

1 Article

2 Information Loss due to the Compression of Sample 3 Data from Discrete Distributions

4 Maryam Moghimi^{*1,2}, H.W. Corley^{1,2}

5 ¹ Center on Stochastic Modeling, Optimization, and Statistics (COSMOS), The University of Texas at
6 Arlington, Arlington, TX, USA

7 ² The authors contributed equally to this paper.

8 * Correspondence: maryam.moghimi@uta.edu; +1-214-971-0904 (M.M.), corley@uta.edu; Tel.: +1-817-272-
9 3092 (H.C.)

10 Received: date; Accepted: date; Published: date

11 **Abstract:** In this paper we study the information lost when a real-valued statistic $T(X_1, \dots, X_n)$ is
12 used to summarize the sample data $\mathbf{x} = (x_1, \dots, x_n)$ of a random sample $\mathbf{X} = (X_1, \dots, X_n)$ from a
13 discrete random variable X with a one-dimensional parameter θ . We compare the probability that
14 the random sample \mathbf{X} yields \mathbf{x} to the probability that the compressed sample $T(\mathbf{X})$ yields $T(\mathbf{x})$.
15 The former probability measures the total information about \mathbf{x} , while the latter measures the
16 compressed information about \mathbf{x} , both of which are expressed here as Shannon information. The
17 difference is the information lost about \mathbf{X} by its compression to $T(\mathbf{X})$. We focus on sufficient
18 statistics for the parameter θ and develop a general formula independent of θ for this lost
19 information as well as for an associated entropy that depends only on T . Our approach would also
20 work for non-sufficient statistics, but the lost information and associated entropy would involve θ .
21 Examples are presented for some standard discrete distributions.

22 **Keywords:** discrete distributions, Shannon information, lost information, sampling, data reduction,
23 data compression, entropy, sufficient statistics, likelihood

24 1. Introduction

25 We consider the data sample $\mathbf{x} = (x_1, \dots, x_n)$ from a random sample $\mathbf{X} = (X_1, \dots, X_n)$ for a
26 discrete random variable X with sample space S and one-dimensional parameter θ . Here a statistic
27 $T(\mathbf{X})$ is a real-valued function of the random sample but not a function of any parameter θ associated
28 with X , though θ may fixed at an arbitrary value. The data sample \mathbf{X} is compressed to the summary
29 statistic $T(\mathbf{X})$, which could be used to characterize \mathbf{X} or to estimate θ . Such data compression is an
30 irreversible process [1] and always involves some information loss. For instance, if $T(\mathbf{X}) = \bar{X}$, the
31 original measurements \mathbf{x} cannot be reconstructed from \bar{x} , and some information about \mathbf{x} is lost.
32 Nonetheless, such data compression is frequently used to make inferences about, for example, the
33 true mean μ of X . Our information-theoretic approach to data compression generalizes the
34 observation in [2] that a binomial random variable loses all the information about the order of
35 successes in the associated sequence of Bernoulli trials.

36 For any real-valued statistic T and the given sample data \mathbf{x} , we decompose the total
37 information about \mathbf{X} available in \mathbf{x} into the sum of (a) the information available in the compressed
38 data $T(\mathbf{x}) = \bar{x}$ and (b) the information lost in the compression. When T is a sufficient statistic for θ
39 this lost information is independent of θ . Moreover, by taking the expected value of this lost
40 information over all possible data sets, we define an associated entropy measure that depends on T
41 but neither \mathbf{x} nor θ . Our approach also works for non-sufficient statistics, but the lost information
42 and associated entropy would then involve θ , and so θ must be estimated to computing these
43 quantities.

44 The paper is organized as follows. In Section 2, we present the necessary definitions, notation,
45 and preliminary results. In Section 3, we decompose the total information available about \mathbf{X} in \mathbf{x}

46 and give various expressions for the Shannon information lost by compressing \mathbf{x} to $T(\mathbf{x})$. In
 47 Section 4, we develop an entropy measure associated with this lost information. In Section 5, we
 48 present examples of our results for some standard discrete distributions and several statistics
 49 sufficient for θ . Conclusions are offered in Section 6.

50 **2. Preliminaries**

51 The following definitions, notation, and results are used here. Further details can be found in [3,4]
 52 and elsewhere. An important class of statistics is first defined.

53 **Definition 2.1 (Sufficient Statistic).** A statistic $T(\mathbf{X})$ is a sufficient statistic (SS) for the parameter
 54 θ if the probability

$$P[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})] \tag{1}$$

55 is independent of θ .

56 Note that P instead of P_θ is used in (1) since this probability is independent of θ . Also observe
 57 that (1) is not a joint conditional distribution for \mathbf{X} since its n condition changes with \mathbf{x} . This
 58 observation becomes significant in Section 4. The fact that (1) does not involve θ is used to prove the
 59 Fisher Factorization Theorem (FFT), which is the usual method for determining if a statistic is an SS for
 60 θ . We use the notation $f(\mathbf{x}|\theta)$ to denote the joint pmf of \mathbf{X} evaluated at the variable \mathbf{x} for a fixed value
 61 of θ .

62 **Result 2.2 (Fisher Factorization Theorem).** The real-valued statistic $T(\mathbf{X})$ is sufficient for θ if and
 63 only if there exist functions $g: R^1 \rightarrow R^1$ and $h: S^n \rightarrow R^1$ such that for any sample data \mathbf{x} and for all
 64 values of θ the joint pmf $f(\mathbf{x}|\theta)$ of \mathbf{X} can be factored as

$$f(\mathbf{x}|\theta) = g[T(\mathbf{x})|\theta] \times h(\mathbf{x}) \tag{2}$$

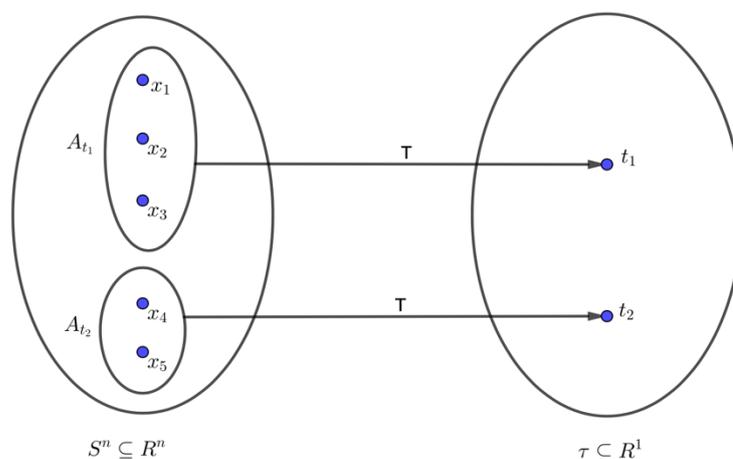
65 for real-valued, nonnegative functions g on R^1 and h on S^n . The function h does not depend on
 66 θ , while g does depend on \mathbf{x} but only through $T(\mathbf{x})$.

