_0 \\ 0 & \text{if } T > t_0 \end{cases} \rightarrow \text{no randomization because our distribution is continuous.}$$

Now choose t_0 s.t. $E_{\theta_0}[\phi(T)] = \alpha$ (level of α)

$$\text{i.e. } \alpha = P_{\theta_0}(T \leq t_0) = P(\theta_0 T \leq \theta_0 t_0) = G(\theta_0 t_0) = \alpha$$

how has mean 1

where $G = \text{Gamma}(n, 1)$ (i.e. sum of $\text{Exp}(1)$).
Thus $\theta_0 t_0 = G^{-1}(\alpha) \Rightarrow t_0 = G^{-1}(\alpha)/\theta_0$.

Then ϕ is UMP.

Right-/Left-/two-sided alternatives:

Right sided alternative: $H_0: \theta \leq \theta_0$, $H_1: \theta > \theta_0$. The UMP test is

$$\phi_R(T) = \begin{cases} 1 & T > t_R \\ q & T = t_R \\ 0 & T < t_R \end{cases} \quad \beta_R(\theta) = E_\theta[\phi_R(T)] \text{ is strictly increasing in } \theta.$$

Left sided alternative: $H_0: \theta \geq \theta_0$, $H_1: \theta < \theta_0$. The UMP test is

$$\phi_L(T) = \begin{cases} 1 & T < t_L \\ q & T = t_L \\ 0 & T > t_L \end{cases} \quad \beta_L(\theta) = E_\theta[\phi_L(T)] \text{ is strictly decreasing in } \theta.$$

Two sided alternative: $H_0: \theta = \theta_0$, $H_1: \theta \neq \theta_0$. ϕ_R is most powerful for $\theta > \theta_0$, and ϕ_L is most powerful for $\theta < \theta_0$.
Thus there exists no UMP test for this case.

Unbiased: We call a test ϕ unbiased, if $\forall \theta \in \Theta_0, \forall \gamma \in \Theta_1$

$$E_\theta[\phi(X)] \leq E_\gamma[\phi(X)]$$

Uniformly most powerful unbiased (UMPU): We call a test ϕ UMPU if

- ϕ has level α
- ϕ is unbiased
- For any ϕ' unbiased test with level α we have

$$E_\theta[\phi'(X)] \leq E_\theta[\phi(X)] \quad \forall \theta \in \Theta_1$$

UMPU and exponential families: Let $\Theta \subseteq \mathbb{R}$ be an interval. Consider testing $H_0: \theta = \theta_0$ against $H_1: \theta \neq \theta_0$. We have a 1-dim. exp. family $p_\theta(x) = \exp[c(\theta)T(x) - d(\theta)] h(x)$, $c(\theta)$ is strictly increasing.

Then

$$\phi(T(x)) = \begin{cases} 1 & \text{if } T(x) < t_L \text{ or } T(x) > t_R \\ q_L & \text{if } T(x) = t_L \\ q_R & \text{if } T(x) = t_R \\ 0 & \text{if } t_L < T(x) < t_R \end{cases}$$

where t_L, t_R, q_L, q_R are chosen such that

$$E_{\theta_0}[\phi(X)] = \alpha \quad \text{and} \quad \frac{d}{d\theta} E_\theta[\phi(X)] \Big|_{\theta=\theta_0} = 0$$

LUMPU

example: mean of $W(\mu, \sigma^2)$, X_1, \dots, X_n i.i.d., $\mu \in \mathbb{R}$ unknown.

we test $H_0: \mu = \mu_0$ against $H_1: \mu \neq \mu_0$. We want to find a UMPU test.
Since our distribution is continuous, our test is non-randomized.
Thus, we have no q_R and q_L from the previous theorem.
 $T = \sum_{i=1}^n X_i$ is a sufficient statistic.
for $t_L < t_R$:

$$E_\mu[\phi(T)] = P_\mu(T > t_R) + 1 \cdot P_\mu(T < t_L)$$

$$= P_\mu\left(\frac{T - n\mu}{\sqrt{n}\sigma_0} > \frac{t_R - n\mu}{\sqrt{n}\sigma_0}\right) + P_\mu\left(\frac{T - n\mu}{\sqrt{n}\sigma_0} < \frac{t_L - n\mu}{\sqrt{n}\sigma_0}\right) \rightarrow \text{std normally distributed.}$$
$$= 1 - \Phi\left(\frac{t_R - n\mu}{\sqrt{n}\sigma_0}\right) + \Phi\left(\frac{t_L - n\mu}{\sqrt{n}\sigma_0}\right) = \alpha$$

$$\text{and } \frac{d}{d\mu} E_\mu[\phi(T)] = \frac{n}{\sqrt{n}\sigma_0} \dot{\Phi}\left(\frac{t_R - n\mu}{\sqrt{n}\sigma_0}\right) - \frac{n}{\sqrt{n}\sigma_0} \dot{\Phi}\left(\frac{t_L - n\mu}{\sqrt{n}\sigma_0}\right) \Big|_{\mu=\mu_0} = 0$$

gives $\dot{\Phi}\left(\frac{t_R - n\mu_0}{\sqrt{n}\sigma_0}\right) = \dot{\Phi}\left(\frac{t_L - n\mu_0}{\sqrt{n}\sigma_0}\right)$

or equivalently: $(t_R - n\mu_0)^2 = (t_L - n\mu_0)^2$. We choose solution $(t_R - n\mu_0) = -(t_L - n\mu_0)$, because otherwise we get a trivial solution.
Plugging this into the first equation gives

$$\Phi\left(\frac{t_R - n\mu_0}{\sqrt{n}\sigma_0}\right) = 1 - \frac{\alpha}{2} \Rightarrow t_R = n\mu_0 - \sqrt{n}\sigma_0 \Phi^{-1}(1 - \frac{\alpha}{2}).$$

Risk of an Estimator: Let $X \sim P_\theta$, $\theta \in \Theta$. $T(X)$ an estimator of a parameter of interest $g(\theta) = g_T$.
A Risk function R measures the loss due to the error of the estimator.

$$R(\theta, T) := E_\theta[L(\theta, T(x))] \text{, with } L(\cdot, \cdot) \text{ a loss-function.}$$

Rao-Blackwell Lemma: Let S be a sufficient estimator for θ . Suppose $\mathcal{A} = \mathbb{R}^p$ is convex, and $L: \mathcal{A} \rightarrow \mathbb{R}$, $a \mapsto L(\theta, a)$ is a convex function for all θ . For any decision $d: \mathcal{X} \rightarrow \mathcal{A}$ we define $d'(S) = E[d(X)|S=s]$. Then

$$R(\theta, d') \leq R(\theta, d) \quad \forall \theta$$

Thus in case of convex loss, an estimator based on the original data X can be replaced by another estimator based only on S .

proof:

Recall Jensen inequality: $E[g(X)] \geq g(E[X])$ for g convex.
Thus $E\left[L(\theta, d(X))|S=s\right] \geq L(\theta, E[d(X)|S=s])$
 $= L(\theta, d'(S))$

And

$$R(\theta, d) = E_\theta[L(\theta, d(X))]$$

$$\xrightarrow{\text{iterated expectation}} = E_\theta\left[E\left[L(\theta, d(X))|S=s\right]\right]$$

$$\geq E_\theta[L(\theta, d'(S))] = R(\theta, d')$$

Location model: $\theta \in \mathbb{R}$ (location parameter). $X_i = \theta + \varepsilon_i$, $i=1, \dots, n$
 ε_i i.i.d. with distribution $P(\cdot)$. We want to estimate θ .

Location-scale model: $\theta = (\mu, \sigma)$, $\mu \in \mathbb{R}$, $\sigma \in (0, \infty)$ (scale parameter).

$$X_i = \mu + \sigma \varepsilon_i, i=1, \dots, n.$$

Location equivariant: We call a statistic $T = T(X)$ (location) equivariant if $\forall c \in \mathbb{R}$, $X = (X_1, \dots, X_n)$ we have

$$T(X_1 + c, \dots, X_n + c) = T(X_1, \dots, X_n) + c$$

(common examples are: \bar{X} , $X_{\lceil \frac{n}{2} \rceil}$ for n odd, ...)

