

## LECTURE 4

## 7. SUFFICIENT STATISTICS

Consider the “usual” statistical setup: the data is  $X$  and the parameter is  $\Theta$ .

To gain information about the parameter we study various functions of the data  $X$ . For instance, if  $X = (X_1, \dots, X_n)$  are IID  $\text{Ber}(\theta)$  given  $\Theta = \theta$ , then we would use  $T(X) = n^{-1}(X_1 + \dots + X_n)$  to get information about the parameter. A function of the data is called a *statistic*.

**Definition 6.** Let  $(\mathcal{T}, \mathcal{C})$  be a measurable space such that the  $\sigma$ -field  $\mathcal{C}$  contains all singletons. A measurable mapping  $T : \mathcal{X} \rightarrow \mathcal{T}$  is called a statistic.

As usual we think of a measurable space as a subspace of  $\mathbb{R}^d$  and the  $\sigma$ -field as the corresponding sub- $\sigma$ -field.

Although, formally a statistic is a mapping from the sample space  $\mathcal{X}$  to some space  $\mathcal{T}$ , we can also think of the composition  $T \circ X : S \rightarrow \mathcal{T}$  (recall that  $S$  is the underlying probability space). This is a random variable taking values in  $\mathcal{T}$  and we often write  $T$  for this random quantity.

In the next sections we will look more closely at different classes of statistics. That is, functions of the data with certain interesting properties. The first such class is the class of *sufficient statistics*.

**7.1. Sufficient statistics (classical).** The idea of sufficiency is to find a function  $T$  of the data  $X$  that summarizes the information about the parameter  $\Theta$ . Above we mentioned the example of IID  $\text{Ber}(\theta)$  random variables,  $X_1, \dots, X_n$ , where we know that we only need to know a function of the data, for instance  $X_1 + \dots + X_n$ , in order to compute an estimate of  $\theta$ . Similarly, we argued for the betting problem that decisions can be based entirely of knowing  $X_1 + \dots + X_n$  and not all the individual  $X_i$ ’s.

**Elementary case:** Let us first see what sufficient statistics is when we have densities. Suppose that the (conditional) distribution of  $X$  and  $T = T(X)$  given  $\Theta = \theta$  both have densities w.r.t. a measure  $\nu$  (think Lebesgue measure or counting measure). Then we say that  $T$  is a sufficient statistic for  $\Theta$  if  $f_{X|T,\Theta}(x | t, \theta)$  does not depend on  $\theta$ .

Note that, with  $t = T(x)$ ,

$$f_{X|T,\Theta}(x | t, \theta) = \frac{f_{X,T|\Theta}(x, t | \theta)}{f_{T|\Theta}(t | \theta)} = \frac{f_{X|\Theta}(x | \theta)}{f_{T|\Theta}(t | \theta)}.$$

Hence  $T$  is sufficient if this ratio does not depend on  $\theta$ .

To see how  $T$  captures the “information” about  $\Theta$  we can write down the likelihood function as

$$f_{X|\Theta}(x | \theta) = f_{X|T,\Theta}(x | T(x), \theta) f_{T|\Theta}(T(x) | \theta)$$

If  $T$  is sufficient, then the first factor on the RHS does not depend on  $\theta$  and the likelihood when observing  $X = x$  is proportional (as a function of  $\theta$ ) to the likelihood when observing  $T = T(x)$ . That is, information about  $\Theta$  comes only through the function  $T$ . If we, for example, want to maximize the likelihood we could maximize  $f_{T|\Theta}(t | \theta)$  instead of maximizing  $f_{X|\Theta}(x | \theta)$ . In this sense, there is no need to know  $x$  itself, it is sufficient to know  $t = T(x)$  to do inference.

**General case:** Formally, sufficient statistics are introduced as follows. First let  $\mu_{T|\Theta}(\cdot | \theta)$  be the conditional distribution of  $T$  given  $\Theta = \theta$ . It is a probability measure on  $\mathcal{C}$  given by  $\mu_{T|\Theta}(C | \theta) = \mu_{X|\Theta}(T^{-1}C | \theta)$ .

