

6. Sufficient, Complete, and Ancillary Statistics

The Basic Statistical Model

Consider again the [basic statistical model](#), in which we have a [random experiment](#) with an observable [random variable](#) X taking values in a set S . Once again, the experiment is typically to sample n objects from a population and record one or more measurements for each item. In this case, the outcome variable has the form

$$\mathbf{X} = (X_1, X_2, \dots, X_n)$$

where X_i is the vector of measurements for the i^{th} item. In general, we suppose that the distribution of \mathbf{X} depends on a parameter θ taking values in a parameter space Θ . The parameter θ may also be vector-valued. We will use subscripts in probability density functions, expected values, etc. to denote the dependence on θ

Sufficient Statistics

Let $U = h(\mathbf{X})$ be a statistic taking values in a set T . Intuitively, U is sufficient for θ if U contains all of the information about θ that is available in the entire data variable \mathbf{X} . Formally, U is **sufficient** for θ if the [conditional distribution](#) of \mathbf{X} given U does not depend on θ .

Sufficiency is related to the concept of **data reduction**. Suppose that \mathbf{X} takes values in \mathbb{R}^n . If we can find a sufficient statistic U that takes values in \mathbb{R}^j , then we can reduce the original data vector \mathbf{X} (whose dimension n is usually large) to the vector of statistics U (whose dimension j is usually much smaller) with no loss of information about the parameter θ .

The following result gives a condition for sufficiency that is equivalent to this definition.

❖ 1. Let $U = h(\mathbf{X})$ be a statistic taking values in T , and let f_θ and g_θ denote the [probability density functions](#) of \mathbf{X} and U respectively. Show that U is sufficient for θ if and only if the function

$$\frac{f_\theta(\mathbf{x})}{g_\theta(h(\mathbf{x}))}, \quad \mathbf{x} \in S$$

is independent of θ . *Hint:* The joint distribution of (\mathbf{X}, U) is concentrated on the set $\{(\mathbf{x}, u) : (\mathbf{x} \in S) \text{ and } (u = h(\mathbf{x}))\} \subseteq S \times T$.

The Factorization Theorem

The definition precisely captures the intuitive notion of sufficiency given above, but can be difficult to apply.

We must know in advance a candidate statistic U , and then we must be able to compute the conditional distribution of X given U . The **factorization theorem** given in the next exercise frequently allows the identification of a sufficient statistic from the form of the probability density function of X .

2. Let f_θ denote the probability density function of X and suppose that $U = h(X)$ is a statistic taking values in T . Show that U is sufficient for θ if and only if there exists $G : T \times \Theta \rightarrow \mathbb{R}$ and $r : S \rightarrow \mathbb{R}$ such that

$$f_\theta(x) = G(h(x), \theta) r(x), \quad x \in S, \theta \in \Theta$$

Note that r depends only on the data x but not on the parameter θ .

3. Show that if U and V are equivalent statistics and U is sufficient for θ then V is sufficient for θ .

Special Distributions

We will determine sufficient statistics for several parametric families of distributions.

4. Suppose that $X = (X_1, X_2, \dots, X_n)$ is a random sample of size n from the **Bernoulli distribution** with success parameter $p \in [0, 1]$. Thus, $X_i = 1$ if trial i is a success, and $X_i = 0$ if trial i is a failure. Let $Y = \sum_{i=1}^n X_i$ denote the number of successes, and recall that Y has the **binomial distribution** with parameters n and p . Show directly from the definition that Y is sufficient for p . Specifically, show that the conditional distribution of X given $Y = k$ is the uniform distribution on the set of points

$$\{(x_1, x_2, \dots, x_n) \in \{0, 1\}^n : x_1 + x_2 + \dots + x_n = k\}$$

The result in the previous exercise is intuitively appealing: in a sequence of Bernoulli trials, all of the information about the probability of success p is contained in the number of successes Y . The particular *order* of the successes and failures provides no additional information. Of course, the sufficiency of Y follows more easily from the **factorization theorem**, but the conditional distribution provides additional insight.

5. Suppose that the distribution of X is a k -parameter **exponential family** with the natural statistic $U = h(X)$. Show that U is sufficient for θ . Because of this result, U is referred to as the **natural sufficient statistic** for the exponential family.