Location invariant: We call a Loss function $L(\theta, a)$ (location) invariant, if $\forall c \in \mathbb{R}$

$$L(\theta + c, a + c) = L(\theta, a) \quad (\theta, a) \in \mathbb{R}^2.$$

Risk of equivariant statistic: If $T(X)$ is equivariant (and $L(\theta, a)$ invariant) then we have that (location model $X = \theta + \varepsilon$)

$$\begin{aligned} R(\theta, T) &= E_\theta [L(\theta, T(X))] = E_\theta [L(0, T(X) - \theta)] \\ &= E_0 [L(0, T(X) - \theta)] = E_0 [L(0, T(\varepsilon))] \quad (\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)) \end{aligned}$$

Thus, the risk does not depend on θ , and we can write

$$R(\theta, T) = R(0, T).$$

Uniform minimum risk equivariant (UMRE): We call an equivariant statistic T UMRE, if

$$R(\theta, T) = \min_{d \text{ equivariant}} R(\theta, d) \quad \forall \theta \quad (\text{minimal risk})$$

$$\text{or equivalently: } R(0, T) = \min_{d \text{ equivariant}} R(0, d)$$

Construction of UMRE estimator: Let $y_i := X_i - X_n$, $i=1, \dots, n$.
and $\bar{X} = (X_1, \dots, X_n)$, $y = (y_1, \dots, y_n)$. Then

$$T \text{ is equivariant} \Leftrightarrow T(\bar{X}) = T(y) + X_n$$

proof:

$$\Rightarrow: T(y) = T(X_1 - X_n, \dots, X_n - X_n) \stackrel{\text{equivariance}}{=} T(\bar{X}) - X_n.$$

\Leftarrow : Replace \bar{X} by $\bar{X} + c$ (y is unchanged, because invariant)
Then $T(\bar{X} + c) = T(y) + X_n + c = T(\bar{X}) + c$. \square

Moreover define

$$T^*(y) := \arg \min_v E[L(0, v + \varepsilon_n) | y]$$

and additionally

$$T^*(\bar{X}) := T^*(y) + X_n$$

Then T^* is UMRE.

$$\begin{aligned} R(0, T) &= E[L_0(T(y) + \varepsilon_n)] = E[E[L_0(T(y) + \varepsilon_n) | y]] \geq E[E[L_0(T^*(y) + \varepsilon_n) | y]] \\ &\quad \text{iterated expectation} \quad \stackrel{\text{def.}}{=} R(0, T^*) \end{aligned}$$

Invariant statistic: We call $T: \mathbb{R}^n \rightarrow \mathbb{R}^n$ maximal invariant, if
 $y(x) = y(x') \Leftrightarrow \exists c \in \mathbb{R}: x = x' + c \quad x, x' \in \mathbb{R}^n$
 $(c \text{ can depend on } x \text{ and } x')$

More general UMRE: Let $d(X)$ be an equivariant estimator.
Let $\gamma = X - d(X)$, and

$$T^*(Y) := \arg \min_Y E [L_0(Y + d(\varepsilon))] | Y]$$

Then

$$T^*(X) := T^*(Y) + d(X) \quad \text{is UMRE.}$$

Basu's Lemma: Let X have distribution P_θ , $\theta \in \Theta$. Suppose T is sufficient and complete, and that the distribution of $Y = Y(X)$ does not depend on θ . Then $\forall \theta \in \Theta$
 T and Y are independent under P_θ .

proof: Let A be a measurable set. $h(T) := P(Y \in A | T) - P(Y \in A)$
Since T is sufficient, and $Y = Y(X)$ only depends on X ,
 $P(Y \in A | T)$ does not depend on θ .

further note that since T is complete, we have that

$$E_\theta[h(T)] = E_\theta[P(Y \in A | T)] - P(Y \in A) = P(Y \in A) - P(Y \in A) = 0 \quad \forall \theta.$$

completeness

$$\Rightarrow h(T) = 0 \quad P_\theta\text{-a.s.} \quad \forall \theta. \quad \Leftrightarrow P(Y \in A | T) = P(Y \in A) \quad P_\theta\text{-a.s.}$$

Since A was chosen arbitrarily, it follows that $T \perp Y$. \square

example: X_1, \dots, X_n i.i.d. $N(\theta, \sigma^2)$, σ^2 known. Then $T = \bar{X}$ is sufficient and complete. Further the distribution of $\gamma := X - \bar{X}$ does not depend on θ .
Thus, by Basu's Lemma, \bar{X} and $X - \bar{X}$ are independent.
Hence, \bar{X} is UMRE.

Note: For an equivariant estimator T we have

$$T \text{ is UMRE} \Leftrightarrow E_\theta[T(X) | X - T(X)] = 0$$

Thus UMRE follows directly from independence, since $E_\theta[T(X)] = 0$

Decisions and Risk: Define \mathcal{A} to be the action space.

$\mathcal{A} = \mathbb{R} \rightarrow$ real-valued parameter, $\mathcal{A} = \{0, 1\} \rightarrow$ hypothesis testing,
 $\mathcal{A} = [0, 1] \rightarrow$ randomized test, $\mathcal{A} = \{\text{intervals}\} \rightarrow$ confidence intervals

A decision $d: \mathcal{X} \rightarrow \mathcal{A}$ is a mapping of observations to actions.
Depending on \mathcal{A} we call d a decision, estimator, test, ...

We define the loss function as $L: \Theta \times \mathcal{A} \rightarrow \mathbb{R}$, $L(\theta, a)$ the loss.

Lastly the Risk of a decision $d(X)$ is given by

$$R(\theta, d) := E_\theta[L(\theta, d(X))] \quad \forall \theta \in \Theta$$

example: (estimation). $g(\theta) \in \mathbb{R}$ parameter of interest. $\Theta = \mathbb{R}$.
Important loss functions are

$$L(\theta, a) := w(\theta) |g(\theta) - a|^r \quad \text{where } w(\cdot) \geq 0, r \geq 0.$$

the Risk is then

$$R(\theta, d) = w(\theta) \cdot E_{\theta} [|g(\theta) - d(X)|^r]$$

For $w=1, r=2$ we receive the quadratic loss and mean squared error.

example: (Tests) $H_0: \theta \in \Theta_0, H_1: \theta \in \Theta_1, \Theta = \{0, 1\}$, with loss

$$L(\theta, a) := \begin{cases} 1 & \text{if } \theta \in \Theta_0 \text{ and } a=1 \\ c & \text{if } \theta \in \Theta_1 \text{ and } a=0 \\ 0 & \text{else} \end{cases} \quad \text{for some } c > 0.$$

Then

$$R(\theta, d) = \begin{cases} P_{\theta}(d(X)=1) & \text{if } \theta \in \Theta_0 \\ c P_{\theta}(d(X)=0) & \text{if } \theta \in \Theta_1 \\ 0 & \text{otherwise} \end{cases} \quad \begin{array}{l} \rightarrow \text{type I error probability} \\ \rightarrow \text{type II error probability} \end{array}$$

Admissibility: Let d be a decision. We call the decision d' strictly better than d , if

$$R(\theta, d') \leq R(\theta, d) \quad \forall \theta \in \Theta$$

and additionally if

$$\exists \theta: R(\theta, d') < R(\theta, d).$$

We call d inadmissible, if there exists a d' strictly better than d .

(2) example: Consider the NP-test Φ_{NP} for $H_0: \theta = \theta_0, H_1: \theta = \theta_1$.
Then Φ_{NP} is admissible iff one of the two cases hold
(1) its power is strictly less than 1
OR (2) it has minimal level among all tests with power 1.