**Definition 7.** Suppose there exist versions of  $\mu_{X|\Theta,T}(\cdot | \theta, t)$  and a function  $r : \mathcal{B} \times \mathcal{T} \rightarrow [0, 1]$  such that

- (i)  $r(\cdot, t)$  is a probability on  $\mathcal{B}$  for each  $t \in \mathcal{T}$ ,
  - (ii)  $r(B, \cdot)$  is measurable  $\mathcal{T}$  for each  $B \in \mathcal{B}$ ,
- and for each  $\theta \in \Omega$  and  $B \in \mathcal{B}$

$$\mu_{X|\Theta,T}(B | \theta, t) = r(B, t), \quad \text{for } \mu_{T|\Theta}(\cdot | \theta) - \text{a.e. } t.$$

Then  $T$  is called a *sufficient statistic for  $\Theta$  (in the classical sense)*.

Note that the function  $r$  satisfies the conditions of a conditional distribution and does not depend on  $\theta$ . Hence,  $T$  is sufficient if  $\mu_{X|\Theta,T}(\cdot | \theta, t)$  is a conditional distribution that does not depend on  $\theta$ .

The simplest example (but not particularly useful) of a sufficient statistic is the data itself. That is  $\mathcal{T} = \mathcal{X}$  and  $T(x) = x$ . Of course,  $\mu_{X|\Theta,X}(B | \theta, x) = I_B(x)$  does not depend on  $\theta$  so the statistic is sufficient. Using this statistic does not help you to summarize information about the parameter as it is as complicated as the data itself. Let's look at some simple cases where there exist simple sufficient statistics.

**Example 10** (c.f. Example 6.2.3 in Casella & Berger). Let  $\{X_n\}$  be IID  $\text{Ber}(\theta)$  given  $\Theta = \theta$  and  $X = (X_1, \dots, X_n)$ . Put  $T(x) = x_1 + \dots + x_n$ . Let us show  $T$  is sufficient. Note that  $T(X)$  is  $\text{Bin}(n, \theta)$  given  $\Theta = \theta$ . For each  $x = (x_1, \dots, x_n)$   $x_i \in \{0, 1\}$  such that  $t = T(x)$

$$f_{X|\Theta,T}(x | \theta, t) = \frac{f_{X,T|\Theta}(x, t)}{f_{T|\Theta}(t | \theta)} = \frac{\theta^t(1-\theta)^{n-t}}{\binom{n}{t}\theta^t(1-\theta)^{n-t}} = \binom{n}{t}^{-1}.$$

Since this does not depend on  $\theta$ ,  $T$  is a sufficient statistic.

**Example 11.** Let  $\{X_n\}$  be IID  $\text{Exp}(\theta)$  given  $\Theta = \theta$  and  $X = (X_1, \dots, X_n)$ . Put  $T(x) = x_1 + \dots + x_n$ . Let us show  $T$  is sufficient. Note that  $T(X)$  is  $\Gamma(n, \theta)$  given  $\Theta = \theta$ . For each  $x_i > 0$  we have with  $t = T(x)$

$$\frac{f_{X|\Theta}(x | \theta)}{f_{T|\Theta}(t | \theta)} = \frac{\prod_{i=1}^n \theta e^{-\theta x_i}}{\frac{\theta^n}{\Gamma(n)} t^{n-1} e^{-\theta t}} = \frac{(n-1)!}{t^{n-1}},$$

which does not depend on  $\theta$ .

**7.2. Sufficient statistics (Bayesian).** In Bayesian statistics there is a slightly different notion of sufficiency, but it often coincides with the classical notion.

**Definition 8.** A statistic  $T$  is called a sufficient statistic for the parameter  $\Theta$  (in the Bayesian sense) if, for every prior  $\mu_\Theta$ , there exist versions of the posterior distributions  $\mu_{\Theta|X}$  and  $\mu_{\Theta|T}$  such that, for every  $A \in \tau$ ,  $\mu_{\Theta|X}(B | x) = \mu_{\Theta|T}(B | T(x))$ ,  $\mu_X$ -a.s., where  $\mu_X$  is the marginal distribution of  $X$ .

Hence, no matter what prior one uses, one only has to consider the sufficient statistic for making inference, because the posterior distribution given  $T = T(x)$  is the same as the posterior given the data  $X = x$ .

