

Sampling Distributions

Sofia Olhede



October 4, 2020

1 Sampling Theory

2 Limiting Distributions

3 Estimation

Sampling Distributions III

- We recall that t_{n-1} denotes the t distribution on $n - 1$ degrees of freedom. The t distribution on k degrees of freedom takes the form

$$f_X(x, k) = \frac{\Gamma(\frac{k+1}{2})}{\Gamma(k/2)\sqrt{k\pi}} \left(1 + \frac{x^2}{k}\right)^{-\frac{k+1}{2}}, \quad x \in \mathbb{R}.$$

- Let us also assume $k > 2$ then the mean and variance of $X \sim t_k$ are

$$\mathbb{E}(X) = 0, \quad \mathbb{V}\text{ar}(X) = \frac{k}{k-2}.$$

- Theorem: let $Y_1 \sim \chi_{d_1}^2$ and let $Y_2 \sim \chi_{d_2}^2$ be independent. Then

$$\frac{Y_1/d_1}{Y_2/d_2} \sim F_{d_1, d_2}.$$

Sampling Distributions IV

- A random variable follows the Fisher F distribution with integer parameters d_1 and d_2 , written as $X \sim F_{d_1, d_2}$ if

$$f_X(x; d_1, d_2) = \frac{1}{B(d_1/2, d_2/2)} \left(\frac{d_1}{d_2}\right)^{d_1/2} x^{d_1/2-1} \left(1 + \frac{d_1}{d_2} x\right)^{-\frac{d_1+d_2}{2}}, \quad x \geq 0.$$

The mean and variance in this case are for $d_2 > 4$

$$\mathbb{E}(X) = \frac{d_2}{d_2 - 2}, \quad \text{Var}(X) = \frac{2d_2^2(d_1 + d_2 - 2)}{d_1(d_2 - 4)(d_2 - 2)^2}.$$

Sampling Distributions V

- Theorem: (Sampling from an Exponential Family).

Let $Y_1, \dots, Y_n \stackrel{\text{iid}}{\sim} f$ where

$$f(y) = \exp \left\{ \sum_{j=1}^k \phi_j T_j(y) - \gamma(\phi_1, \dots, \phi_k) + S(y) \right\}, \quad y \in \mathcal{Y}, \quad \phi \in \Phi.$$

This is a density of a k parameter exponential family form. If Φ is open then

- (1) The minimal sufficient statistic for ϕ is τ where

$$\tau_j = \sum_{i=1}^n T_j(y_i).$$

- (2) The function γ is infinitely differentiable in all k of its variables and

$$\mathbb{E}(\tau) = n \nabla_{\phi} \gamma(\phi) \quad \text{and} \quad \text{Cov}(\tau) = n \nabla_{\phi}^2 \gamma(\phi).$$

Sampling Distributions VI

- Unfortunately the sampling distribution of T is not always available in closed form.
- It may become easier to work with an approximation valid for large n .
- These approximations will be understood as a form of convergence of F .
- Definition: Convergence in Distribution (Weak Convergence). Let $\{F_n\}_{n \geq 1}$ be a sequence of distribution functions and let G be a distribution function on \mathbb{R} . We say that F_n converges weakly or in distribution to G and write $F_n \xrightarrow{\mathcal{L}} G$ whenever

$$F_n(y) \xrightarrow{n \rightarrow \infty} G(y),$$

for all y constituting continuity points of G .

- A stronger notion of convergence corresponds to convergence in probability.

Sampling Distributions VI

- Definition (convergence in probability): When a sequence of random variables satisfies $\Pr\{\|Y_n - Y\| > \epsilon\} \rightarrow 0$ for all $\epsilon > 0$ and a given (random variable) Y , then we say that Y_n converges in probability to Y , and write $Y_n \xrightarrow{P} Y$.
- $\xrightarrow{\mathcal{L}}$ relates distribution functions. It says that the probabilistic behaviour of a sequence Y_n becomes more and more alike that of the limit Y .
- \xrightarrow{P} relates random variables. It says that the actual realisations of Y_n can be progressively approximated with high probability by those of Y .
- Theorem: (a) $Y_n \xrightarrow{P} Y \Rightarrow Y_n \xrightarrow{\mathcal{L}} Y$.
(b) $Y_n \xrightarrow{\mathcal{L}} c \Rightarrow Y_n \xrightarrow{P} c$ for $c \in \mathbb{R}$.
- Theorem (The Continuous Mapping Theorem)
Let $g : \mathbb{R} \rightarrow \mathbb{R}$ be continuous on the range of Y . Then
(a) $Y_n \xrightarrow{P} Y \Rightarrow g(Y_n) \xrightarrow{P} g(Y)$,
(b) $Y_n \xrightarrow{\mathcal{L}} Y \Rightarrow g(Y_n) \xrightarrow{\mathcal{L}} g(Y)$.