Suppose: $R(\theta_1, \Phi_{NP}) > 0$. Then Φ_{NP} is admissible.

proof: Assume $R(\theta_0, \phi) \leq R(\theta_0, \Phi_{NP})$.

$$\text{then } R(\theta_1, \phi) = \underbrace{R(\theta_1, \phi)}_{\geq c \cdot [R(\theta_0, \Phi_{NP}) - R(\theta_0, \phi)]} - R(\theta_1, \Phi_{NP}) + R(\theta_1, \Phi_{NP}) \geq R(\theta_1, \Phi_{NP})$$

may be proof
in script is
better
(Lemma 10.2.2)

□

Minimaxity: We call a decision d minimax, if

$$\sup_{\theta \in \Theta} R(\theta, d) = \inf_{\substack{d' \\ \text{decision}}} \left[\sup_{\theta \in \Theta} R(\theta, d') \right]$$

i.e. the best decision in the worst possible case.

example: A Neyman-Pearson test Φ_{NP} is minimax if and only if
 $R(\theta_0, \Phi_{NP}) = R(\theta_1, \Phi_{NP})$

don't yet understand →

proof:
" \Rightarrow " Let ϕ be a test. Assume $R(\theta_0, \Phi_{NP}) = R(\theta_1, \Phi_{NP})$, but Φ_{NP} is not minimax. Then for some test ϕ

$$\max_{i \in \{0, 1\}} R(\theta_i, \phi) < \max_{i \in \{0, 1\}} R(\theta_i, \Phi_{NP})$$

Then $R(\theta_0, \phi) < R(\theta_0, \Phi_{NP})$ and $R(\theta_1, \phi) < R(\theta_1, \Phi_{NP})$. But this contradicts the NP-Lemma.

" \Leftarrow " $S = \{(R(\theta_0, \phi), R(\theta_1, \phi)) | \phi\}$ is convex. Thus for $R(\theta_0, \Phi_{NP}) < R(\theta_1, \Phi_{NP})$ $\exists \phi$ with $R(\theta_0, \Phi_{NP}) < R(\theta_0, \phi) < R(\theta_1, \Phi_{NP})$ and $R(\theta_1, \phi) < R(\theta_1, \Phi_{NP})$. Then clearly Φ_{NP} is not minimax. Similarly for $R(\theta_0, \Phi_{NP}) > R(\theta_1, \Phi_{NP})$

Bayes risk: We now think of $\Theta \in \mathbb{H}$ as a random variable.

It is a probability measure on \mathbb{H} with density w ("prior density")
We define the Bayes risk to be $= w(\theta) d\mu(\theta)$

$$r(\Pi, d) := \int_{\mathbb{H}} R(\vartheta, d) d\Pi(\vartheta) = \sum_{\vartheta} R(\vartheta, d) \cdot w(\vartheta)$$

where $R(\vartheta, d) = E_{\vartheta}[L(\vartheta, d(x))]$. (\rightarrow weighted risks)

We call a decision d Bayes (wrt. Π), if

$$r_w(d) = r(\Pi, d) = \inf_{d'} r(\Pi, d') . (\rightarrow \text{minimal Bayes risk})$$

We write $d = d_{\text{Bayes}}$.

Bayes test: Consider testing $H_0: \Theta = \Theta_0$ against $H_1: \Theta = \Theta_1$, $\mathcal{A} = [0, 1]$.

Let $L(\Theta_0, a) = a$, $L(\Theta_1, a) = 1-a$, $w(\Theta_0) = w_0$, $w(\Theta_1) = w_1 = 1-w_0$.

Then

$$r_w(\phi) = w_0 R(\Theta_0, \phi) + w_1 R(\Theta_1, \phi) \quad (\rightarrow \text{want to minimize})$$

The Bayes test is $\Phi_{\text{Bayes}}(x) = \begin{cases} 1 & \text{if } p_1/p_0 > w_0/w_1 \\ q & \text{if } p_1/p_0 = w_0/w_1 \\ 0 & \text{if } p_1/p_0 < w_0/w_1 \end{cases}$.

Proof: The Bayes risk is given by

$$\begin{aligned} r_w(\phi) &= w_0 \cdot E_{\Theta_0}[L(\Theta_0, \phi(x))] + w_1 \cdot E_{\Theta_1}[L(\Theta_1, \phi(x))] \\ &= w_0 \cdot E_{\phi}[q] + w_1 \cdot E_{\phi}[1-q] \\ &= w_0 \int \phi p_0 + w_1 \cdot (1 - \int \phi p_1) \\ &= \int \phi \cdot (w_0 p_0 - w_1 p_1) + w_1 \end{aligned}$$

We can choose ϕ to minimize the above term as follows

$$\phi = \begin{cases} 1 & w_0 p_0 - w_1 p_1 < 0 \rightarrow (*) \text{ is negative. Thus maximize } \phi \\ q & w_0 p_0 - w_1 p_1 = 0 \rightarrow \text{doesn't really matter (integral is zero anyways)} \\ 0 & w_0 p_0 - w_1 p_1 > 0 \rightarrow (*) \text{ is positive. Thus minimize value of } \phi \end{cases}$$

□
q arbitrary
(0,1)

We can also show that

$$2. r_w(\Phi_{\text{Bayes}}) = 1 - \int |w_1 p_1 - w_0 p_0|.$$

Construction of Bayes estimators: We now think of p_θ as the density of X given the value of Θ , i.e.

$$p_\theta(x) = p(x|\theta) \quad x \in \mathcal{X}$$

and marginal density $p(\cdot) = \int_{\mathbb{H}} p(\cdot|\vartheta) w(\vartheta) d\mu(\vartheta)$.

The a posteriori density of Θ is

$$w(\vartheta|x) = p(x|\vartheta) \cdot \frac{w(\vartheta)}{p(x)} \quad \vartheta \in \mathbb{H}, x \in \mathcal{X}$$

Now, given data $X=x$, consider Θ as a random variable with density $w(\vartheta|x)$. For

$$L(x,a) := E_{\Theta}[L(\Theta,a)|X=x] = \int_{\mathbb{H}} L(\vartheta,a) w(\vartheta|x) d\mu(\vartheta)$$

$d(x) := \arg \min_a L(x,a)$ is the Bayes decision d_{Bayes} .

construction of bayes estimator

proof: Let d' be any decision

$$r_w(d') = \int_{\Theta} R(\theta, d') w(\theta) d\mu(\theta)$$

$$= \int_{\Theta} \left[\int_X L(\theta, d'(x)) p(x|\theta) d\nu(x) \right] w(\theta) d\mu(\theta)$$

switch integrals
and vs. a
posteriori
density

$$= \int_X \left[\int_{\Theta} L(\theta, d'(x)) w(\theta|x) d\mu(\theta) \right] p(x) d\nu(x)$$

$$= \int_X L(X, d'(x)) p(x) d\nu(x)$$

$$\geq \int_X L(X, d(x)) p(x) d\nu(x) = r_w(d).$$

from above □

example: given $\theta \sim \text{Poisson}(\theta)$, $\theta \sim \text{Gamma}(k, \lambda)$. Thus

$$w(\theta) = \lambda^k \theta^{k-1} e^{-\lambda\theta} / \Gamma(k), E[\theta] = k/\lambda.$$

The a posteriori density is

$$w(\theta|x) = p(x|\theta) \frac{w(\theta)}{p(x)} = \dots \propto e^{-\theta(1+\lambda)} \theta^{k+x-1} \sim \text{Gamma}(k+x, 1+\lambda)$$

With quadratic loss, the Bayes estimator is $\frac{k+x}{1+\lambda} = E[\theta|x]$

The maximum a posteriori estimator is $\frac{k+x-1}{1+\lambda}$

Extended Bayes: We call a statistic T extended Bayes, if there exists a sequence of prior densities $\{w_m\}_{m=1}^{\infty}$,

such that

$$r_{wm}(T) - \inf T' r_{wm}(T') \rightarrow 0, \text{ as } m \rightarrow \infty.$$

(Minmax) Lemma 11.1.1: Suppose T is a statistic with risk $R(\theta, T) = R(T)$ independent of θ . Then

(1) T admissible $\Rightarrow T$ minimax

(2) T Bayes $\Rightarrow T$ minimax

(3) T extended Bayes $\Rightarrow T$ minimax

proof:

(1) T is admissible. Thus $\forall T'$ there is either a θ with $R(\theta, T') > R(T)$, or $R(\theta, T') \geq R(T) \forall \theta$. Hence $\sup_{\theta} (R(\theta, T')) \geq R(T)$. T takes minimal risk across all decisions \rightarrow it is minimax.