67 We focus on a sufficient statistic T for θ in Section 3, where we need the notion of a partition [5]
 68 as defined next.

69 **Definition 2.3 (Partition).** Let S be the denumerable sample space of the discrete random variable
 70 X , and thus let S^n be the denumerable sample space of the random sample \mathbf{X} . For any statistic
 71 $T: S^n \rightarrow R^1$, let τ_T be the denumerable set $\tau_T = \{t | \exists \mathbf{x} \in S^n \text{ for which } t = T(\mathbf{x})\}$, which is the range of
 72 T . Then T partitions the sample space S^n into the mutually exclusive and collectively exhaustive
 73 partition sets $A_t = \{\mathbf{x} \in S^n | T(\mathbf{x}) = t\}, \forall t \in \tau_T$.

74 Figure 2.1 below illustrates the situation.

75 Figure 2.1



76
 77

We also need the well-known likelihood function.

78 **Definition 2.4 (Likelihood Function).** Let \mathbf{x} be sample data from a random sample \mathbf{X} from a
 79 discrete random variable X with sample space S and real-valued parameter θ , and let $f(\mathbf{x}|\theta)$
 80 denote the joint pmf of the random sample \mathbf{X} . For any sample data \mathbf{x} , the likelihood function of θ is
 81 defined as

$$L(\theta|\mathbf{x}) = f(\mathbf{x}|\theta). \quad (3)$$

82 The likelihood function in (3) is a function of the variable θ for given data \mathbf{x} . However, the joint
 83 pmf $f(\mathbf{x}|\theta)$ as a function of \mathbf{x} for fixed θ is frequently called the likelihood function as well. In this
 84 case we also write the joint pmf as $L(\mathbf{x}|\theta)$. We distinguish the two cases since $L(\theta|\mathbf{x})$ is not a statistic
 85 but $L(\mathbf{x}|\theta)$ is one that incorporates all available information about \mathbf{X} . Moreover, $L(\mathbf{x}|\theta)$ is an SS for θ
 86 [4] and uniquely determines an associated SS called the likelihood kernel to be used in subsequent
 87 examples.

88 **Definition 2.5 (Likelihood kernel).** Let S be the sample space of \mathbf{X} . For fixed θ , suppose that
 89 $L(\mathbf{x}|\theta)$ can be factored as

$$L(\mathbf{x}|\theta) = K(\mathbf{x}|\theta) \times R(\mathbf{x}), \quad \forall \mathbf{x} \in S^n, \quad (4)$$

90 where $K: S^n \rightarrow \mathbb{R}^1$ and $R: S^n \rightarrow \mathbb{R}^1$ have the following properties:

- 91 (a) every nonnumerical factor of $K(\mathbf{x}|\theta)$ contains θ ;
- 92 (b) $R(\mathbf{x})$ does not contain θ ;
- 93 (c) for $\forall \mathbf{x} \in S^n$, both $K(\mathbf{x}|\theta) \geq 0$ and $R(\mathbf{x}) \geq 0$; and
- 94 (d) $K(\mathbf{x}|\theta)$ is not divisible by any positive number except 1.

95 Then $K(\mathbf{x}|\theta)$ is defined as the likelihood kernel of $L(\mathbf{x}|\theta)$ and $R(\mathbf{x})$ the residue of $L(\mathbf{x}|\theta)$.

96 **Theorem 2.6.** The likelihood kernel $K(\mathbf{x}|\theta)$ has the following properties.

- 97 (i) $K(\mathbf{x}|\theta)$ uniquely exists.
- 98 (ii) $K(\mathbf{x}|\theta)$ is an SS for θ .
- 99 (iii) For any θ_1 and θ_2 , the likelihood ratio $\frac{L(\mathbf{x}|\theta_1)}{L(\mathbf{x}|\theta_2)}$ equals $\frac{K(\mathbf{x}|\theta_1)}{K(\mathbf{x}|\theta_2)}$.

100 **Proof.** To prove (i), for fixed θ we first show that the likelihood kernel $K(\mathbf{x}|\theta)$ of **Definition 2.5**
 101 exists by construction. Since the formula for $L(\theta|\mathbf{x}) = f(\mathbf{x}|\theta)$ must explicitly contain θ , the parameter
 102 θ cannot appear only in the range of \mathbf{x} . Hence $L(\mathbf{x}|\theta)$ as a function of \mathbf{x} can be factored into
 103 $K(\mathbf{x}|\theta) \times R(\mathbf{x})$ satisfying (a) and (b) of **Definition 2.5**, where $K(\mathbf{x}|\theta) \geq 0, \forall \mathbf{x} \in S^n$, and the numerical
 104 factor of $K(\mathbf{x}|\theta)$ is either $+1$ or -1 . Then $R(\mathbf{x}) \geq 0, \forall \mathbf{x} \in S^n$, since $K(\mathbf{x}|\theta) \geq 0, \forall \mathbf{x} \in S^n$, and
 105 $K(\mathbf{x}|\theta) \times R(\mathbf{x}) = f(\mathbf{x}|\theta) \geq 0$. Thus (c) is satisfied. Finally, the only positive integer that evenly divides
 106 $+1$ or -1 is 1, so (d) holds. It follows that the likelihood kernel $K(\mathbf{x}|\theta)$ and its associated $R(\mathbf{x})$ in
 107 **Definition 2.5** are well defined and exist.