Let's see how this looks like with densities. If both  $\mu_{\Theta|X}(\cdot | x)$  and  $\mu_{\Theta|T}(\cdot | t)$  have densities w.r.t. the prior  $\mu_\Theta$  then

$$\begin{aligned}\mu_{\Theta|X}(A | x) &= \int_A f_{\Theta|X}(\theta | x) \mu_\Theta(d\theta), \\ \mu_{\Theta|T}(A | t) &= \int_A f_{\Theta|T}(\theta | t) \mu_\Theta(d\theta),\end{aligned}$$

holds for any  $A \in \tau$  and hence  $T$  is sufficient if and only if  $f_{\Theta|X}(\theta | x) = f_{\Theta|T}(\theta | T(x))$   $\mu_X$ -a.s.

One way to check that  $T$  is sufficient in the Bayesian sense is to check that  $\mu_{\Theta|X}(A | \cdot)$  is a function of  $T(x)$ . We have the following result.

**Lemma 2.** *Let  $T$  be a statistic and  $\mathcal{B}_T$  the sub- $\sigma$ -field of  $\mathcal{B}$  generated by  $T$ .  $T$  is sufficient in the Bayesian sense if and only if, for every prior  $\mu_\Theta$ , there exists a version of  $\mu_{\Theta|X}$  such that for each  $A \in \tau$ ,  $\mu_{\Theta|X}(A | \cdot)$  is measurable  $\mathcal{B}_T$  (In other words, it is a function of  $T(x)$ ).*

*Proof.* 'only if' part: If  $T$  sufficient in the Bayesian sense then for every prior and each  $A \in \tau$ ,  $\mu_{\Theta|X}(A | x) = \mu_{\Theta|T}(A | T(x))$  holds  $\mu_X$ -a.e. Since  $\mu_{\Theta|T}(A | T(\cdot))$  is measurable  $\mathcal{B}_T$  it follows that so is  $\mu_{\Theta|X}(A | \cdot)$ .

'if' part: Suppose that for every prior and each  $A \in \tau$ ,  $\mu_{\Theta|X}(A | \cdot)$  is measurable  $\mathcal{B}_T$ . We want to show  $T$  sufficient in Bayesian sense. That is, that  $\mu_{\Theta|X}(A | x) = \mu_{\Theta|T}(A | T(x))$   $\mu_X$ -a.s. We use the fact (e.g. Schervish, Proposition A.49 (4) p. 588) that for two functions  $f$  and  $g$  that are measurable w.r.t. a  $\sigma$ -field  $\mathcal{F}$  and a measure  $\mu$

$$\int_B f d\mu = \int_B g d\mu \text{ for each } B \in \mathcal{F} \text{ implies } f(x) = g(x), \mu - a.e.$$

Hence, in our case it is sufficient to show that for each  $B \in \mathcal{B}_T$

$$\int_B \mu_{\Theta|X}(A | x) \mu_X(dx) = \int_B \mu_{\Theta|T}(A | T(x)) \mu_X(dx).$$

The LHS is  $\Pr(\Theta \in A, X \in B)$ . Since  $B \in \mathcal{B}_T$  there is a set  $C \in \mathcal{C}$  such that  $B = T^{-1}C$ . The RHS becomes

$$\begin{aligned}\int_B \mu_{\Theta|T}(A | T(x)) \mu_X(dx) &= \{\text{change of variables}\} \\ &= \int_C \mu_{\Theta|T}(A | t) \mu_T(dt) \\ &= \Pr(\Theta \in A, T(X) \in C) \\ &= \Pr(\Theta \in A, X \in B).\end{aligned}$$

Hence, we have the desired equality and the proof is complete.  $\square$

**Example 12.** Let  $\{X_n\}$  be conditionally IID  $\text{Exp}(\theta)$  given  $\Theta = \theta$  and  $X = (X_1, \dots, X_n)$ . Put  $T(x) = x_1 + \dots + x_n$ . Let us show  $T$  is sufficient in the Bayesian sense. Let  $\mu_\Theta$  be the prior (which is arbitrary). Then the posterior distribution has density (Bayes theorem)

$$f_{\Theta|X}(\theta | x) = \frac{\prod_{i=1}^n \theta e^{-\theta x_i}}{\int \prod_{i=1}^n \psi e^{-\psi x_i} \mu_\Theta(d\psi)} = \frac{\theta e^{-\theta \sum_{i=1}^n x_i}}{\int \psi^n e^{-\psi \sum_{i=1}^n x_i} \mu_\Theta(d\psi)} = .$$