(2) This follows from (3), since Bayes implies extended bayes.

(3) Assume that for each prior w_m there exists a bayes decision: $r_{wm}(T_m) = \inf_{T'} r_{wm}(T')$ $m=1, \dots$

By the extended bayes property: $\forall \varepsilon > 0 \exists m$ large enough such that

$$R(T) = r_{wm}(T) \leq r_{wm}(T_m) + \varepsilon \leq r_{wm}(T') + \varepsilon \leq \sup_{\theta} R(\theta, T') + \varepsilon$$

This holds for any arbitrary T' . Further since ε can be chosen arbitrarily small it follows that T is minimax.
(since $R(T) \leq \sup_{\theta} R(\theta, T')$)

□

example: $X \sim \text{Binomial}(n, \theta)$, prior $\Theta \sim \text{Beta}(r, s)$. The bayes estimator for quadratic loss is

$$T = \frac{X+r}{n+r+s}$$

Its risk is

$$R(\theta, T) = E_\theta[(T-\theta)^2] = \text{var}_\theta(T) + \text{bias}_\theta^2(T) \\ = \frac{\dots \theta^2 + \dots \theta + \dots}{\dots}$$

To apply the lemma, $R(T)$ must not depend on θ . Thus set the coefficients of θ^2 and θ to be zero.

$$\Rightarrow r=s=\sqrt{n}/2$$

$$\Rightarrow T = \frac{X + \frac{\sqrt{n}}{2}}{n + \sqrt{n}}$$

is minimax. (since it is bayes).

Admissibility Lemmas

Now assume Θ is an open set. And $R(\theta, T) < \infty \forall T$.

Lemma 11.2.1: Suppose the statistic T is Bayes with prior density w . Then T is admissible if (1) or (2) hold:

(1) The statistic T is the unique Bayes decision

$$r_w(T) = r_w(T') \Rightarrow \forall \theta: T=T' P_\theta\text{-a.s.}$$

(2) $\forall T': R(\theta, T')$ is continuous in θ , and also, $\forall U \subseteq \Theta$ open the prior probability $\Pi(U) := \int_U w(\theta) d\Pi(\theta)$ of U is strictly positive.

proof: (1):

Assume T is inadmissible. Then for some T' : $R(\theta, T') \leq R(\theta, T) \forall \theta$. Then also $r_w(T') \leq r_w(T)$. But because T is bayes, T attains the minimum bayes risk: $r_w(T') = r_w(T)$.

By assumption, then $\forall \theta: T=T' P_\theta\text{-a.s.}$, and thus

$\forall \theta: R(\theta, T') = R(\theta, T)$. But this is a contradiction to the second part of the admissibility definition.

(2) Suppose T is inadmissible. Then for some $T': R(\theta, T') < R(\theta, T) \forall \theta$.

And for some $\theta_0: R(\theta_0, T') < R(\theta_0, T)$. This implies that for some $\varepsilon > 0$ and an open neighborhood $U \subseteq \Theta$ of θ_0 we have

$$R(\theta, T') \leq R(\theta, T) - \varepsilon, \quad \theta \in U$$

But then

$$r_w(T') = \int_U R(\theta, T') w(\theta) d\Pi(\theta) + \int_{U^c} \underbrace{R(\theta, T')}_{\leq R(\theta, T)} w(\theta) d\Pi(\theta)$$

$$\leq \int_U R(\theta, T) w(\theta) d\Pi(\theta) - \varepsilon \Pi(U) + \int_{U^c} R(\theta, T) w(\theta) d\Pi(\theta)$$

$$= r_w(T) - \varepsilon \Pi(U) < r_w(T) \quad \text{contradiction}$$

□

Lemma 11.2.2: Suppose that T is extended Bayes, and that $\forall T': R(\theta, T')$ is continuous in θ . Assume $\forall U \subseteq \Theta$ open

$$\frac{r_{w_m}(T) - \inf_{T'} r_{w_m}(T')}{\Pi_m(U)} \rightarrow 0 \quad \text{as } m \rightarrow \infty$$

Here $\Pi_m(U) := \int_U w_m(\theta) d\Pi_m(\theta)$ is the probability of U under the prior Π_m .

Then T is admissible.

+ (proof)

↓
proof: Suppose T is inadmissible. Then $\exists T': R(\theta, T') \leq R(\theta, T)$
and $\exists \theta_0: R(\theta_0, T') < R(\theta_0, T)$. Then $\exists \varepsilon > 0 \exists U \subseteq \Theta$ open s.t.
 $R(\theta, T') \leq R(\theta, T) - \varepsilon, \forall \theta \in U$

Then $r_{W_m}(T') \leq r_{W_m}(T) - \varepsilon \Pi_m(U)$

Further suppose that a Bayes decision T_m for W_m exists K_m ,
i.e. $r_{W_m}(T_m) = \inf_{T'} r_{W_m}(T')$, $m=1,2,\dots$

Then K_m :

$$r_{W_m}(T_m) \leq r_{W_m}(T') \leq r_{W_m}(T) - \varepsilon \Pi_m(U)$$

$$\Rightarrow \frac{r_{W_m}(T) - r_{W_m}(T_m)}{\Pi_m(U)} \geq \varepsilon > 0 \quad \rightarrow \text{contradiction} \quad \square$$

(In-)admissible estimators for the normal mean: $X \sim N(\theta, 1)$, $\theta \in \Theta = \mathbb{R}$.

Let $R(\theta, T) = E_\theta[(T-\theta)^2]$ be the quadratic risk.

Consider estimators of the form $T = aX + b$, $a > 0$, $b \in \mathbb{R}$.

Then T is admissible iff either (1) or (2) holds

(1) $a < 1$

(2) $a = 1$ and $b = 0$.

proof:

" \Leftarrow " (1): We apply Lemma 11.2.1. First we show T is Bayes
for some prior. Then we show cond. (1) of the Lemma,
that T is unique. This implies that T is admissible.
First, choose the prior $\theta \sim W(c, c^2)$. From a previous
example we know that

$$T_{\text{Bayes}} = E[\theta | X] = \frac{c^2 X + c}{c^2 + 1}$$

Thus choose $a = \frac{c^2}{c^2 + 1}$, $b = \frac{c}{c^2 + 1}$. Then $T = T_{\text{Bayes}}$.

Now show uniqueness. From a past example we also know

$$r_W(T') = E[\text{var}[\theta | X]] + E[(T-T')^2]$$

Thus if $r_W(T') = r_W(T)$, then $E[(T-T')^2] = 0$.

We can write $X = \theta + \varepsilon$ with $\varepsilon \sim N(0, 1)$ indep. of θ .

Then $X \sim W(c, c^2 + 1)$ corresponds to $p(x) = \int p_{\theta|x}(x) w(\theta) d\theta$.
 $E[(T-T')^2] = 0 \Rightarrow T = T'$ P-a.s.

Since P dominates all P_θ , we have $T = T'$ P-a.s..

Thus T is unique and admissible.

" \Leftarrow " (2): In this case we have $T = X$. We use Lemma 11.2.2.

Because $R(\theta, T) = 1 \forall \theta$, also $r_W(T) = 1$ for any prior w .