108 We next show that $K(\mathbf{x}|\theta)$ as constructed above is unique. Let $K_1(\mathbf{x}|\theta)$ with residue $R_1(\mathbf{x})$ and
 109 $K_2(\mathbf{x}|\theta)$ with $R_2(\mathbf{x})$ both satisfy **Definition 2.5**. Thus for $j = 1, 2$, $R_j(\mathbf{x})$ does not contain θ while
 110 every nonnumerical factor of $K_j(\mathbf{x}|\theta)$ does contain θ . It follows that $K_1(\mathbf{x}|\theta) \geq 0$ and $K_2(\mathbf{x}|\theta) \geq 0$
 111 must be identical or else be a positive multiple of one another. Assume that $K_2(\mathbf{x}|\theta) = \lambda K_1(\mathbf{x}|\theta)$ for
 112 some $\lambda > 0$. If $\lambda \neq 1$, $K_2(\mathbf{x}|\theta)$ is divisible by a positive number other than 1 to avoid (d). Thus,
 113 $K(\mathbf{x}|\theta)$ is unique.

114 To prove (ii) we show that this unique $K(\mathbf{x}|\theta)$ is an SS for θ . For $L(\theta|\mathbf{x}) = f(\mathbf{x}|\theta)$, let $g[z] = z$
 115 and $h(\mathbf{x}) = R(\mathbf{x})$ in (2). Then $L(\theta|\mathbf{x}) = f(\mathbf{x}|\theta) = g[K(\mathbf{x}|\theta)] \times h(\mathbf{x}) = K(\mathbf{x}|\theta) \times R(\mathbf{x})$. Thus $K(\mathbf{x}|\theta)$ is an
 116 SS by the FFT of **Result 2.2**.

117 Finally, (iii) follows immediately from **Definition 2.5** and the fact that $L(\mathbf{x}|\theta_2) \neq 0$ for $\mathbf{x} \in S^n$. ■

118 We next discuss the notion of information to be used here. Actually, probability itself is a measure
 119 of information in the sense that it captures the surprise level of an event. An observer obtains more
 120 information, i.e., surprise, if an unlikely event occurs than if a likely one does. Instead of probability,
 121 however, we use the additive measure known as Shannon information [6, 7] defined as follows.

122 **Definition 2.7 (Shannon Information).** Let \mathbf{x} be sample data for the random sample \mathbf{X} from the
 123 discrete random variable X with a one-dimensional parameter θ , and let $f(\mathbf{x}|\theta)$ be the joint pmf of
 124 \mathbf{X} at \mathbf{x} . The Shannon information obtained from the sample data \mathbf{x} is defined as

$$I(\mathbf{x}|\theta) = -\log f(\mathbf{x}|\theta), \quad (5)$$

125 where the units of $I(\mathbf{x}|\theta)$ is bits if the base of the logarithm is 2, which is to be used here.

126 The expected information over $\forall \mathbf{x} \in S^n$ will also be used.

127 **Definition 2.8 (Entropy).** Under the conditions of **Definition 2.7**, the entropy $H(\mathbf{X}|\theta)$ is defined
128 as the expected value of $I(\mathbf{X}|\theta)$; i.e.,

$$H(\mathbf{X}|\theta) = \sum_{\mathbf{x}} f(\mathbf{x}|\theta)I(\mathbf{x}|\theta). \quad (6)$$

129 Since entropy is the expected information over all possible random samples, it measures the
130 available information about \mathbf{X} better than would a single data set \mathbf{x} , which might not be typical [8].
131 We next give a method to obtain the information loss about \mathbf{X} that occurs when a data set \mathbf{x} is
132 compressed to $T(\mathbf{x})$. In our approach, we focus on a sufficient statistic T so there will be no θ in (5)
133 for the lost information below. However, our approach is applicable to a non-sufficient statistic as
134 well if θ is estimated from the data.

135 3. Information Decomposition under Data Compression by a Real-Valued Statistic

136 We now develop a procedure to determine how much information about \mathbf{X} contained in a data
137 set \mathbf{x} is lost when the data is compressed to $T(\mathbf{x})$ by the sufficient statistic T . Consider the joint
138 conditional probability

$$P_{\theta}[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})], \quad (7)$$

139 which is identified with the probabilistic information lost about the event $\mathbf{X} = \mathbf{x}$ by the data
140 compression of \mathbf{x} to $T(\mathbf{x})$. The notation P_{θ} refers to the fact that the discrete probability (7) in
141 general involves the parameter θ . We next express (7) using the definition of conditional probability
142 to obtain the basis of our development. **Result 3.1** is given in [3, p. 273] and proved below to illustrate
143 the reasoning.

144 **Result 3.1.** Let \mathbf{x} be sample data for a random sample \mathbf{X} from a discrete random variable X
145 with sample space S and real-valued parameter θ , and let $T(\mathbf{X})$ be any real-valued statistic. Then

$$P_{\theta}[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})] = \frac{P_{\theta}[\mathbf{X} = \mathbf{x}]}{P_{\theta}[T(\mathbf{X}) = T(\mathbf{x})]}. \quad (8)$$

146 **Proof.** Using the definition of conditional probability, rewrite (7) as

$$\frac{P_{\theta}[\mathbf{X} = \mathbf{x}; T(\mathbf{X}) = T(\mathbf{x})]}{P_{\theta}[T(\mathbf{X}) = T(\mathbf{x})]}. \quad (9)$$

147 But $T(\mathbf{X}) = T(\mathbf{x})$ whenever $\mathbf{X} = \mathbf{x}$, so (8) follows. ■

148 Observe that if T is an SS for θ , the right side of (8) is independent of θ and hence so is the left.
149 Now taking the negative logarithm of (8) and rearranging terms gives

$$-\log P_{\theta}[\mathbf{X} = \mathbf{x}] = -\log P_{\theta}[T(\mathbf{X}) = T(\mathbf{x})] - \log P_{\theta}[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})]. \quad (10)$$

150 From (8) note that $P_{\theta}[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})] \geq P_{\theta}[\mathbf{X} = \mathbf{x}]$ since $P_{\theta}[T(\mathbf{X}) = T(\mathbf{x})] \leq 1$, so
151 $-\log P_{\theta}[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})] \leq -\log P_{\theta}[\mathbf{X} = \mathbf{x}]$. Similarly, $-\log P_{\theta}[T(\mathbf{X}) = T(\mathbf{x})] \leq -\log P_{\theta}[\mathbf{X} = \mathbf{x}]$.
152 These facts suggest that the left side of (10) is the total Shannon information in bits about \mathbf{X} contained
153 in the sample data \mathbf{x} . On the right side of (10), the term $-\log P_{\theta}[T(\mathbf{X}) = T(\mathbf{x})]$ is considered the
154 information about \mathbf{X} contained in the compressed data summary $T(\mathbf{x})$, and the term
155 $-\log P_{\theta}[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})]$ is identified as the information about \mathbf{X} that has been lost as the result
156 of the data compression by $T(\mathbf{x})$.