Since  $T(X)$  is  $\Gamma(n, \theta)$  given  $\Theta = \theta$  it follows that

$$f_{\Theta|T}(\theta | t) = \frac{\frac{\theta^n}{\Gamma(n)} t^{n-1} e^{-\theta t}}{\int \frac{\psi^n}{\Gamma(n)} t^{n-1} e^{-\psi t} \mu_\Theta(d\psi)} = \frac{\theta^n e^{-\theta t}}{\int \psi^n e^{-\psi t} \mu_\Theta(d\psi)}.$$

Hence  $f_{\Theta|X}(\theta | x) = f_{\Theta|T}(\theta | T(x))$  so  $T$  is sufficient in the Bayesian sense.

It is satisfying to know that in most situations one may encounter the classical and Bayesian notion of sufficiency are the same.

**Theorem 6.** *Let  $(\mathcal{T}, \mathcal{C})$  be a measurable space and  $T$  a statistic. Suppose there exists a  $\sigma$ -finite measure  $\nu$  such that  $\mu_{X|\Theta}(\cdot | \theta) \ll \nu$  for all  $\theta \in \Omega$ . Then  $T$  is sufficient in the classical sense if and only if it is sufficient in the Bayesian sense.*

*Density proof.* Suppose all relevant densities exists.

Let  $\mu_\Theta$  be an arbitrary prior. If  $T$  sufficient in the classical sense, then  $f_{X|\Theta, T}(x | \theta, t) = f_{X|T}(x | t)$ . Hence the posterior density is (with  $t = T(x)$ )

$$\begin{aligned} \frac{d\mu_{\Theta|X}}{d\mu_\Theta}(\theta | x) &= \frac{f_{X|\Theta}(x | \theta)}{\int_\Omega f_{X|\Theta}(x | \theta) \mu_\Theta(d\theta)} \\ &= \frac{f_{X|\Theta, T}(x | \theta, t) f_{T|\Theta}(t | \theta)}{\int_\Omega f_{X|\Theta, T}(x | \theta, t) f_{T|\Theta}(t | \theta) \mu_\Theta(d\theta)} \\ &= \frac{f_{X|T}(x | t) f_{T|\Theta}(t | \theta)}{\int_\Omega f_{X|T}(x | t) f_{T|\Theta}(t | \theta) \mu_\Theta(d\theta)} \\ &= \frac{f_{T|\Theta}(t | \theta)}{\int_\Omega f_{T|\Theta}(t | \theta) \mu_\Theta(d\theta)} \\ &= \frac{d\mu_{\Theta|T}}{d\mu_\Theta}(\theta | t). \end{aligned}$$

For the converse suppose that  $T$  is sufficient in the Bayesian sense so that  $\frac{d\mu_{\Theta|X}}{d\mu_\Theta}(\theta | x) = \frac{d\mu_{\Theta|T}}{d\mu_\Theta}(\theta | T(x))$ . Then, with  $t = T(x)$  we have

$$\begin{aligned} f_{X|T, \Theta}(x | t, \theta) &= \frac{f_{X|\Theta}(x | \theta)}{f_{T|\Theta}(t | \theta)} \\ &= \frac{\frac{d\mu_{\Theta|X}}{d\mu_\Theta}(\theta | x) \int_\Omega f_{X|\Theta}(x | \theta) \mu_\Theta(d\theta)}{\frac{d\mu_{\Theta|T}}{d\mu_\Theta}(\theta | t) \int_\Omega f_{T|\Theta}(t | \theta) \mu_\Theta(d\theta)} \\ &= \frac{f_X(x)}{f_T(t)} \end{aligned}$$

which does not depend on  $\theta$ . Hence,  $T$  is sufficient in the classical sense.  $\square$

**7.3. How to find a sufficient statistic?** Suppose someone hands you a parametric family  $\mathcal{P}_0 = \{f_{X|\Theta}(\cdot | \theta), \theta \in \Omega\}$  of densities w.r.t. a measure  $\nu$ . How do you come up with a sufficient statistic  $T$ ? Further, if you have come up with a suggestion of a statistic  $T$ , how do you check if it is a sufficient statistic? The next theorem gives the answer.