Let w_m be the density of $W(0, m)$. From previous examples
we know $T_{\text{Bayes}} = T_m = \frac{m}{m+1} X$

By bias-variance decomposition, it has risk

$$R(\theta, T_m) = \frac{m^2}{(m+1)^2} + \left(\frac{m}{m+1} - 1 \right)^2 \theta^2 = \frac{m^2}{(m+1)^2} + \frac{\theta^2}{(m+1)^2}$$

As $E[\theta^2] = 0$: $r_{W_m}(T_m) = m/(m+1)$

Thus $r_{W_m}(T) - r_{W_m}(T_m) = \frac{1}{m+1}$

Then T is extended Bayes. It is left to show the
requirement of Lemma 11.2.2

↓(proof)

Instead of an open neighborhood, it suffices to consider open intervals $U = (u, u + h)$ with $u, h > 0$ fixed. We have

$$T_m(U) = \Phi\left(\frac{u+h}{\sqrt{m}}\right) - \Phi\left(\frac{u}{\sqrt{m}}\right) = \frac{1}{\sqrt{m}} \Phi'\left(\frac{u}{\sqrt{m}}\right) + o\left(\frac{1}{\sqrt{m}}\right)$$

$$\text{For } m \text{ large } \Phi\left(\frac{u}{\sqrt{m}}\right) \approx \Phi(0) = \frac{1}{\sqrt{2\pi}} > \frac{1}{4} \text{ (say)}$$

Thus, for m sufficiently large

$$\frac{T_m(U) - T_m(U_m)}{T_m(U)} \leq \frac{4}{h\sqrt{m}} \rightarrow 0 \quad \text{as } m \rightarrow \infty$$

This proves admissibility.

"=>" (1+2): We show that if (1) or (2) doesn't hold, then T is inadmissible.

case 1: $a > 1$. We have $R(\theta, aX+b) \geq \text{var}(aX+b) \geq 1 = R(\theta, X)$

So $aX+b$ is inadmissible

case 2: $a=1, b \neq 0$. The bias term makes $aX+b$ inadmissible:

$$R(\theta, aX+b) = 1 + b^2 > 1 = R(\theta, X)$$

□

Least squares estimator (LSE): We now regard our observations now as a matrix $X \in \mathbb{R}^{n \times p}$ (with non-random entries). Given each row, we look for the best linear approximation of y_i (row label). If X has rank p , the LSE is given by:

$$\hat{\beta} := \arg \min_{b \in \mathbb{R}^p} \|y - Xb\|_2^2 \quad \begin{aligned} &\rightarrow X \text{ has a column of 1's} \\ &\begin{pmatrix} 1 & x_1 & \dots & x_n \end{pmatrix} = X \end{aligned}$$

Then $X\hat{\beta}$ corresponds exactly to the projection of y onto the column space of X .

Distribution of LSE: For $f = E[y]$ define $\beta^* := (X^T X)^{-1} X^T f$. We call $X\beta^*$ the best linear approximation of f .

Lemma 12.3.1: Suppose $E[\varepsilon \varepsilon^T] = \sigma^2 I$ ($= \text{Var}(\varepsilon)$) where $\varepsilon := y - f = y - E[y]$. Then

$$(1) E[\hat{\beta}] = \beta^*, \quad \text{Cov}(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$$

$$(2) E[\|X(\hat{\beta} - \beta^*)\|_2^2] = \sigma^2 p$$

$$(3) E[\|X\hat{\beta} - f\|_2^2] = \underbrace{\sigma^2 p}_{\text{estimation error}} + \underbrace{\|X\beta^* - f\|_2^2}_{\text{misspecification error}}$$

proof:

$$(1) \text{ compute directly: } E[\hat{\beta} - \beta^*] = E[(X^T X)^{-1} X^T (y - f)] = A \underbrace{E[\varepsilon]}_{=0} = 0$$

$$\begin{aligned} \text{and } \text{Cov}(\hat{\beta}) &= \text{Cov}(\hat{\beta} - \beta^*) = \text{Cov}(A\varepsilon) = A \underbrace{\text{Cov}(\varepsilon)}_{\sigma^2 I} A^T \\ &= \sigma^2 A A^T = \sigma^2 (X^T X)^{-1} \end{aligned}$$

↓(2+3)

b(proof)

(2): Define the projection $PPT := X(X^T X)^{-1} X^T$. Then we can write

$$\|X(\hat{\beta} - \beta)\|_2^2 = \|PPT\epsilon\|_2^2 = \sum_{j=1}^p V_j$$

where $V = PTE$ (since: $\|PPTE\|_2^2 = E^T P P T E = (PTE)^T PTE = \|PTE\|_2^2$)

Note that $E[V] = PTE[\epsilon] = 0$ and $Cov(V) = PTE[Cov(\epsilon)]P = \sigma^2 I$

Then

$$E\left[\sum_{j=1}^p V_j^2\right] = \sum_{j=1}^p E[V_j^2]$$

diagonal elements of Cov(V)

(3): We have for any b

$$\|Xb - f\|_2^2 = \|Xb - X\beta^* + X\beta^* - f\|_2^2 = \|X(b - \beta^*)\|_2^2 + \|X\beta^* - f\|_2^2$$

because $X\beta^* - f$ is orthogonal to X. Now choose $b = \hat{\beta}$, apply (2), and take the Expectation on both sides.

The statement follows. \square

Asymptotic theory

We now regard the i.i.d. $X_1, \dots, X_n \in \mathbb{X}$ as the first n random variables of an infinite sequence. They have distribution P.

Write $P = P_x P_x \dots$ to be the distribution of the sequence $\{X_i\}_{i=1}^{\infty}$. We still have that $P = P_\theta$, $\theta = g(\theta) \in \mathbb{R}^P$ and $\Gamma = \{g(\theta) \mid \theta \in \Theta\}$.

$T_n = T_n(X_1, \dots, X_n)$ is an estimator.

Almost sure convergence: The random variables $Z_n \in \mathbb{R}^P$ converge almost surely to Z, if

$$P(\{w \in \Omega \mid \lim_{n \rightarrow \infty} Z_n(w) = Z(w)\}) = 1$$

we write $Z_n \xrightarrow{\text{a.s.}} Z$.

Convergence in probability: Let $\{Z_n\}_{n=1}^{\infty}$ and Z be \mathbb{R}^P -valued r.v.

We say that Z_n converges in probability to Z, if

$$\forall \varepsilon > 0: \lim_{n \rightarrow \infty} P(\|Z_n - Z\| > \varepsilon) = 0$$

we write $Z_n \xrightarrow{P} Z$.

Note: Chebychev's inequality can be used to prove convergence in probability. For all $\psi: [0, \infty) \rightarrow [0, \infty)$ increasing functions

$$P(\|Z_n - Z\| \geq \varepsilon) \leq \frac{E[\psi(\|Z_n - Z\|)]}{\psi(\varepsilon)}.$$

Convergence in distribution: Let $\{Z_n\}_{n=1}^{\infty}$ and Z be \mathbb{R}^P -valued r.v.

We say that Z_n converges in distribution to Z if

$$\forall f \text{ continuous and bounded} \quad \lim_{n \rightarrow \infty} E[f(Z_n)] = E[f(Z)]$$

we write $Z_n \xrightarrow{D} Z$.

(\Rightarrow implies bounded in probability)

Convergence implications:

- conv. in probability \Rightarrow conv. in distribution
- conv. in distribution to const. $c \in \mathbb{R}$ \Rightarrow conv. in probability
- almost sure conv. \Rightarrow conv. in probability
- conv. in probability \Rightarrow \exists subsequence with almost sure conv.

Central Limit theorem: (CLT) Let X_1, \dots, X_n be i.i.d. in \mathbb{R} .
with $E[X_i] = \mu$, $\text{Var}[X_i] = \sigma^2$. Then

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{D} N(0, \sigma^2) \quad (\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i)$$

cdf version:

$$P\left(\sqrt{n} \frac{(\bar{X}_n - \mu)}{\sigma} \leq z\right) \rightarrow \Phi(z) \quad \forall z$$

Cramér-Wold device: Let (Z_n) be r.v. then

$$Z_n \xrightarrow{D} Z \iff a^T Z_n \xrightarrow{D} a^T Z \quad \forall a \in \mathbb{R}^p$$

Stochastic order symbols: Let $\{Z_n\}$ be \mathbb{R}^p -valued random variables.

And $\{r_n\}$ strictly positive random variables. We write

$$Z_n = O_p(1) \quad \text{if} \quad \lim_{M \rightarrow \infty} \limsup_{n \rightarrow \infty} P(\|Z_n\| > M) = 0$$

We say that Z_n is bounded in probability. (Uniform tightness of $\{Z_n\}$).
Further we write $Z_n = O_p(r_n)$ if $\frac{Z_n}{r_n} = O_p(1)$.