157 In particular, this lost information represents a combinatorial loss in the sense that multiple \mathbf{x} 's
158 may give the same value $T(\mathbf{x}) = t$ as depicted in Figure 2.1 above. In other words, the lost information
159 $-\log P_{\theta}[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})]$ is a measure of the knowledge unavailable about the data sample \mathbf{x}

160 when only the compressed data summary $T(\mathbf{x})$ is known and not \mathbf{x} itself. For a sufficient statistic
 161 $T(\mathbf{X})$ for θ , this lost information is independent of θ . It is a characteristic of $T(\mathbf{X})$ for the given data
 162 sample \mathbf{x} .

163 In terms of Figure 2.1 above, the situation may be described as follows. On the left is the sample
 164 space $S^n \subseteq \mathbf{R}^n$ over which probabilities on \mathbf{X} are computed. On the right is the range $\tau_T \subseteq \mathbf{R}^1$ of T
 165 over which the probability of $T(\mathbf{X})$ are computed. T compresses the data sample \mathbf{x} into $T(\mathbf{x})$,
 166 where multiple \mathbf{x} 's may give the same $T(\mathbf{x}) = t$. In Figure 2.1 the distinct data samples \mathbf{x}_1 , \mathbf{x}_2 , and
 167 \mathbf{x}_3 are all compressed into the same value t_1 . But knowing that $T(\mathbf{x}) = t_1$ for some data sample \mathbf{x}
 168 does not provide sufficient information to know unequivocally, for example, that $\mathbf{x} = \mathbf{x}_1$. Information
 169 is lost in the compression. One can also say that the total information $-\log P_\theta[\mathbf{X} = \mathbf{x}]$ deriving from the
 170 left side of Figure 2.1 is compressed to $-\log P_\theta[T(\mathbf{X}) = T(\mathbf{x})]$ deriving from the right. The reduction
 171 of information from the left to the right side is precisely the lost information
 172 $-\log P_\theta[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})]$. For fixed t , it is lost due to the ambiguity as to which data sample on
 173 the left actually gave t when only t is known. There is no ambiguity when T is one-to-one.

174 The general decomposition of information in (10) is next summarized in **Definition 3.2**, where
 175 T does not need to be sufficient for θ .

176 **Definition 3.2 ($I_{\text{total}}, I_{\text{comp}}, I_{\text{lost}}$).** Let \mathbf{x} be sample data for a random sample \mathbf{X} from a discrete
 177 random variable X with sample space S and real-valued parameter θ . For any real-valued statistic.
 178 $T(\mathbf{X})$, the Shannon information about \mathbf{X} obtained from the sample data \mathbf{x} can be decomposed as

$$I_{\text{total}}(\mathbf{x}|\theta) = I_{\text{comp}}(\mathbf{x}|\theta, T) + I_{\text{lost}}(\mathbf{x}|\theta, T), \tag{11}$$

179 where

$$I_{\text{total}}(\mathbf{x}|\theta) = -\log P_\theta[\mathbf{X} = \mathbf{x}], \tag{12}$$

$$I_{\text{comp}}(\mathbf{x}|\theta, T) = -\log P_\theta[T(\mathbf{X}) = T(\mathbf{x})], \tag{13}$$

180 and

$$I_{\text{lost}}(\mathbf{x}|\theta, T) = -\log P_\theta[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})]. \tag{14}$$

181 Both **Result 3.1** and **Definition 3.2** are valid for any real-valued statistic for \mathbf{X} . The notation
 182 $I_{\text{total}}(\mathbf{x}|\theta)$ indicates that I_{total} is a function of the sample data \mathbf{x} for a fixed but arbitrary parameter
 183 value θ . Similarly, both $I_{\text{comp}}(\mathbf{x}|\theta, T)$ and $I_{\text{lost}}(\mathbf{x}|\theta, T)$ are functions of \mathbf{x} for fixed θ and T .
 184 However, in this paper we focus on sufficient statistics, which provide a simpler expression for
 185 $I_{\text{lost}}(\mathbf{x}|\theta, T)$ that does not involve θ . For a sufficient statistic T for θ , we use the notation $I_{\text{lost}}(\mathbf{x}|T)$
 186 for the lost information, though $I_{\text{total}}(\mathbf{x}|\theta)$ and $I_{\text{comp}}(\mathbf{x}|\theta, T)$ still require θ . The next result is an
 187 application of the FFT of **Result 2.2**.

188 **Theorem 3.3 (Lost Information for an SS).** Let \mathbf{x} be sample data for a random sample \mathbf{X} from
 189 a discrete random variable X with sample space S and real-valued parameter θ . Let T be an SS for
 190 θ , $f(\mathbf{x}|\theta)$ be the joint pmf of \mathbf{X} , and $f(\mathbf{x}|\theta) = g[T(\mathbf{x})|\theta] \times h(\mathbf{x})$ as in **Result 2.2**. Then for all $\mathbf{x} \in S^n$

$$I_{\text{lost}}(\mathbf{x}|T) = -\log \frac{h(\mathbf{x})}{\sum_{\mathbf{y} \in A_{T(\mathbf{x})}} h(\mathbf{y})}, \tag{15}$$

191 where $A_{T(\mathbf{x})}$ is defined in **Definition 2.3** for $t = T(\mathbf{x})$.