Limiting Distributions I

- Theorem (Slutsky's theorem): Let $X_n \xrightarrow{\mathcal{L}} X$ and let $Y_n \xrightarrow{\mathcal{L}} c$ where $c \in \mathbb{R}$. Then
 - (a) $X_n + Y_n \xrightarrow{\mathcal{L}} X + c$.
 - (b) $X_n Y_n \xrightarrow{\mathcal{L}} Xc$.
- Theorem (General version of Slutsky's theorem): Let $g : \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ be continuous and assume that $X_n \xrightarrow{\mathcal{L}} X$ and let $Y_n \xrightarrow{\mathcal{L}} c$ where $c \in \mathbb{R}$. Then $g(X_n, Y_n) \xrightarrow{\mathcal{L}} g(X, c)$ with increasing n .
- By applying the continuous mapping theorem and by applying Slutsky's theorem we can obtain new approximations. But they require a starting point.
- Theorem (Law of Large Numbers): let Y_n be independent random variables with $\mathbb{E} Y_k = \mu$ and $\mathbb{E} |Y_k| < \infty$ for all k . Then $n^{-1}(Y_1 + \dots + Y_n) \xrightarrow{P} \mu$.

Limiting Distributions II

- Theorem (Central Limit Theorem). Let $\{Y_n\}$ be a sequence of iid random variables with mean μ and variance σ^2 which is assumed finite. Then

$$\sqrt{n}\left(\frac{1}{n} \sum_{i=1}^n (Y_i - \mu)\right) \xrightarrow{\mathcal{L}} N(0, \sigma^2).$$

- Theorem (Delta method): Let $Z_n = a_n(X_n - \theta) \xrightarrow{\mathcal{L}} Z$ where $a_n \in \mathbb{R}^+$ and $\theta \in \mathbb{R}$ for all n and assume $a_n \rightarrow \infty$. Let $g(\cdot)$ be continuously differentiable at θ . Then

$$a_n\{g(X_n) - g(\theta)\} \xrightarrow{\mathcal{L}} g'(\theta)Z.$$

- Proof: Taylor expansion around θ gives

$$g(X_n) = g(\theta) + g'(\theta^*)\{X_n - \theta\}, \quad \|\theta - \theta^*\| < \|\theta - X_n\|.$$

(This is the Lagrange form of the remainder).

Limiting Distributions III

- Thus we may deduce

$$\|\theta - \theta^*\| < \|\theta - X_n\| = a_n^{-1}|a_n(\theta - X_n)| = a_n^{-1}|Z_n| \xrightarrow{P} 0. \quad (1)$$

The latter equation uses Slutsky's theorem. Thus $\theta^* \xrightarrow{P} 0$. By the continuous mapping theorem $g'(\theta^*) \xrightarrow{P} g'(\theta)$.

Therefore it follows

$$a_n\{g(X_n) - g(\theta)\} = a_n\{g(\theta) + g'(\theta^*)(X_n - \theta) - g(\theta)\} \quad (2)$$

$$= g'(\theta^*)a_n(X_n - \theta) \xrightarrow{\mathcal{L}} g'(\theta)Z. \quad (3)$$



Limiting Distributions IV

- We can apply these methods to derive sampling distributions.
- Corollary: Let $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} f$ where

$$f(x) = \exp\{\phi T(x) - \gamma(\phi) + S(x)\}, \quad x \in \mathcal{X},$$

with $\phi \in \Phi \subset \mathbb{R}$ and

$$\bar{T}_n = \frac{1}{n} \sum_{i=1}^n T(X_i) = n^{-1} \tau(X_1, \dots, X_n).$$

If Φ is open then γ is infinitely differentiable and so

$$\sqrt{n}(\bar{T}_n - \gamma'(\phi)) \xrightarrow{\mathcal{L}} N(0, \gamma''(\phi)).$$

Limiting Distributions V

- Theorem (Weighted sum CLT): Let $\{W_n\}$ be an iid sequence of real random variables, with common mean 0 and variance 1. Let $\{\gamma_n\}$ be a sequence of real constants. Then

$$\sup_{1 \leq j \leq n} \frac{\gamma_j^2}{\sum_{i=1}^n \gamma_i^2} \rightarrow 0 \implies \frac{1}{\sqrt{\sum_{i=1}^n \gamma_i^2}} \sum_{j=1}^n \gamma_j W_j \xrightarrow{\mathcal{L}} N(0, 1).$$

- For joint convergence we need to consider random vectors.