Also if $Z_n \xrightarrow{P} 0$ then write $Z_n = o_p(1)$

Z_n is of small order r_n in probability, if $\frac{Z_n}{r_n} = o_p(1)$.
We write $Z_n = o_p(r_n)$.

Note: $Z_n = o_p(1) \Rightarrow Z_n = O_p(1)$

$$Z_n = O_p(r_n) \Rightarrow Z_n = O_p(r_n).$$

Note: Z_n converges in distribution $\Rightarrow Z_n = O_p(1)$

Slutsky's Theorem: Let $(\{Z_n, A_n\}, Z)$ be a collection of \mathbb{R}^p -valued random variables. $a \in \mathbb{R}^p$ a vector of constants.
If $Z_n \xrightarrow{D} Z$ and $A_n \xrightarrow{P} a$, then

$$A_n^T Z_n \xrightarrow{D} a^T Z$$

proof:

Take f Lipschitz and bounded.

$$|f| \leq C_B, |f(z) - f(\bar{z})| \leq C_L \|z - \bar{z}\|$$

$$\text{Write } |E[f(A_n^T Z_n)] - E[f(a^T Z)]|$$

$$\leq \underbrace{|E[f(A_n^T Z_n)] - E[f(a^T Z_n)]|}_{(*)} + \underbrace{|E[f(a^T Z_n)] - E[f(a^T Z)]|}_{\rightarrow 0}$$

→ also by
Cramér
World device

Z_n is bounded and A_n is as well.
Thus we can define f on a compact set. Thus, since f is continuous, it is automatically Lipschitz

Because f is bounded and Lipschitz, the second term goes to 0.
Now show that $(*)$ goes to 0. Let $\varepsilon > 0$, $M > 0$ be arbitrary.

Define $S_n := \{ \|Z_n\| \leq M, \|A_n - a\| \leq \varepsilon \}$. Then

$$|E[f(A_n^T Z_n)] - E[f(a^T Z_n)]| \leq E[|f(A_n^T Z_n) - f(a^T Z_n)|]$$

$$= E[|f(A_n^T Z_n) - f(a^T Z_n)|] \cdot \mathbb{P}(S_n) + E[|f(A_n^T Z_n) - f(a^T Z_n)|] \mathbb{P}(\bar{S}_n)$$

$$\leq C_L \varepsilon M + 2C_B \mathbb{P}(\bar{S}_n) \rightarrow 0$$

And because $\mathbb{P}(S_n) \leq \mathbb{P}(\|Z_n\| > M) + \mathbb{P}(\|A_n - a\| > \varepsilon) \rightarrow 0$ (for M large enough)
(for ε arbitrarily small). ◻

consistency: We call estimators $\{T_n\}$ of $\gamma = g(\theta)$ consistent, if

$$T_n \xrightarrow{\text{P}_0} \gamma \quad (\text{for } \theta \text{ fixed}).$$

Asymptotic normality: We call estimators $\{T_n\}$ of $\gamma = g(\theta)$ asymptotically normal with asymptotic covariance matrix V_θ , if

$$\sqrt{n}(T_n - \gamma) \xrightarrow{\text{D}_0} N(0, V_\theta) \quad \text{and if } \{T_n\} \text{ is consistent.}$$

Asymptotic linearity: We call estimators $\{T_n\}$ of $\gamma = g(\theta) \in \mathbb{R}^p$ asymptotically linear, if for a $l_\theta: \mathcal{X} \rightarrow \mathbb{R}^p$ with

$$E_\theta[l_\theta(x)] = 0 \quad \text{and} \quad E_\theta[l_\theta(x) l_\theta^T(x)] =: V_\theta < \infty$$

it holds that

$$T_n - \gamma = \frac{1}{n} \sum_{i=1}^n l_\theta(x_i) + \sigma_{P_\theta} \left(\frac{1}{\sqrt{n}} \right) \quad \begin{matrix} \rightarrow \text{corresponds to a} \\ \text{taylor expansion.} \end{matrix}$$

We call l_θ the influence function of $\{T_n\}$. It approximately measures the influence of an additional observation x .

Remark: asymptotically linear \Rightarrow asymptotically normal.

$$\text{because } \sqrt{n}(T_n - \gamma) = \frac{1}{\sqrt{n}} \sum_{i=1}^n l_\theta(x_i) + \sigma_{P_\theta} \xrightarrow{\text{D}_0} N(0, V_\theta).$$

example: Assume entries of X have finite variance, $\bar{T}_n := \bar{X}_n$ is a linear, and thus asymptotically linear, estimator of mean M with influence function $l_\theta(x) = x - M$. (no θ -term in $T_n - \gamma$)

example: mean M , $\gamma = \mu^2$, $T_n = \bar{X}_n^2$. Then we have

$$T_n - \gamma = \bar{X}_n^2 - \mu^2 = (\bar{X}_n - M)(\bar{X}_n + M) = (\bar{X}_n - M) \cdot 2M - \underbrace{(\bar{X}_n - M)^2}_{\sigma_{P_\theta}^2 \left(\frac{1}{n} \right)} = \sigma_{P_\theta} \left(\frac{1}{n} \right)$$

Thus $l_\theta(x) = 2M(x - M)$
 $\Rightarrow V_\theta = 4\mu^2 \sigma^2$.

Same steps for $\gamma = \sigma^2$ and $T_n = \hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$.

S-technique: Suppose we have estimators $\{T_n\}$ of $\gamma = g(\theta) \in \mathbb{R}^p$, that fulfill for a $l_\theta: \mathcal{X} \rightarrow \mathbb{R}^p$

$$E_\theta[l_\theta(x)] = 0 \quad \text{and} \quad \text{Cov}(l_\theta(x)) = V_\theta < \infty$$

$$\text{with} \quad T_n - g(\theta) = \frac{1}{n} \sum_{i=1}^n l_\theta(x_i) + \sigma_{P_\theta} \left(\frac{1}{\sqrt{n}} \right)$$

Let $h: \mathbb{R}^p \rightarrow \mathbb{R}$ be differentiable at $\gamma = g(\theta)$. Write $\dot{h}(\gamma) \in \mathbb{R}^p$.
Then $(h(T_n))$ is also asymptotically linear

$$h(T_n) - h(g(\theta)) = \frac{1}{n} \sum_{i=1}^n \dot{h}(\gamma)^T l_\theta(x_i) + \sigma_{P_\theta} \left(\frac{1}{\sqrt{n}} \right)$$

Note: $\sqrt{n}(h(T_n) - h(\gamma)) \xrightarrow{\text{D}_0} N(0, \dot{h}(\gamma)^T V_\theta \dot{h}(\gamma))$

example: Bernoulli(θ), $X \in \{0, 1\}$, $P_\theta(X=1) = 1 - P_\theta(X=0) = \theta$. $X_1, \dots, X_n \sim X$ i.i.d.
 $T = \bar{X}_n$, $\bar{X}_n - \theta = \frac{1}{n} \sum_{i=1}^n (X_i - \theta)$, $\sqrt{n}(\bar{X}_n - \theta) \xrightarrow{\text{D}_0} N(0, \theta(1-\theta))$.

$$\text{Let } h(\theta) = \log(\theta/(1-\theta)), \quad \dot{h}(\theta) = 1/\theta(1-\theta)$$

Thus by the S-technique we have

$$\sqrt{n} \left(\log \left(\frac{\bar{X}_n}{1-\bar{X}_n} \right) - \log \left(\frac{\theta}{1-\theta} \right) \right) \xrightarrow{\text{D}_0} N(0, \frac{1}{\theta(1-\theta)})$$

M-estimators: We define a new general class of estimators. ($\mu = g(\theta)$)
 But first we need to define some quantities. For each $\theta \in \Gamma$ (parameter of interest) we define a loss function $p_\theta(x) := L(\theta, d(x))$, for some decision d .