**Theorem 7** (Factorization Theorem, c.f. Theorem 6.2.6 in Casella & Berger). *Suppose  $\mathcal{P}_0 = \{P_\theta : \theta \in \Omega\}$  is a parametric family and there exists a  $\sigma$ -finite  $\nu$  such*

that  $P_\theta \ll \nu$  for all  $\theta \in \Omega$  with  $dP_\theta/d\nu(x) = f_{X|\Theta}(x | \theta)$ . Then  $T(X)$  is sufficient for  $\Theta$  (in either sense) if and only if there exist functions  $h$  and  $g$  such that

$$f_{X|\Theta} = h(x)g(\theta, T(x)).$$

*Density proof.* Supposing all the relevant densities exist.

It is sufficient to check the equivalence in the Bayesian sense. If  $f_{X|\Theta}(x | \theta) = h(x)g(\theta, T(x))$ , the by Bayes' theorem

$$\begin{aligned} \frac{d\mu_{\Theta|X}}{d\mu_\Theta}(\theta | x) &= \frac{f_{X|\Theta}(x | \theta)}{\int_\Omega f_{X|\Theta}(x | \theta) \mu_\Theta(d\theta)} \\ &= \frac{h(x)g(\theta, T(x))}{\int_\Omega h(x)g(\theta, T(x)) \mu_\Theta(d\theta)} \\ &= \frac{g(\theta, T(x))}{\int_\Omega g(\theta, T(x)) \mu_\Theta(d\theta)}, \end{aligned}$$

which is a function of  $T(x)$ . Hence it is sufficient in the Bayesian sense (and also in the classical sense). Conversely, suppose  $T(X)$  is sufficient in the Bayesian sense so that  $f_{\Theta|X}(\theta | x) = f_{\Theta|T}(\theta | T(x))$ . Then

$$f_{X|\Theta}(x | \theta) = f_{\Theta|X}(\theta | x) f_X(x) = \underbrace{f_X(x)}_{h(x)} \underbrace{f_{\Theta|T}(\theta | T(x))}_{g(\theta, T(x))}$$

□

**Example 13** (Exponential families). If we put  $T(X) = (t_1(X), \dots, t_k(X))$  then by the factorization theorem it follows that  $T(X)$  is sufficient. Indeed,

$$f_{X|\Theta}(x | \theta) = \underbrace{h(x)}_{h(x)} \underbrace{c(\theta) \exp\left\{\sum_{i=1}^k \theta_i t_i(x)\right\}}_{g(\theta, T(x))}.$$

Hence, a sufficient statistic always exists. We can compute the density of the sufficient statistics.

## 8. A FIRST GLANCE AT DECISION THEORY

Many statistical problems can be phrased in the language of decision theory. Suppose as usual that we have data  $X$  whose distribution depend on a parameter  $\Theta$ . Based on observing  $X = x$  we want to take some action. Let  $\mathbb{N}$  be a set of possible actions. On  $\mathbb{N}$  we take a  $\sigma$ -field  $\alpha$ . The result of our action depend of course on the chosen action, but also on the parameter  $\Theta$ . We say that every action induces a loss. A *loss function* is a function  $L : \Omega \times \mathbb{N} \rightarrow \mathbb{R}$ . We interpret  $L(\theta, a)$  as the incurred loss if we took action  $a$  and  $\Theta = \theta$ .

*One could let the loss depend on some other unobserved quantity  $V$  but we will not need this higher generality right now.*

**Definition 9.** A deterministic decision rule is a measurable function  $\delta$  from  $\mathcal{X}$  to  $\mathbb{N}$ . We interpret  $\delta(x)$  as the action to take if we observe  $X = x$ .

A randomized decision rule is a mapping from  $\mathcal{X}$  to a probability measure on  $(\mathbb{N}, \alpha)$  such that  $x \mapsto \delta(A; x)$  is measurable for each  $A \in \alpha$ .