Limiting Distributions VI

- Definitions. Let $\{Y_n\}$ be a sequence of random vectors of \mathbb{R}^d and Y a random vector of \mathbb{R}^d with $Y_n = \begin{pmatrix} Y_n^{(1)} & \dots & Y_n^{(d)} \end{pmatrix}^T$ and $Y = \begin{pmatrix} Y^{(1)} & \dots & Y^{(d)} \end{pmatrix}^T$. Also define the CDFs F_{Y_n} and F_Y . We say that Y_n converges in distribution to Y as $n \rightarrow \infty$ (and write $Y_n \xrightarrow{\mathcal{L}} Y$) if for every continuity point of F_Y we have

$$F_{Y_n}(y) \xrightarrow{n \rightarrow \infty} F_Y(y).$$

- Theorem (Cramér–Wold Device). Let $\{Y_n\}$ be a sequence of random variables of \mathbb{R}^d and let Y be a random vector of \mathbb{R}^d . Then

$$Y_n \xrightarrow{\mathcal{L}} Y \Leftrightarrow u^T Y_n \xrightarrow{\mathcal{L}} u^T Y, \quad \forall u \in \mathbb{R}^d.$$

Limiting Distributions VII

- Continuous mapping theorem and Slutsky's lemma generalize to the vector case.
- In either case
 - (a) Continuous mapping: if $g : \mathbb{R}^p \rightarrow \mathbb{R}^d$ is continuous on the range of U and if $U_n \xrightarrow{\mathcal{L}} U$ in \mathbb{R}^p then $g(U_n) \xrightarrow{\mathcal{L}} g(U)$ in \mathbb{R}^d .
 - (b) Slutsky: if $g : \mathbb{R}^p \times \mathbb{R}^q \rightarrow \mathbb{R}^d$ is continuous and $U_n \xrightarrow{\mathcal{L}} U$ in \mathbb{R}^p as well as $W_n \xrightarrow{\mathcal{L}} u \in \mathbb{R}^q$ for deterministic u then $g(U_n, W_n) \xrightarrow{\mathcal{L}} g(U, u)$.

Limiting Distributions VIII

- Convergence in probability also easily generalizes to the vector case:
- Definition (Convergence in Probability (vectors)): Given a sequence of random vectors $\{Y_n\}$ in \mathbb{R}^d satisfies $\Pr\{\|Y_n - Y\| > \epsilon\} \xrightarrow{n \rightarrow \infty} 0$ for any $\epsilon > 0$ and a given random vector Y we say that Y_n converges in probability to Y and write $Y_n \xrightarrow{P} Y$.
- Theorem (Multivariate law of large numbers). Let $\{Y_n\}$ be iid random vectors with expectation μ and if $\mathbb{E}\|Y_n\| < \infty$, for all $k \geq 1$,

$$\frac{1}{n} \sum_{k=1}^n Y_k \xrightarrow{P} \mu.$$

- Theorem (Multivariate CLT). Let $\{X_n\}$ be a sequence of random vectors in \mathbb{R}^d with mean μ and covariance Ω and define \bar{X}_n as the mean of the vectors X_1, \dots, X_n . Then

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{\mathcal{L}} Z \sim \mathcal{N}_d(0, \Omega).$$

Limiting Distributions VIII

- Theorem: (Delta Method, vector case). Set $Z_n \equiv a_n(X_n - u) \xrightarrow{\mathcal{L}} Z$ in \mathbb{R}^d where $a_n \in \mathbb{R}$, $u \in \mathbb{R}^d$ and $a_n \rightarrow \infty$. Let $g : \mathbb{R}^d \rightarrow \mathbb{R}^p$ be continuously differentiable at u . Then

$$a_n\{g(X_n) - g(u)\} \xrightarrow{\mathcal{L}} J_g(u)Z,$$

where $J_g(y)$ is the $p \times d$ Jacobian matrix of g ,

$$J_g(y) = \begin{pmatrix} \frac{\partial}{\partial x_1} g_1(y) & \cdots & \frac{\partial}{\partial x_d} g_1(y) \\ \vdots & \ddots & \vdots \\ \frac{\partial}{\partial x_1} g_p(y) & \cdots & \frac{\partial}{\partial x_d} g_p(y) \end{pmatrix}. \quad (4)$$