Further the theoretical risk

$$R(c) := E_\theta [p_c(x)] \quad \text{such that} \\ \mu = \arg \min_{c \in \Gamma} R(c) = \arg \min_{c \in \Gamma} E_\theta [p_c(x)]. \quad (\rightarrow \text{given } p_c, \text{ this also defines } \mu)$$

We can write $R(c) = 0$, if $c \mapsto p_c(x)$ is diff.able. Then write
 $\psi_c(x) := p'_c(x) = \frac{\partial}{\partial c} p_c(x)$ and thus $\dot{R}(c) = E_\theta [\psi_c(x)]$.

The empirical risk is defined as

$$\hat{R}_n(c) := \frac{1}{n} \sum_{i=1}^n p_c(x_i), \quad c \in \Gamma.$$

Finally, an M-estimator $\hat{\mu}_n$ of μ is defined as

$$\hat{\mu}_n := \arg \min_{c \in \Gamma} \frac{1}{n} \sum_{i=1}^n p_c(x_i) = \arg \min_{c \in \Gamma} \hat{R}_n(c).$$

If $\forall x \ p_c(x)$ is differentiable in c , then we define the Z-estimator $\hat{\mu}_n$ of μ to be defined by

$$\dot{\hat{R}}_n(\hat{\mu}_n) = 0 \quad \text{where } \dot{\hat{R}}_n(c) = \frac{1}{n} \sum_{i=1}^n \psi_c(x_i).$$

example: $X \in \mathbb{R}$, $\Gamma = \mathbb{R}$, $\mu = \mu = E_\theta[X]$, $p_c(x) = (x-c)^2$

$$\text{then } \mu = \arg \min_c E_\theta [(x-c)^2]$$

by bias variance decomposition, we have $E_\theta [(x-c)^2] = \sigma^2 + (\mu - c)^2$

Then clearly the term $\frac{1}{n} \sum_{i=1}^n (x_i - c)^2$ is minimized at \bar{x}_n .

This is the M-estimator.

example: $X \in \mathbb{R}$, $\Gamma = \mathbb{R}$, $\mu = \mu = E_\theta[X]$, $p_c(x) = |x-c| \rightarrow$ not diff.able at c .

First show that $R(c) = \int_{x \leq c} F(x) dx + \int_{x > c} (1-F(x)) dx$
 for $F(x) = P(X=x)$ the density. $= - \int_{x > c} (x-c) d(1-F(x))$

$$\Rightarrow R(c) = E[p_c(x)] = \int_{x \leq c} (c-x) dF(x) + \int_{x > c} (x-c) dF(x)$$

$$= \underbrace{[(c-x) F(x)]_{-\infty}^c}_{=0} - \underbrace{[(x-c)(1-F(x))]_c^\infty}_{=0} + \int_{x \leq c} F(x) dx + \int_{x > c} (1-F(x)) dx$$

$$\text{Then } \dot{R}(c) = F(c) - (1-F(c)) = 2F(c) - 1 \stackrel{!}{=} 0 \Rightarrow \mu = F^{-1}\left(\frac{1}{2}\right)$$

Assuming F is strictly increasing and continuous at $F^{-1}\left(\frac{1}{2}\right)$
 We can write the same for the empirical risk

$$\hat{R}_n(c) = \int_{x \leq c} \hat{F}_n(x) dx + \int_{x > c} (1-\hat{F}_n(x)) dx \Rightarrow \hat{\mu}_n = \hat{F}_n^{-1}\left(\frac{1}{2}\right)$$

Thus $\hat{\mu}_n = \begin{cases} X_{\lceil \frac{n+1}{2} \rceil} & n \text{ odd} \\ \frac{X_{\lfloor \frac{n}{2} \rfloor} + X_{\lceil \frac{n}{2} \rceil}}{2} & n \text{ even} \end{cases}$

MLE as special case of M-estimation: $X \sim P_\theta$, $\Theta \subseteq \mathbb{R}^p$
 suppose $\{x \mid p_\theta(x) > 0\}$ does not depend on θ (to not divide by 0 later)
 For the MLE we have:

$$\hat{\theta}_n = \arg \max_{\theta \in \Theta} \sum_{i=1}^n \log(p_\theta(x_i)) = \arg \min_{\theta \in \Theta} \hat{R}_n(\theta)$$

→ then $p_\theta(x)$ must be $p_\theta(x) = -\log p_{\tilde{\theta}}(x)$
 Define $R(\theta) := E_{\tilde{\theta}}[p_\theta(x)] = -E_{\tilde{\theta}}[\log(p_\theta(x))]$ (for $\tilde{\theta}$ the true value of the parameter)

We call $K(\theta | \tilde{\theta}) = R(\theta) - R(\tilde{\theta}) = E_{\tilde{\theta}}\left[\frac{p_\theta(x)}{p_{\tilde{\theta}}(x)}\right]$

the Kullback-Leibler information.

Lemma 14.1.1: The risk $R(\theta) = -E_{\tilde{\theta}}[\log(p_\theta(x))]$ is minimized at the true parameter $\tilde{\theta}$.

i.e. $K(\theta | \tilde{\theta}) = R(\theta) - R(\tilde{\theta}) \geq 0$

Proof: $R(\theta) - R(\tilde{\theta}) = -E_{\tilde{\theta}}\left[\log\left(\frac{p_\theta(x)}{p_{\tilde{\theta}}(x)}\right)\right] \stackrel{\text{jensen}}{\geq} -\log\left(E_{\tilde{\theta}}\left[\frac{p_\theta(x)}{p_{\tilde{\theta}}(x)}\right]\right)$
 $= -\log\left[\int \frac{p_\theta(x)}{p_{\tilde{\theta}}(x)} p_{\tilde{\theta}}\right] = -\log(1) = 0 \quad \square$

Further, if $p_\theta(x)$ is differentiable, then $\psi_\theta(x) = -\frac{p'_\theta(x)}{p_\theta(x)} = -S_\theta(x)$
 and $\dot{R}(\theta) = -E_{\tilde{\theta}}[S_\theta(x)] = 0$

Theorem 14.2.1: (Consistency of M-estimators)

Suppose the uniform convergence

$$\sup_{C \in \Gamma} |\hat{R}_n(C) - R(C)| \xrightarrow{\text{IP}} 0$$

\hat{R} : empirical risk
 R : theoretical risk

Then

$$R(\hat{\theta}_n) \xrightarrow{\text{IP}} R(\theta)$$

proof: $0 \leq R(\hat{\theta}_n) - R(\theta) = [R_n(\hat{\theta}_n) - R(\hat{\theta}_n)] + [\hat{R}_n(\theta) - R(\theta)] + [\hat{R}_n(\hat{\theta}_n) - \hat{R}_n(\theta)]$
 $\leq 2 \cdot \sup_C |\hat{R}_n(C) - R(C)| \xrightarrow{\text{IP}} 0 \quad (\text{by assumption}) \quad \square$

Theorem 14.2.2: $X \sim P_\theta$, $\Pi \subseteq \mathbb{R}^p$, $p_\theta: \mathcal{X} \rightarrow \mathbb{R}$, $c \in \Pi$, $\psi_c := \frac{2}{3c} p_\theta$
 remember then $R(c) = P p_c = E[p_\theta(x)]$, $\hat{R}_n(c) = \hat{P}_n p_c = \frac{1}{n} \sum_{i=1}^n p_\theta(x_i)$

$$\dot{R}(c) = P \psi_c, \quad \hat{R}_n(c) = \hat{P}_n \psi_c$$

(Note that here $P f := E_\theta[f(x)]$ and $\hat{P}_n f = \frac{1}{n} \sum_{i=1}^n f(x_i)$)
 Assume the Z-estimator, i.e. $\hat{R}_n(\hat{\theta}_n) = 0$ for $R(\theta) = 0$.

Theorem: Assume (i) $\Pi \subseteq \mathbb{R}$, (ii) $c \mapsto \psi_c(x)$ is continuous $\forall x$,

(iii) $E[|\psi_c(x)|] < \infty \forall c$

(iv) $\exists \delta \text{ s.t. } \dot{R}(c) > 0 \quad \forall c \in (x, x+\delta) \text{ and } \dot{R}(c) < 0 \quad \forall c \in (x-\delta, x)$
 and $\dot{R}(x) = 0$.