192 **Proof.** Let $\mathbf{x} \in S^n$. Then $f(\mathbf{x}|\theta) > 0$ since \mathbf{x} is a realization of \mathbf{X} . Since T is an SS, we write (7)
 193 without θ . Then it suffices to establish that

$$P[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})] = \frac{h(\mathbf{x})}{\sum_{\mathbf{y} \in A_{T(\mathbf{x})}} h(\mathbf{y})}, \tag{16}$$

194 from which (15) immediately follows. Rewrite (8) as

$$P[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})] = \frac{P_\theta[\mathbf{X} = \mathbf{x}]}{P_\theta[T(\mathbf{X}) = T(\mathbf{x})]} = \frac{f(\mathbf{x}|\theta)}{\sum_{\mathbf{y} \in A_{T(\mathbf{x})}} f(\mathbf{y}|\theta)}, \quad (17)$$

195 so from (17) and (2), then

$$P[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})] = \frac{g[T(\mathbf{x})|\theta] \times h(\mathbf{x})}{\sum_{\mathbf{y} \in A_{T(\mathbf{x})}} g[T(\mathbf{y})|\theta] \times h(\mathbf{y})}. \quad (18)$$

196 But $T(\mathbf{y}) = T(\mathbf{x}), \forall \mathbf{y} \in A_{T(\mathbf{x})}$ in (18), so

$$P[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})] = \frac{g[T(\mathbf{x})|\theta] \times h(\mathbf{x})}{g[T(\mathbf{x})|\theta] \times \sum_{\mathbf{y} \in A_{T(\mathbf{x})}} h(\mathbf{y})}, \forall \mathbf{x} \in S^n. \quad (19)$$

197 Since $f(\mathbf{x}|\theta) > 0$ and hence $g[T(\mathbf{x})|\theta] \neq 0$, this term can be canceled on the right side of (19) to yield
 198 (16). Taking $-\log$ of (16) completes the proof. ■

199 Now consider **Theorem 3.3** when each A_t is a singleton in (16), i.e., when T is a one-to-one
 200 function. In this extreme case, $P[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})] = 1$ since $\sum_{\mathbf{y} \in A_{T(\mathbf{x})}} h(\mathbf{y}) = h(\mathbf{x})$ in the
 201 denominator of the right side of (16). Thus $I_{\text{lost}}(\mathbf{x}|T) = 0$ from which $I_{\text{comp}}(\mathbf{x}|\theta, T) = I_{\text{total}}(\mathbf{x}|\theta)$ for
 202 all \mathbf{x} in S^n . Thus the special case of a one-to-one T justifies the identification of the lost information
 203 as $I_{\text{lost}}(\mathbf{x}|\theta, T) = -\log P_\theta[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})]$. In other words, for all data samples $\mathbf{x}, \mathbf{y} \in S^n$, if $\mathbf{x} \neq \mathbf{y}$
 204 whenever $T(\mathbf{x}) \neq T(\mathbf{y})$, then $P[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})]$ is not diminished by the compression of the
 205 singleton $A_{T(\mathbf{x})}$ to the number $T(\mathbf{x})$.

206 More generally, it is also true that $I_{\text{lost}}(\mathbf{x}|\theta, T) = 0$ when T is one-to-one but not sufficient for
 207 θ . In this case, write $P_\theta[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})] = \frac{P_\theta[\mathbf{X}=\mathbf{x}]}{P_\theta[T(\mathbf{X})=T(\mathbf{x})]} = \frac{f(\mathbf{x}|\theta)}{\sum_{\mathbf{y} \in A_{T(\mathbf{x})}} f(\mathbf{y}|\theta)}$. But since T is one-to-one,
 208 $\sum_{\mathbf{y} \in A_{T(\mathbf{x})}} f(\mathbf{y}|\theta) = f(\mathbf{x}|\theta)$, $P_\theta[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})] = 1$, and again $I_{\text{lost}}(\mathbf{x}|\theta, T) = 0$.

209 Now consider the other extreme case where $T(\mathbf{x}) = c$ is constant on S^n . Then
 210 $P_\theta[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = c] = \frac{P_\theta[\mathbf{X}=\mathbf{x}]}{P_\theta[T(\mathbf{X})=c]}$. But $P_\theta[T(\mathbf{X}) = c] = 1$, so $P_\theta[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = c] = P_\theta[\mathbf{X} = \mathbf{x}]$ and
 211 $I_{\text{lost}}(\mathbf{x}|\theta, T) = I_{\text{total}}(\mathbf{x}|\theta, T)$ on S^n . In this case, $I_{\text{comp}}(\mathbf{x}|\theta, T) = 0$ because the event $T(\mathbf{x}) = c$ gives no
 212 information about \mathbf{x} .

213 Next, in the following corollary we show that (16) can be simplified when T is the likelihood
 214 function.

215 **Corollary 3.4 (Information Loss for Likelihood Function).** Under the assumptions of **Theorem**
 216 **3.3**, if $T(\mathbf{x}) = L(\mathbf{x}|\theta)$, then

$$I_{\text{lost}}(\mathbf{x}|L) = -\log \frac{1}{|A_{L(\mathbf{x}|\theta)}|}, \quad (20)$$

217 where $|A_{L(\mathbf{x}|\theta)}|$ is the cardinality of the partition set A_t for $t = L(\mathbf{x}|\theta)$.

218 **Proof.** For $T(\mathbf{x}) = L(\mathbf{x}|\theta) = f(\mathbf{x}|\theta)$ in (2), let g be the identity function and $h(\mathbf{x}) = 1$. Then
 219 substituting $h(\mathbf{x}) = 1$ into (16) gives the denominator $\sum_{\mathbf{y} \in A_{L(\mathbf{x}|\theta)}} 1 = |A_{L(\mathbf{x}|\theta)}|$ to yield (20). ■

220 We next state a reproductive property of a statistic T' that is a one-to-one function of a sufficient
 221 statistic T for θ .

222 **Theorem 3.5.** If there is a one-to-one function between a sufficient statistic T for θ and an
 223 arbitrary real-valued statistic T' on S^n , the following hold.

- 224 (i) T' is also an SS.
- 225 (ii) T and T' partition the sample space S into the same partition sets.
- 226 (iii) $I_{\text{lost}}(\mathbf{x}|T) = I_{\text{lost}}(\mathbf{x}|T'), \forall \mathbf{x} \in S^n$.