Estimation I

- What is estimation (\equiv “learning” in machine learning)?
- Imagine you assume Y is distributed according to $F(y_1, \dots, y_n; \theta)$ where $y \in \mathcal{Y}^n$.
- Assume you know the form of $F(y_1, \dots, y_n; \theta)$ but not the value of θ .
- Guessing θ on having observed y_1, \dots, y_n is estimation.
- Point estimation corresponds to producing a decent estimator of θ given we have observed y_1, \dots, y_n .
- Or more formally we define a point estimator as a statistic with codomain Θ or $T : \mathcal{Y}^n \mapsto \Theta$.
- We usually write this as $\hat{\theta}(Y_1, \dots, Y_n)$. θ is deterministic and $\hat{\theta}$ is stochastic.

Estimation II

- Not that whenever we realise a different set of Y_1, \dots, Y_n then we realise a different $\hat{\theta}(Y_1, \dots, Y_n)$.
- How do we design an estimator $\hat{\theta}(Y_1, \dots, Y_n)$?
- A good estimator would normally produce a value of $\hat{\theta}(Y_1, \dots, Y_n)$ near θ .
- We usually address this in terms of the mean and variance of $\hat{\theta}(Y_1, \dots, Y_n)$.
- The first interpretation of this is that “on average” we get the right value from $\hat{\theta}(Y_1, \dots, Y_n)$, and the spread of values obtained is not significant.
- This brings us round to the notion of the mean square error.

Estimation III

- Definition (mean square error): assume that $\hat{\theta}$ is an estimator of the parameter θ corresponding to the model $F(y; \theta)$, where $\theta \in \Theta \subset \mathbb{R}^d$. The mean square error of $\hat{\theta}$ is then defined as

$$\text{MSE}\{\hat{\theta}, \theta\} = \mathbb{E}\left[\|\hat{\theta} - \theta\|^2\right] \quad (5)$$

- Lemma: (Mean Square Error Decomposition) The mean square error admits the decomposition

$$\mathbb{E}\left[\|\hat{\theta} - \theta\|^2\right] = \|\mathbb{E}\hat{\theta} - \theta\|^2 + \mathbb{E}\|\hat{\theta} - \mathbb{E}\hat{\theta}\|^2.$$

This decomposition can colloquially be described as the mean square error of an estimator, is the bias square plus its variance.

Estimation IV

- The proof of this result is straightforward and just corresponds to

$$\begin{aligned}\mathbb{E}\left[\|\hat{\theta} - \theta\|^2\right] &= \mathbb{E}\left[\|\hat{\theta} - \mathbb{E}\hat{\theta} + \mathbb{E}\hat{\theta} - \theta\|^2\right] \\ &= \mathbb{E}\left[\|\hat{\theta} - \mathbb{E}\hat{\theta}\|^2\right] + 0 + 0 + \left[\|\mathbb{E}\hat{\theta} - \theta\|^2\right].\end{aligned}\quad (6)$$

- Why is the MSE important? Basically it transpires that the concentration of $\hat{\theta}$ around θ can always be phrased in terms of the MSE.
- Lemma: Let $\hat{\theta}$ be an estimator of $\theta \in \mathbb{R}^p$. For any $\epsilon > 0$

$$\Pr\left\{\|\hat{\theta} - \theta\| > \epsilon\right\} \leq \frac{\text{MSE}\{\hat{\theta}, \theta\}}{\epsilon^2}.$$

From this relationship we can note that $\text{MSE}\{\hat{\theta}, \theta\} \rightarrow 0 \Rightarrow \hat{\theta} \xrightarrow{P} \theta$.

Measures of Performance

- When an estimator has this property, then it is called consistent for θ .
- Definition: An estimator $\hat{\theta}_n$ is consistent for parameter θ in terms of sample size n if $\hat{\theta}_n \xrightarrow{P} \theta$.
- So once the MSE vanishes, consistency follows, but the converse is more tricky.
- Can we then always design a consistent estimator? Not necessarily!
- Definition: (Identifiability). A given model $F(y; \theta)$ is identifiable if given any $\theta_1 \in \Theta$ and $\theta_2 \in \Theta$ if it holds that

$$\theta_1 \neq \theta_2 \Rightarrow F(y; \theta_1) \neq F(y; \theta_2).$$

- If the model is not identifiable then you can get the same model with different values of the parameters.