6. Suppose that $X = (X_1, X_2, \dots, X_n)$ is a random sample of size n from the **normal distribution** with mean $\mu \in \mathbb{R}$ and variance $\sigma^2 \in (0, \infty)$.

- Show that (Y, V) is sufficient for (μ, σ^2) where $Y = \sum_{i=1}^n X_i$ and $V = \sum_{i=1}^n X_i^2$.
- Show that (M, S^2) is sufficient for (μ, σ^2) where M is the sample mean of X and S^2 is the sample variance of X . *Hint:* Use part (a) and equivalence.

7. Suppose that $X = (X_1, X_2, \dots, X_n)$ is a random sample of size n from the **Poisson distribution** with

mean $a \in (0, \infty)$. Show that $Y = \sum_{i=1}^n X_i$ is sufficient for a .

8. Suppose that $\mathbf{X} = (X_1, X_2, \dots, X_n)$ is a random sample from the [gamma distribution](#) with shape parameter $k \in (0, \infty)$ and scale parameter $b \in (0, \infty)$
- Show that (Y, V) is sufficient for (k, b) where $Y = \sum_{i=1}^n X_i$ and $V = \prod_{i=1}^n X_i$.
 - Show that (M, U) is sufficient for (k, b) where M is the sample (arithmetic) mean of \mathbf{X} and U is the sample geometric mean of \mathbf{X} . *Hint:* Use part (a) and equivalence.
9. Suppose that $\mathbf{X} = (X_1, X_2, \dots, X_n)$ is a random sample from the [beta distribution](#) with left parameter $a \in (0, \infty)$ and right parameter $b \in (0, \infty)$. Show that (U, V) is sufficient for (a, b) where $U = \prod_{i=1}^n X_i$ and $V = \prod_{i=1}^n (1 - X_i)$.
10. Suppose that $\mathbf{X} = (X_1, X_2, \dots, X_n)$ is a random sample from the [Pareto distribution](#) with shape parameter $a \in (0, \infty)$. Show that $U = \prod_{i=1}^n X_i$ is sufficient for a .
11. Suppose that $\mathbf{X} = (X_1, X_2, \dots, X_n)$ is a random sample from the [uniform distribution](#) on the interval $[0, a]$ where $a \in (0, \infty)$ is the unknown parameter. Show that $X_{n,n} = \max \{X_1, X_2, \dots, X_n\}$ (the n^{th} [order statistic](#)) is sufficient for a .

Minimal Sufficient Statistics

The entire data variable \mathbf{X} is trivially sufficient for θ . However, as noted above, there usually exists a statistic U that is sufficient for θ and has smaller dimension, so that we can achieve real data reduction. Naturally, we would like to find the statistic U that has the smallest dimension possible. In many cases, this smallest dimension j will be the same as the dimension k of the parameter vector θ . However, as we will see, this is not necessarily the case; j can be smaller or larger than k .

Formally, suppose that a statistic U is sufficient for θ . Then U is **minimally sufficient** if U is a function of any other statistic V that is sufficient for θ . Once again, the definition precisely captures the notion of minimal sufficiency, but is hard to apply. The following exercise gives an equivalent condition.

12. Let f_θ denote the probability density function of \mathbf{X} corresponding to the parameter value θ and suppose that $U = h(\mathbf{X})$ is a statistic taking values in T . Show that U is minimally sufficient for θ if the following condition holds: for $\mathbf{x} \in S$ and $\mathbf{y} \in S$

$$\frac{f_\theta(\mathbf{x})}{f_\theta(\mathbf{y})} \text{ is independent of } \theta \text{ if and only if } h(\mathbf{x}) = h(\mathbf{y})$$

Hint: If $V = g(\mathbf{X})$ is another sufficient statistic, use the [factorization theorem](#) and the condition above to show that $g(\mathbf{x}) = g(\mathbf{y})$ implies $h(\mathbf{x}) = h(\mathbf{y})$ for $\mathbf{x} \in S$ and $\mathbf{y} \in S$. Then conclude that U is a function of V .