227 **Proof.** To prove (i), let u be a real-valued one-to-one function of T' such that

$$T(\mathbf{x}) = u[T'(\mathbf{x})]. \quad (21)$$

228 Since T is an SS, by equation (2) there are real-valued functions g on \mathbb{R}^1 and h on S^n for which

$$f(\mathbf{x}|\theta) = g[T(\mathbf{x})|\theta] \times h(\mathbf{x}). \quad (22)$$

229 By substituting $T(\mathbf{x})$ from (21) in (22), we get

$$f(\mathbf{x}|\theta) = g(u[T'(\mathbf{x})|\theta]) \times h(\mathbf{x}), \quad (23)$$

230 which can be rewritten as

$$f(\mathbf{x}|\theta) = (g \circ u)[T'(\mathbf{x})|\theta] \times h(\mathbf{x}). \quad (24)$$

231 Since T' in (24) satisfies the condition of **Result 2.2** for $g' = g \circ u$, T' is an SS.

232 To prove (ii), we use **Definition 2.3**. Let T partition the sample space S^n into the mutually
 233 exclusive and collectively exhaustive sets $A_t = \{\mathbf{x}|T(\mathbf{x}) = t\}, \forall t \in \tau_T$. By equation (21) we can also write
 234 A_t as

$$A_t = \{\mathbf{x}|u[T'(\mathbf{x})] = t\}, \forall t \in \tau_T. \quad (25)$$

235 Since u is a one-to-one function, it has an inverse u^{-1} . Letting $u^{-1}(t) = t'$, we apply u^{-1} to the right
 236 side of (25) and get

$$A_t = \{\mathbf{x}|T'(\mathbf{x}) = t'\}, \forall t' \in u(\tau_T). \quad (26)$$

237 But $u(\tau_T) = \tau_{T'}$ and the cardinalities $|\tau_T| = |\tau_{T'}|$, so the right side of (26) is $A_{t'}$ and

$$A_t = A_{t'}. \quad (27)$$

238 Finally, to get (iii) we use **Theorem 3.3** to calculate information lost over two statistics T and T' .
 239 Since $h(\mathbf{x})$ is the same in (22) and (24) and since equation (27) holds, we sum $h(\mathbf{x})$ over the same sets
 240 in the denominator of equation (16) for both T and T' to give

$$I_{\text{lost}}(\mathbf{x}|T) = I_{\text{lost}}(\mathbf{x}|T') \quad (28)$$

241 and complete the proof. ■

242 We next compare the information loss of the sufficient statistic $L(\mathbf{x}|\theta)$ to other sufficient statistics.
 243 For the sufficient statistic $K(\mathbf{x}|\theta)$, a lemma is needed.

244 **Lemma 3.6.** Let \mathbf{x} be any data sample for a random sample \mathbf{X} from the discrete random variable
 245 X with real-valued parameter θ . Then $K(\mathbf{x}|\theta)$ is a function of $L(\mathbf{x}|\theta)$ and $\tau_L \geq \tau_K$.

246 **Proof.** From [3, p. 280], $K(\mathbf{x}|\theta)$ is a function of $L(\mathbf{x}|\theta)$ if and only if $K(\mathbf{x}|\theta) = K(\mathbf{y}|\theta)$ whenever
 247 $L(\mathbf{x}|\theta) = L(\mathbf{y}|\theta)$. For all data samples \mathbf{x} and \mathbf{y} , we thus prove that if $L(\mathbf{x}|\theta) = L(\mathbf{y}|\theta)$, then $K(\mathbf{x}|\theta) =$
 248 $K(\mathbf{y}|\theta)$. Thus suppose that $L(\mathbf{x}|\theta) = L(\mathbf{y}|\theta)$. By **Definition 2.5** we can decompose $L(\mathbf{x}|\theta)$ and $L(\mathbf{y}|\theta)$
 249 into $K(\mathbf{x}|\theta)R(\mathbf{x})$ and $K(\mathbf{y}|\theta)R(\mathbf{y})$, respectively. Note that $K(\mathbf{y}|\theta) \neq 0$. Otherwise $L(\mathbf{y}|\theta) = 0$ in
 250 contradiction to \mathbf{y} being sample data with a nonzero probability of occurring. Write

$$\frac{K(\mathbf{x}|\theta)}{K(\mathbf{y}|\theta)} = \frac{R(\mathbf{y})}{R(\mathbf{x})}. \quad (29)$$

251 Suppose that $K(\mathbf{x}|\theta) \neq K(\mathbf{y}|\theta)$ so that $\frac{K(\mathbf{x}|\theta)}{K(\mathbf{y}|\theta)} = \frac{R(\mathbf{y})}{R(\mathbf{x})} \neq 1$ in (29). From **Definition 2.5**, every
 252 nonnumerical factor of $K(\mathbf{x}|\theta)$ and $K(\mathbf{y}|\theta)$ contains θ , and neither $K(\mathbf{x}|\theta)$ nor $K(\mathbf{y}|\theta)$ is divisible by
 253 any positive number except the number 1. Hence, since $\frac{R(\mathbf{y})}{R(\mathbf{x})}$ does not contain θ , the nonnumerical
 254 factors of $K(\mathbf{x}|\theta)$ and $K(\mathbf{y}|\theta)$ must cancel in (29) and the remaining numerical factors could not be
 255 identical. Thus at least one of these factors would be divisible by a positive number other than 1 in
 256 contradiction to **Definition 2.5**. It now follows that $K(\mathbf{x}|\theta) = K(\mathbf{y}|\theta)$, so $K(\mathbf{x}|\theta)$ is some function u of
 257 $L(\mathbf{x}|\theta)$. Finally, $\tau_L \geq \tau_K$ since this function u is surjective from S^n onto its image $u(S^n)$. ■

258 **Lemma 3.7.** Under the conditions of **Lemma 3.6**, the sufficient statistics L and K satisfy

$$I_{\text{comp}}(\mathbf{x}|\theta, L) \geq I_{\text{comp}}(\mathbf{x}|\theta, K), \forall \mathbf{x} \in S^n. \quad (30)$$

259 **Proof.** Let $\mathbf{x} \in S^n$ and suppose that $\mathbf{y} \in A_{L(\mathbf{x})}$. Then $L(\mathbf{y}|\theta) = L(\mathbf{x}|\theta)$, so it follows from
 260 **Lemma 3.6** that $K(\mathbf{y}|\theta) = K(\mathbf{x}|\theta)$ and thus $\mathbf{y} \in A_{K(\mathbf{x})}$. Hence $A_{L(\mathbf{x})} \subseteq A_{K(\mathbf{x})}$, and so

$$P_\theta[L(\mathbf{X}|\theta) = L(\mathbf{x}|\theta)] = \sum_{\mathbf{y} \in A_{L(\mathbf{x})}} f(\mathbf{x}|\theta) \leq \sum_{\mathbf{y} \in A_{K(\mathbf{x})}} f(\mathbf{x}|\theta) = P_\theta[K(\mathbf{X}|\theta) = K(\mathbf{x}|\theta)], \forall \mathbf{x} \in S^n \quad (31)$$