Measures of Performance

- As an example consider $N(\mu_1 + \mu_2, \sigma^2)$. From observations we can only estimate $\mu_1 + \mu_2$.
- We will always assume identifiability unless we explicitly specify otherwise.
- We can use $\text{MSE}\{\hat{\theta}, \theta\}$ to determine which estimators are “better”.
- But what is the best MSE given a problem? This is a difficult problem.
- A easier problem might be, among unbiased estimators (bias zero), can we make the MSE arbitrarily small?

Measures of Performance

- Assume that Y_1, \dots, Y_n have a joint density or probability mass function $f(\mathbf{y}; \theta)$ depending on an unknown $\theta \in \mathbb{R}$.
- Let us try to find a function $\Lambda(\theta)$ such that

$$\text{Var}\{\hat{\theta}\} \geq \Lambda(\theta) \forall \theta \in \Theta,$$

with the choice of any $\hat{\theta}$ such that $\mathbb{E}\{\hat{\theta}\} = \theta$.

- To be able to study this question, we shall start with studying expectations with respect to y .
- We start from (with some regularity conditions thrown in):

$$\begin{aligned} \mathcal{V}(\theta) &\equiv \mathbb{E}\{S(Y)\} \\ &= \int_{y^n} S(y) f(y; \theta) dy \\ \frac{\partial}{\partial \theta} \mathcal{V}(\theta) &= \frac{\partial}{\partial \theta} \int_{y^n} S(y) f(y; \theta) dy = \int_{y^n} S(y) \frac{\partial}{\partial \theta} f(y; \theta) dy. \end{aligned}$$

Measures of Performance

- Now with some further manipulation we note that

$$\begin{aligned}\frac{\partial}{\partial \theta} \mathcal{Y}(\theta) &\equiv \frac{\partial}{\partial \theta} \mathbb{E}\{S(Y)\} = \int_{\mathcal{Y}^n} S(y) \frac{\partial}{\partial \theta} f(y; \theta) dy \\ &= \int_{\mathcal{Y}^n} S(y) \frac{f(y; \theta)}{f(y; \theta)} \frac{\partial}{\partial \theta} f(y; \theta) dy \\ &= \int_{\mathcal{Y}^n} S(y) f(y; \theta) \frac{\partial}{\partial \theta} \log(f(y; \theta)) dy.\end{aligned}\tag{7}$$

- Now we can use this to derive further results, taking $U(y; \theta) = \frac{\partial}{\partial \theta} \log(f(y; \theta))$ and $S(y) = 1$, we get

$$\begin{aligned}\mathbb{E}\{U\} &= \int_{\mathcal{Y}^n} f(y; \theta) \frac{\partial}{\partial \theta} \log(f(y; \theta)) dy \\ &= \frac{\partial}{\partial \theta} \int_{\mathcal{Y}^n} f(y; \theta) dy = \frac{\partial}{\partial \theta} 1 = 0.\end{aligned}\tag{8}$$

Measures of Performance

- Given $\mathbb{E}\{U\} = 0$ it follows that $\text{Var}\{U\} = \mathbb{E}\{U^2\}$. We note

$$\text{Var}\{U\} = \mathbb{E}\left\{\left(\frac{\partial}{\partial\theta} \log(f(Y; \theta))\right)^2\right\}. \quad (9)$$

- If we instead take $S(Y) = \hat{\theta}(Y)$ then we arrive at

$$\begin{aligned} \text{Cov}\{\hat{\theta}(Y), U(Y)\} &= \mathbb{E}\{\hat{\theta}(Y)U(Y)\} - \mathbb{E}\{\hat{\theta}(Y)\} \mathbb{E}\{U(Y)\} \\ &= \mathbb{E}\{\hat{\theta}(Y)U(Y)\}. \end{aligned} \quad (10)$$

Furthermore, we note that

$$\begin{aligned} \mathbb{E}\{\hat{\theta}(Y)U(Y)\} &= \int_{\mathcal{Y}^n} \hat{\theta}(y)f(y; \theta) \frac{\partial}{\partial\theta} \log(f(y; \theta)) dy \\ &= \frac{\partial}{\partial\theta} \int_{\mathcal{Y}^n} \hat{\theta}(y)f(y; \theta) dy \stackrel{\text{unbiased}}{=} \frac{\partial}{\partial\theta} \theta = 1. \end{aligned} \quad (11)$$

- What can we learn from this?