13. Show that if U and V are equivalent statistics and U is minimally sufficient for θ then V is minimally sufficient for θ .
14. Suppose that the distribution of \mathbf{X} is a k -parameter exponential family with natural sufficient statistic $U = h(\mathbf{X})$. Show that U is a minimally sufficient for θ .
15. Show that the sufficient statistics given above for the Bernoulli, Poisson, normal, gamma, and beta families are minimally sufficient for the given parameters.
16. Suppose that $\mathbf{X} = (X_1, X_2, \dots, X_n)$ is a random sample from the [uniform distribution](#) on the interval $[a, a + 1]$ where $a \in (0, \infty)$ is the unknown parameter. Show that $(X_{n,1}, X_{n,n})$, the vector consisting of the first and last [order statistics](#), is minimally sufficient for a . Note that we have a single parameter, but the minimally sufficient statistic is a vector of dimension 2.

Properties of Sufficient Statistics

Sufficiency is related to several of the methods of constructing estimators that we have studied.

17. Suppose that U is sufficient for θ and that there exists a [maximum likelihood estimator](#) of θ . Show that there exists a MLE V that is a function of U . *Hint:* Use the [factorization theorem](#).

In particular, suppose that V is the unique MLE of θ and that V is sufficient for θ . If U is sufficient for θ then V is a function of U by the previous exercise. Hence it follows that V is minimally sufficient for θ .

18. Suppose that the statistic U is sufficient for the parameter θ and that V is a [Bayes' estimator](#) of θ . Show that V is a function of U . *Hint:* Use the [factorization theorem](#).

The following exercise gives the [Rao-Blackwell theorem](#), named for [CR Rao](#) and [David Blackwell](#). The theorem shows how a sufficient statistic can be used to improve an unbiased estimator.

19. Suppose that U is sufficient for θ and that V is an unbiased estimator of a real parameter $\lambda = \lambda(\theta)$. Use sufficiency and properties of [conditional expectation](#) and conditional variance to show that
- $\mathbb{E}_\theta(V|U)$ is a valid statistic. That is, it does not depend on θ , in spite of the formal dependence on θ in the expected value.
 - $\mathbb{E}(V|U)$ is a function of U .
 - $\mathbb{E}(V|U)$ is an unbiased estimator of λ .
 - $\text{var}_\theta(\mathbb{E}(V|U)) \leq \text{var}_\theta(V)$ for any $\theta \in \Theta$ so $\mathbb{E}(V|U)$ is uniformly better than V .

Complete Statistics

Suppose that $U = h(\mathbf{X})$ is a statistic taking values in a set T . Then U is a [complete](#) statistic for θ if for any real-valued function g on T

$$\mathbb{E}_\theta(g(U)) = 0 \text{ for all } \theta \in \Theta \Rightarrow \mathbb{P}_\theta(g(U) = 0) = 0 \text{ for all } \theta \in \Theta$$

To understand this rather strange looking condition, suppose that $g(U)$ is a statistic constructed from U that is being used as an estimator of 0 (thought of as a function of θ). The completeness condition means that the only such unbiased estimator is the statistic that is 0 with probability 1.

20. Show that if U and V are equivalent statistics and U is complete for θ then V is complete for θ .

Special Distributions

21. Suppose that $X = (X_1, X_2, \dots, X_n)$ is a random sample of size n from the Bernoulli distribution with success parameter $p \in (0, 1)$. Show that the number of successes $Y = \sum_{i=1}^n X_i$ is complete for p . *Hint:* Note that $\mathbb{E}_p(g(Y))$ can be written as a polynomial in $t = \frac{p}{1-p}$. If this polynomial is 0 for all t in an open interval, then the coefficients must be 0.

22. Suppose that $X = (X_1, X_2, \dots, X_n)$ is a random sample of size n from the Poisson distribution with parameter $a \in (0, \infty)$. Show that the sum of the sample values $Y = \sum_{i=1}^n X_i$ is complete for a . *Hint:* Note that $\mathbb{E}_a(g(Y))$ can be written as a power series in a . If this series is 0 for all a in an open interval, then the coefficients must be 0.

23. Suppose that $X = (X_1, X_2, \dots, X_n)$ is a random sample of size n from the exponential distribution with scale parameter $b \in (0, \infty)$. Show that the sum of the sample values $Y = \sum_{i=1}^n X_i$ is complete for b . *Hint:* Note that $\mathbb{E}_b(g(Y))$ is the Laplace transform of a certain function. If this transform is 0 for all b in an open interval, then the function must be 0.