261 Taking the Shannon information of both sides of the inequality in (31) and using (13) gives (30). ■

262 **Theorem 3.8.** Let \mathbf{x} be sample data for a random sample \mathbf{X} from a discrete random variable X
 263 with the real-valued parameter θ . Then for all $\mathbf{x} \in S^n$,

$$I_{\text{lost}}(\mathbf{x}|L) \leq I_{\text{lost}}(\mathbf{x}|K). \quad (32)$$

264 **Proof.** Let $\mathbf{x} \in S^n$. Note that $I_{\text{total}}(\mathbf{x}|\theta)$ in (12) does not depend on the arbitrary sufficient statistic
 265 T of (11). Hence

$$I_{\text{total}}(\mathbf{x}|\theta) = I_{\text{comp}}(\mathbf{x}|\theta, L) + I_{\text{lost}}(\mathbf{x}|L) = I_{\text{comp}}(\mathbf{x}|\theta, K) + I_{\text{lost}}(\mathbf{x}|K). \quad (33)$$

266 Then (32) follows immediately from (30) and (33). ■

267 As a consequence of **Theorem 3.5**, **Theorem 3.8** has an immediate corollary.

268 **Corollary 3.9.** Under the conditions of **Theorem 3.8**, let T be a sufficient statistic for θ for
 269 which there is a one-to-one function between T and K . Then for all $\mathbf{x} \in S^n$,

$$I_{\text{lost}}(\mathbf{x}|L) \leq I_{\text{lost}}(\mathbf{x}|T). \quad (34)$$

270 **Corollary 3.9** raises the question whether (34) holds for all sufficient statistics T for θ or even
 271 for all real-valued statistics T . It is conjectured that the first conclusion is false and hence so is the
 272 second, but the question remains open. It is conceivable that notion of a minimal sufficient statistic [3]
 273 is relevant. Regardless, the proofs of **Lemma 3.7** and **Theorem 3.8** illustrate the fact that the relation
 274 between the lost information for two statistics T and T' is determined by the relation between their
 275 partition sets $A_t = \{\mathbf{x}|T(\mathbf{x}) = t\}$ and $B_{t'} = \{\mathbf{x}|T'(\mathbf{x}) = t'\}$. For example, if for every A_t there exists a
 276 $B_{t'}$ for which $A_t \subset B_{t'}$, then the partition of S^n by the $B_{t'}$ of T' is said to be coarser than the
 277 partition by the A_t of T . In that case, $I_{\text{lost}}(\mathbf{x}|\theta, T) \leq I_{\text{lost}}(\mathbf{x}|\theta, T')$ because each $\mathbf{x} \in S^n$ has more $\mathbf{y} \in$
 278 S^n with $T'(\mathbf{y}) = T'(\mathbf{x})$ than there are with $T(\mathbf{y}) = T(\mathbf{x})$. In words, $T'(\mathbf{y}) = t'$ is at least as
 279 ambiguous as $T(\mathbf{y}) = t$ in determining the data sample giving the value of the respective statistics.

281 4. Entropic Loss for an SS

282 For a sufficient statistic T for θ we now propose an entropy measure to characterize T by the
 283 expected lost information incurred by compressing the random sample \mathbf{X} into $T(\mathbf{X})$. This
 284 expectation is taken over all possible data sets \mathbf{x} . This nonstandard entropy measure is called
 285 entropic loss, and it depends on neither a particular data set \mathbf{x} nor the value of θ . Before defining
 286 this measure, we need to determine the appropriate pmf to use in taking an expectation. The
 287 following results are used.

288 **Result 4.1.** Under the assumptions of **Theorem 3.3**, for any data sample let $t = T(\mathbf{x})$ and
 289 consider the partition set A_t . Then

$$\sum_{\mathbf{x} \in A_t} P[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = t] = 1. \quad (35)$$

290 **Proof.** Summing (16) over $\mathbf{x} \in A_t$ yields

$$\sum_{\mathbf{x} \in A_t} P[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = t] = \frac{\sum_{\mathbf{x} \in A_t} h(\mathbf{x})}{\sum_{\mathbf{y} \in A_t} h(\mathbf{y})} = 1. \quad (36)$$

291 to give (35). ■

292 **Result 4.2.** Under the assumptions of **Theorem 3.3**, the sum

$$\sum_{\mathbf{x} \in S^n} P[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})] = |\tau_T|. \quad (37)$$

293 **Proof.** We perform the sum on the left of (37) by first summing over $\mathbf{x} \in A_t$ for fixed t and then
 294 summing over each $t \in \tau_T$ to give

$$\sum_{\mathbf{x} \in S^n} P[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})] = \sum_{t \in \tau_T} \sum_{\mathbf{x} \in A_t} P[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = t], \quad (38)$$

295 The inner series on the right side of (38) sums to one by **Result 4.1**. Hence the outer sum yields
 296 $|\tau_T|$ for $\tau_T = \{t | \exists \mathbf{x} \in S^n \text{ for which } t = T(\mathbf{x})\}$ ■

297 From (37) it follows that the left side of (37) is not a probability distribution on S^n unless $|\tau_T| =$
 298 1. Moreover, $P[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})]$ is not a conditional probability distribution even if $|\tau_T| = 1$
 299 since the condition $T(\mathbf{X}) = T(\mathbf{x})$ varies with \mathbf{x} . However, we use **Result 4.2** to normalize
 300 $P[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})]$ and obtain the appropriate pmf for calculating the expectation of $I_{\text{lost}}(\mathbf{X}|T)$.

301 **Definition 4.3 (Entropic Loss).** Under the assumptions of **Theorem 3.3**, the entropic loss
 302 resulting from the data compression by T is defined as

$$H_{\text{lost}}(\mathbf{X}, T) = \frac{-1}{|\tau_T|} \sum_{\mathbf{x} \in S^n} P[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})] \log P[\mathbf{X} = \mathbf{x} | T(\mathbf{X}) = T(\mathbf{x})], \quad (39)$$

303 which from (15) and (16) can be rewritten as

$$H_{\text{lost}}(\mathbf{X}, T) = \frac{-1}{|\tau_T|} \sum_{\mathbf{x} \in S^n} \frac{h(\mathbf{x})}{\sum_{\mathbf{y} \in A_{T(\mathbf{x})}} h(\mathbf{y})} \log \frac{h(\mathbf{x})}{\sum_{\mathbf{y} \in A_{T(\mathbf{x})}} h(\mathbf{y})}. \quad (40)$$

304 Note that (39) and (40) are independent of both \mathbf{x} and θ . Also, as noted in Section 3 for
 305 $I_{\text{lost}}(\mathbf{x}|T)$, if each $A_{T(\mathbf{x})}$ is a singleton in (40), then $H_{\text{lost}}(\mathbf{X}, T) = 0$. We now compute $H_{\text{lost}}(T)$ for the
 306 sufficient statistic $T(\mathbf{X}) = L(\mathbf{X}|\theta)$.