Measures of Performance

- The Cauchy–Schwarz inequality states that if we have any two random variables A and B then it follows that

$$\mathbb{V}\text{ar}\{A\} \mathbb{V}\text{ar}\{B\} \geq \mathbb{C}\text{ov}^2\{A, B\}. \quad (12)$$

Taking $A = \hat{\theta}(Y)$ and $B = U(Y)$ we therefore arrive at:

$$\begin{aligned} \mathbb{V}\text{ar}\{\hat{\theta}(Y)\} \mathbb{V}\text{ar}\{U(Y)\} &\geq \mathbb{C}\text{ov}^2\{\hat{\theta}(Y), U(Y)\} \\ \Rightarrow \mathbb{V}\text{ar}\{\hat{\theta}(Y)\} &\geq \frac{1^2}{\mathbb{V}\text{ar}\{U(Y)\}} \\ &= \frac{1}{\mathbb{E}\left\{\left(\frac{\partial}{\partial \theta} \log(f(Y; \theta))\right)^2\right\}}. \end{aligned} \quad (13)$$

- This yields the Cramér–Rao lower bound on the variance of an unbiased estimator.

Measures of Performance

- We define the Fisher information to be

$$\mathcal{I}_n(\theta) = \mathbb{E} \left\{ \left(\frac{\partial}{\partial \theta} \log(f(Y; \theta)) \right)^2 \right\}.$$

- Then the Cramér–Rao lower bound on the variance of unbiased estimator $\hat{\theta}(y)$ states that

$$\mathbb{V}\text{ar} \left\{ \hat{\theta}(Y) \right\} \geq \frac{1}{\mathcal{I}_n(\theta)}.$$

- If $Y = (Y_1, \dots, Y_n)$ has n iid entries then $f(y; \theta) = f(y_1; \theta)f(y_2; \theta) \dots f(y_n; \theta)$ and so we find:

$$\mathcal{I}_n(\theta) = n \cdot \mathcal{I}_1(\theta).$$

- Unless the density is very peculiar further simplifications can be made.

Measures of Performance

- If we play around further using Fubini's theorem to exchange the order of integration and differentiation then we arrive at

$$\mathcal{I}_n(\theta) = \mathbb{E} \left\{ - \left(\frac{\partial^2}{\partial \theta^2} \log(f(Y; \theta)) \right) \right\}.$$

- The real reason why this hold, and why we have a negative sign, will become apparent later in the course.
- The next mathematical question is—can we achieve the Cramér–Rao bound?
- OK let us start backwards!
- Assume that

$$\text{var} \{ \hat{\theta} \} = \frac{1}{\mathcal{I}_n(\theta)}.$$

Measures of Performance

- Then the bound becomes an equality. We then have

$$\text{var}\{\hat{\theta}\} = \frac{\text{Cov}^2\left\{\hat{\theta}, \frac{\partial}{\partial\theta} \log(f(Y; \theta))\right\}}{\text{Var}\left\{\frac{\partial}{\partial\theta} \log(f(Y; \theta))\right\}}.$$

- This is true if and only if $\hat{\theta}$ is a linear function of $\frac{\partial}{\partial\theta} \log(f(Y; \theta))$ or

$$\hat{\theta} = a \frac{\partial}{\partial\theta} \log(f(Y; \theta)) + b,$$

for some constant a and b .

- Solving this equation yields that we can achieve this if and only if the density (frequency) of Y is a one-parameter exponential family with sufficient statistic $\hat{\theta}$.

Measures of Performance

- Sufficiency is key to this result.
- Theorem (Rao-Blackwell Theorem): Let $\hat{\theta}$ be an unbiased estimator of θ that has finite variance. Assume T is sufficient for θ . In this case $\hat{\theta}^* = \mathbb{E}\{\hat{\theta} | T\}$ is an unbiased estimator and

$$\text{Var}\{\hat{\theta}^*\} \leq \text{Var}\{\hat{\theta}\}.$$

- Equality is attained if and only if $\Pr\{\hat{\theta}^* = \hat{\theta}\} = 1$.
- Getting rid of irrelevant information improves estimation performance.
- The new estimator $\hat{\theta}^*$ is called a “Rao-Blackwellised” version of $\hat{\theta}$.