The results in the previous exercises generalize to exponential families, although the proof is complicated. Specifically, if the distribution of X is a k -parameter exponential family with the natural sufficient statistic $U = h(X)$ then U is complete for θ (as well as minimally sufficient for θ). This applies to random samples from the Bernoulli, Poisson, normal, gamma, and beta distributions discussed above.

The notion of completeness depends very much on the parameter space.

24. Suppose that $X = (X_1, X_2, X_3)$ is a random sample of size 3 from the Bernoulli distribution with success parameter $p \in \{\frac{1}{3}, \frac{1}{2}\}$. Show that $Y = X_1 + X_2 + X_3$ is not complete for p .

The Lehmann-Scheffé Theorem

The next exercise shows the importance of complete sufficient statistics; it is known as the **Lehmann-Scheffé theorem**, named for **Erich Lehmann** and **Henry Scheffé**.

25. Suppose that U is sufficient and complete for θ and that $T = r(U)$ is an unbiased estimator of a real parameter $\lambda = \lambda(\theta)$. Show that T is a uniformly minimum variance unbiased estimator of λ . The proof is

based on the following steps:

- Suppose that V is an unbiased estimator of λ . By the Rao-Blackwell theorem, $\mathbb{E}(V|U)$ is also an unbiased estimator of λ and is uniformly better than V .
- Since $\mathbb{E}(V|U)$ is a function of U , use completeness to conclude that $T = \mathbb{E}(V|U)$ with probability 1.

26. Suppose that $\mathbf{X} = (X_1, X_2, \dots, X_n)$ is a random sample of size n from the **Bernoulli distribution** with success parameter $p \in [0, 1]$. As usual, let $Y = \sum_{i=1}^n X_i$ denote the number of successes. Show that an UMVUE for $p(1-p)$, the variance of the sample distribution, is

$$\frac{Y}{n-1} \left(1 - \frac{Y}{n}\right)$$

27. Suppose that $\mathbf{X} = (X_1, X_2, \dots, X_n)$ is a random sample of size n from the **Poisson distribution** with parameter μ . Let $Y = \sum_{i=1}^n X_i$. Show that an UMVUE for $\mathbb{P}(X = 0) = e^{-\mu}$ is

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

Hint: Use the **probability generating function** of Y .

Ancillary Statistics

Suppose that $V = r(\mathbf{X})$ is a statistic taking values in a set T . If the distribution of V does not depend on θ , then V is called an **ancillary** statistic for θ . Thus, the notion of an ancillary statistic is complementary to the notion of a sufficient statistics (which contains all information about the parameter that is contained in the sample). Thus, the result in the following exercise, known as **Basu's Theorem** and named for Debabrata Basu, makes this point more precisely.

28. Suppose that U is complete and sufficient for a parameter θ and that V is an ancillary statistic. Show that U and V are independent. The following steps sketch the proof:

- Let g denote the probability density function of V and let $v \mapsto g(v|U)$ denote the conditional probability density function of V given U .
- Use properties of conditional expected value to show that $\mathbb{E}(g(v|U)) = g(v)$ for $v \in T$.
- Use completeness to conclude that $g(v|U) = g(v)$ with probability 1.

29. Show that if U and V are equivalent statistics and U is ancillary for θ then V is ancillary for θ .

30. Suppose that $\mathbf{X} = (X_1, X_2, \dots, X_n)$ is a random sample from a **scale family** with scale parameter $b \in (0, \infty)$. Show that V is an ancillary statistic for b if V is a function of

$$\left(\frac{X_1}{X_n}, \frac{X_2}{X_n}, \dots, \frac{X_{n-1}}{X_n}\right)$$

31. Suppose that $X = (X_1, X_2, \dots, X_n)$ is a random sample of size n from the [gamma distribution](#) with shape parameter $k \in (0, \infty)$ and scale parameter $b \in (0, \infty)$. Let M denote the sample arithmetic mean of X and let U denote the sample geometric mean of X . Show that $\frac{M}{U}$ is ancillary for b , and thus conclude that M and $\frac{M}{U}$ are independent. *Hint:* Use the previous exercise.

[Virtual Laboratories](#) > [7. Point Estimation](#) > [1](#) [2](#) [3](#) [4](#) [5](#) **[6](#)**

[Contents](#) | [Applets](#) | [Data Sets](#) | [Biographies](#) | [External Resources](#) | [Key words](#) | [Feedback](#) | ©