We think of executing a randomized decision rule as follows. Given  $X = x$  we “throw a coin” according to  $\delta(\cdot; x)$ . This gives us an action  $a \in \aleph$  which is the decision we take. A deterministic decision rule can be thought of as a special case of a randomized decision rule where all the probability mass is placed at a single action. In this case the action  $\delta(x)$  is identified with the probability measure on  $\aleph$  given by  $\delta(A; x) = I_A(\delta(x))$ .

If  $\delta$  is a deterministic rule, then we associate the loss  $L(\theta, \delta(x))$ . If  $\delta$  is a randomized decision rule we associate the loss  $L(\theta, \delta(\cdot; x)) = \int_{\aleph} L(\theta, a) \delta(da; x)$ . That is, the average loss when we draw the action from  $\delta(\cdot; x)$ .

In the Bayesian case, one introduces the *posterior risk function*

$$r(\delta | x) = \int_{\Omega} L(\theta, \delta(x)) \mu_{\Theta|X}(d\theta | x).$$

That is, the average loss for decision rule  $\delta$  given the observation  $X = x$ . One would like to find a decision rule that minimizes the posterior risk simultaneously for all  $x \in \mathcal{X}$ .

**Definition 10.** If  $\delta_0$  is a decision rule such that for all  $x$ ,  $r(\delta_0 | x) < \infty$  and for all  $x$  and all decision rules  $\delta$ ,  $r(\delta_0 | x) \leq r(\delta | x)$ , then  $\delta_0$  is called a *formal Bayes rule*.

If  $\delta_0$  is a decision rule and there exists a subset  $B \subset \mathcal{X}$  such that for all  $x \in B$ ,  $r(\delta_0 | x) < \infty$  and for all  $x \in B$  and all decision rules  $\delta$ ,  $r(\delta_0 | x) \leq r(\delta | x)$ , then  $\delta_0$  is called a *partial Bayes rule*.

In classical decision theory we condition on  $\Theta = \theta$  and introduce the *risk function*

$$R(\theta, \delta) = \int_{\mathcal{X}} L(\theta, \delta(x)) \mu_{X|\Theta}(dx | \theta).$$

That is, the conditional mean of the loss, given  $\Theta = \theta$ . Here we would like to find a rule  $\delta$  that minimizes the risk function simultaneously for all values of  $\theta$ .

**8.1. A coin tossing example.** Consider the following situation. You have an amount of  $m$  dollars to bet on the outcome of a Bernoulli random variable  $X_{n+1}$ . You observe  $X = (X_1, \dots, X_n)$ . Suppose  $X_1, \dots, X_{n+1}$  are conditionally iid  $\text{Ber}(\theta)$  random variables given  $\Theta = \theta$ . Based on the observations in  $X$  you have to make a decision whether to bet on  $X_{n+1} = 0$  or  $X_{n+1} = 1$ . If you win, you win the amount  $m$  and otherwise you lose  $m$ .

Formulate this as a Bayesian decision problem. Write down the sample space  $\mathcal{X}$ , the parameter space  $\Omega$ , and the action space  $\aleph$ . Choose an appropriate prior distribution and an appropriate loss function of your choice. Then find the best decision rule, i.e. the decision rule  $\delta$  that minimizes the posterior risk simultaneously for all  $x$ .

**8.2. Sufficient statistics in decision theory.** If  $T$  is a sufficient statistic we would expect that we can base our decisions on  $T$  and do not need all the information in  $X$  since  $T$  contains all information about the unknown parameter  $\Theta$ . In the Bayesian setting we have the following theorem that supports this.