Then $\mathbb{P}(\exists \hat{\theta}_n : \hat{R}_n(\hat{\theta}_n) = 0) \rightarrow 1 \quad (\rightarrow \text{with high probability there is a solution.})$

and $\|\hat{\theta}_n - x\| = o_{\text{IP}}(1)$. (\rightarrow the solution $\hat{\theta}_n$ is in the neighborhood of interest) /29

↓ proof: Let ε be arbitrary with $0 < \varepsilon < \tilde{\varepsilon}$. From (iv) it follows that $\dot{R}(\gamma + \varepsilon) \geq 2\tilde{\varepsilon}$ and $\dot{R}(\gamma - \varepsilon) \leq -2\tilde{\varepsilon}$ (*)

Let $A = \{\dot{R}_n(\gamma + \varepsilon) - \dot{R}(\gamma + \varepsilon) > -\tilde{\varepsilon}, \text{ and } \dot{R}_n(\gamma - \varepsilon) - \dot{R}(\gamma - \varepsilon) \leq \tilde{\varepsilon}\}$
then $P(A^c) \rightarrow 0$ i.e. for n large enough (*) also holds for \dot{R}_n .

So $\dot{R}_n(\gamma + \varepsilon) = \underbrace{\dot{R}_n(\gamma + \varepsilon) - \dot{R}(\gamma + \varepsilon)}_{\geq -\tilde{\varepsilon}} + \underbrace{\dot{R}(\gamma + \varepsilon)}_{\geq 2\tilde{\varepsilon}} \geq \tilde{\varepsilon}.$ } on A

and $\dot{R}_n(\gamma - \varepsilon) = \underbrace{\dot{R}_n(\gamma - \varepsilon) - \dot{R}(\gamma - \varepsilon)}_{\leq \tilde{\varepsilon}} + \underbrace{\dot{R}(\gamma - \varepsilon)}_{\leq -2\tilde{\varepsilon}} \leq -\tilde{\varepsilon}$ } because \dot{R}_n is continuous

Thus by (ii) $\exists \hat{\gamma}_n \in (\gamma - \varepsilon, \gamma + \varepsilon)$ such that $\dot{R}_n(\hat{\gamma}_n) = 0$
And since ε was chosen arbitrarily it follows that

$$\|\hat{\gamma}_n - \gamma\| \xrightarrow{P} 0$$

□

Influence function of Z-estimator: Suppose the following

(i) $\hat{\gamma}_n \xrightarrow{P} \gamma$
(ii) $\dot{R}_n(\hat{\gamma}_n) = 0, \dot{R}(\gamma) = 0$

(iii) $V_n(c) := \sqrt{n}(\hat{P}_n - P)\Psi_0 = \sqrt{n}(\dot{R}_n(c) - \dot{R}(c)) \quad c \in \mathbb{R}$ ("empirical process")
is asymptotically continuous at γ , i.e.

if sequences $\{\gamma_n\} \in \Gamma$ with $\|\gamma_n - \gamma\| = o_{P_0}(1)$ it holds $|V_n(\gamma_n) - V_n(\gamma)| = o_{P_0}(1)$

(iv) $M_\theta = \frac{\partial}{\partial c} \dot{R}(c) > 0$ and particularly exists.

(v) $J_\theta := P\Psi_n\Psi_n^T < \infty$ (covariance matrix)

Then $\hat{\gamma}_n$ is asymptotically linear with influence function

$$L_\theta = -M_\theta^{-1}\Psi_\theta$$

Note: $\sqrt{n}(\hat{\gamma}_n - \gamma) \xrightarrow{D} N(0, \underbrace{M_\theta^{-1} J_\theta M_\theta^{-1}}_{\text{"sandwich formula"}}$)

"sandwich formula"

Asymptotic relative efficiency: $\gamma \in \Gamma \subseteq \mathbb{R}$. Let $T_{n,1}, T_{n,2}$ be estimators of γ s.t.

$$\sqrt{n}(T_{n,j} - \gamma) \xrightarrow{D} N(0, V_{\theta,j}) \quad \text{for } j=1,2.$$

Then we call $e_{2:1} := \frac{V_{\theta,1}}{V_{\theta,2}}$

the asymptotic relative efficiency of $T_{n,2}$ with respect to $T_{n,1}$.

If $e_{2:1} > 1$, then estimator $T_{n,2}$ is asymptotically more efficient than $T_{n,1}$. Since an asymptotic $(1-\alpha)$ -confidence interval for γ based on $T_{n,2}$ is then narrower than one based on $T_{n,1}$.

↓ asymptotic relative efficiency

example: $X = \mu + \varepsilon$, $\varepsilon \sim F_0$ symmetric around 0, $F_0 = F_0'$, $F_0 > 0$
 $\text{var}(\varepsilon) := \sigma^2 < \infty$.
 X_1, \dots, X_n i.i.d. $\sim X$. $T_{n,1} = \bar{X}_n$, $T_{n,2}$ = sample median.
 $V_{\theta,1} = \sigma^2$, $V_{\theta,2} = 1/(4\sigma^2)$. Thus

$$e_{2:1} = 4\sigma^2 F_0^2(\theta)$$

case: for F_0 the std. normal distribution: $e_{2:1} = \frac{2}{\pi} \approx 0.64 < 1$.
→ $T_{n,1}$ is asymptotically more efficient.

Asymptotic pivots (based on MLE): $\Theta \subseteq \mathbb{R}^p$ and remember the MLE

$$\hat{\theta}_n = \arg \max_{\theta \in \Theta} \sum_{i=1}^n \log(P_\theta(x_i))$$

Further remember the loss function $p_{\theta,2} = -\log(p_\theta)$, $q_{\theta,2} = \dot{p}_\theta/\theta = -S_\theta(x)$
where $S_{\theta,2} := p_{\theta,2}/p_{\theta,1}$.

And the asymptotic variance of the MLE is $I^{-1}(\theta)$ where
 $I(\theta) = P_\theta S_\theta S_\theta^\top$ is the Fisher information:

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{D} N(0, I^{-1}(\theta)) \quad \forall \theta.$$

Lemma 14.8.1: Define

$$2\Delta_n(\hat{\theta}_n) - 2\Delta_n(\theta) := 2 \sum_{i=1}^n [\log(P_{\hat{\theta}_n}(x_i)) - \log(P_\theta(x_i))]$$

Under regularity conditions, the term above is an asymptotic pivot for θ , i.e. ($\gamma = g(\theta)$), its asymptotic distribution does not depend on θ .

Especially, its asymptotic distribution is χ^2 -distrib. with p degrees of freedom:

$$2\Delta_n(\hat{\theta}_n) - 2\Delta_n(\theta) \xrightarrow{D} \chi^2_p \quad \forall \theta.$$

Likelihood ratio tests and asymptotics: Consider the hypothesis

$H_0: R(\theta) = 0$ where $R(\theta) = \begin{pmatrix} R_1(\theta) \\ \vdots \\ R_q(\theta) \end{pmatrix}$ restrictions on $\theta \in \mathbb{R}^p$

We define the unrestricted MLE

$$\hat{\theta}_n = \arg \max_{\theta \in \Theta} \sum_{i=1}^n \log(P_\theta(x_i))$$

and the restricted MLE:

$$\hat{\theta}_n^0 = \arg \max_{\substack{\theta \in \Theta \\ R(\theta)=0}} \sum_{i=1}^n \log(P_\theta(x_i))$$

Assume $\dot{R}(\theta) = \frac{\partial}{\partial \theta} R(\theta) \mid_{\theta=\theta_0}$ has rank q .

$$\Delta_n(\hat{\theta}_n) - \Delta_n(\hat{\theta}_n^0) = \sum_{i=1}^n [\log(P_{\hat{\theta}_n}(x_i)) - \log(P_{\hat{\theta}_n^0}(x_i))]$$

is the log-likelihood ratio for testing $H_0: R(\theta) = 0$

Lemma 14.10.1: Under regularity conditions and if $H_0: R(\theta) = 0$ holds, we have

$$2\Delta_n(\hat{\theta}_n) - 2\Delta_n(\hat{\theta}_n^0) \xrightarrow{D} \chi^2_q.$$