**Theorem 8.** *If there is a formal Bayes rule and  $T$  is a sufficient statistics (in the Bayesian sense) then there is a formal Bayes rule which is a function of  $T$ .*

*Proof.* Let  $\delta$  be a formal Bayes rule and take  $x \in \mathcal{X}$ . Since  $T$  is sufficient we have

$$r(\delta | x) = \int_{\Omega} L(\theta, \delta) \mu_{\Theta|X}(d\theta | x) = \int_{\Omega} L(\theta, \delta) \mu_{\Theta|T}(d\theta | T(x)).$$

We claim that for each  $y$  such that  $T(x) = T(y)$  it follows that  $r(\delta | x) = r(\delta | y)$ . If not, suppose without loss of generality that  $r(\delta | x) < r(\delta | y)$  for some  $y \in \mathcal{X}$  with  $T(y) = T(x)$ . Let  $\delta'$  be a decision rule such that  $\delta'(z) = \delta(x)$  for all  $z$  such that  $T(z) = T(x)$ . Then it follows that

$$\begin{aligned} r(\delta' | y) &= \int_{\Omega} L(\theta, \delta') \mu_{\Theta|X}(d\theta | y) \\ &= \int_{\Omega} L(\theta, \delta') \mu_{\Theta|T}(d\theta | T(y)) \\ &= \int_{\Omega} L(\theta, \delta) \mu_{\Theta|T}(d\theta | T(x)) \\ &= r(\delta | x) < r(\delta | y), \end{aligned}$$

which contradicts that  $\delta$  is a formal Bayes rule. We conclude that the claim is true.

The decision rule  $\delta'$  just defined is a function of  $T(x)$  and satisfies  $r(\delta' | T(x)) = r(\delta | x)$  for each  $x$ . Hence, it is a formal Bayes rule that is a function of  $T$ .  $\square$

Note that (in the proof above) the formal Bayes rule  $\delta$  that we started with do not have to be a function of  $T$ . For instance, it may be the case that  $T(x) = T(y)$ ,  $\delta(x) \neq \delta(y)$ , and  $L(\theta, \delta(x)) = L(\theta, \delta(y))$  for each  $\theta$ . Then  $r(\delta | x) = r(\delta | y)$  although,  $\delta(x) \neq \delta(y)$ .

In the classical setting we have the following.

**Theorem 9.** *If  $\delta_0$  is a (randomized) decision rule and  $T$  is sufficient statistic (in classical sense), then there exists a decision rule  $\delta_1$  which is a function of  $T$  and  $R(\theta, \delta_0) = R(\theta, \delta_1)$  for all  $\theta$ .*

In the theorem, if  $\delta_0$  is deterministic we interpret it as the randomized rule  $\delta$  by  $\delta(A; x) = I_A(\delta_0(x))$ . That is, the probability measure on  $(\mathbb{N}, \alpha)$  that puts all its mass on  $\delta_0(x)$ .

*Proof.* Let  $A \in \alpha$  and take

$$\delta_1(A; x) = E_{\theta}[\delta_0(A; X) | T = t].$$

Since  $T$  is sufficient the expectation does not depend on  $\theta$ . We claim that for any  $\delta_0$ -integrable function  $h : \mathbb{N} \rightarrow \mathbb{R}$

$$E \left[ \int_{\mathbb{N}} h(a) \delta_0(da; X) | T = t \right] = \int_{\mathbb{N}} h(a) \delta_1(da; t).$$

To see this, start with  $h$  as indicator, then simple function, and finally measurable function. Then we see that

$$\begin{aligned}
R(\theta, \delta_1) &= \int_{\mathcal{X}} \int_{\mathbb{N}} L(\theta, a) \delta_1(da; T(x)) \mu_{X|\Theta}(dx \mid \theta) \\
&= \int_{\mathcal{X}} E \left[ \int_{\mathbb{N}} L(\theta, a) \delta_0(da; X) \mid T = T(x) \right] \mu_{X|\Theta}(dx \mid \theta) \\
&= E_{\theta} \left[ E \left[ \int_{\mathbb{N}} L(\theta, a) \delta_0(da; X) \mid T \right] \right] \\
&= E_{\theta} \left[ \int_{\mathbb{N}} L(\theta, a) \delta_0(da; X) \right] \\
&= \int_{\mathcal{X}} \int_{\mathbb{N}} L(\theta, a) \delta_0(da; x) \mu_{X|\Theta}(dx \mid \theta) \\
&= R(\theta, \delta_0).
\end{aligned}$$

□

One should note that even if  $\delta_0$  is a deterministic rule, the resulting rule  $\delta_1(A; t) = E_{\theta}[\delta_0(A; X) \mid T = t]$  may be randomized.