307 **Theorem 4.4 (Entropic Loss for Likelihood Function).** Under the assumptions of **Theorem 3.3**,
 308 the entropic loss resulting from data compression by $T(\mathbf{x}) = L(\mathbf{x}|\theta)$ is

$$H_{\text{lost}}(\mathbf{X}, L) = \frac{-1}{|\tau_L|} \sum_{t \in \tau_L} \log \frac{1}{|A_t|}. \quad (41)$$

309 **Proof.** From (20) write

$$H_{\text{lost}}(\mathbf{X}, L) = \frac{-1}{|\tau_L|} \sum_{\mathbf{x} \in S^n} \frac{1}{|A_{L(\mathbf{x})}|} \log \frac{1}{|A_{L(\mathbf{x})}|}. \quad (42)$$

310 We decompose the sum over $\mathbf{x} \in S^n$ in (42) to consecutive sums over $\mathbf{x} \in A_t$ and then $t \in \tau_T$ to get

$$H_{\text{lost}}(\mathbf{X}, L) = \frac{-1}{|\tau_L|} \sum_{t \in \tau_L} \sum_{\mathbf{x} \in A_t} \frac{1}{|A_t|} \log \frac{1}{|A_t|} = \frac{-1}{|\tau_L|} \sum_{t \in \tau_L} \frac{|A_t|}{|A_t|} \log \frac{1}{|A_t|}. \quad (43)$$

311 Equation (41) now follows from (43). ■

312 Since $H_{\text{lost}}(\mathbf{X}, T)$ has been defined only for a sufficient statistic T for θ and is independent of
 313 θ , as well as the data sample \mathbf{x} . $H_{\text{lost}}(\mathbf{X}, T)$ could thus be used to compare sufficient statistics. In
 314 particular, if the sufficient statistics $T_1(\mathbf{X})$ and $T_2(\mathbf{X})$ are considered as estimators for θ , then
 315 entropic loss could serve as a metric for regarding, say, T_1 as a better estimator for θ than T_2 if
 316 $H_{\text{lost}}(\mathbf{X}, T_1) < H_{\text{lost}}(\mathbf{X}, T_2)$.

317 **Result 4.5.** If there is a one-to-one function between two sufficient statistics T and T' for θ , then
 318 they have the same entropic loss for a random sample \mathbf{X} ; i.e.,

$$H_{\text{lost}}(\mathbf{X}, T) = H_{\text{lost}}(\mathbf{X}, T'). \quad (44)$$

319 **Proof.** For all $\mathbf{x} \in S^n$, $I_{\text{lost}}(\mathbf{x}|T) = I_{\text{lost}}(\mathbf{x}|T')$ from **Theorem 3.5**, so

$$-\log \frac{h(\mathbf{x})}{\sum_{\mathbf{y} \in A_{T(\mathbf{x})}} h(\mathbf{y})} = -\log \frac{h(\mathbf{x})}{\sum_{\mathbf{y} \in A_{T'(\mathbf{x})}} h(\mathbf{y})} \quad (45)$$

320 from which

$$\frac{h(\mathbf{x})}{\sum_{\mathbf{y} \in A_{T(\mathbf{x})}} h(\mathbf{y})} = \frac{h(\mathbf{x})}{\sum_{\mathbf{y} \in A_{T'(\mathbf{x})}} h(\mathbf{y})} \tag{46}$$

321 Thus from (45) and (46)

$$\frac{h(\mathbf{x})}{\sum_{\mathbf{y} \in A_{T(\mathbf{x})}} h(\mathbf{y})} \log \frac{h(\mathbf{x})}{\sum_{\mathbf{y} \in A_{T(\mathbf{x})}} h(\mathbf{y})} = \frac{h(\mathbf{x})}{\sum_{\mathbf{y} \in A_{T'(\mathbf{x})}} h(\mathbf{y})} \log \frac{h(\mathbf{x})}{\sum_{\mathbf{y} \in A_{T'(\mathbf{x})}} h(\mathbf{y})}. \tag{47}$$

322 Now summing (47) over $\mathbf{x} \in S^n$ yields

$$\sum_{\mathbf{x} \in S^n} \frac{h(\mathbf{x})}{\sum_{\mathbf{y} \in A_{T(\mathbf{x})}} h(\mathbf{y})} \log \frac{h(\mathbf{x})}{\sum_{\mathbf{y} \in A_{T(\mathbf{x})}} h(\mathbf{y})} = \sum_{\mathbf{x} \in S^n} \frac{h(\mathbf{x})}{\sum_{\mathbf{y} \in T'(\mathbf{x})} h(\mathbf{y})} \log \frac{h(\mathbf{x})}{\sum_{\mathbf{y} \in A_{T'(\mathbf{x})}} h(\mathbf{y})}. \tag{48}$$

323 But from **Theorem 3.5**, $|\tau_T| = |\tau_{T'}|$. Thus dividing the left side of (48) by $-|\tau_T|$ and the right side by
 324 $-|\tau_{T'}|$ yields (44). ■

325 Given (32), it might be anticipated that

$$H_{\text{lost}}(\mathbf{X}, L) \leq H_{\text{lost}}(\mathbf{X}, K). \tag{49}$$

326 However, we conjecture that (49) is not always true, but we have no counterexample. If this conjecture
 327 is true, then $L(\mathbf{x}|\theta)$ would not in general have the minimum entropic loss among sufficient statistics
 328 for θ .

329 5. Examples and Computational Issues

330 In this section we present examples involving the discrete Poisson, binomial, and geometric
 331 distributions [9]. For each distribution, three sufficient statistics for some parameter θ are analyzed.
 332 Thus the right side of (8) is independent of θ , as well as the information $I_{\text{lost}}(\mathbf{x}|T)$ conveyed by the
 333 data sample \mathbf{x} about \mathbf{X} . Even for sufficient statistics, calculating the information quantities of this
 334 paper may present computational issues, some of which are discussed in this section. Our examples
 335 are therefore simple in order to focus on the definitions and results of Sections 3 and 4.

336 **Example 5.1 (Poisson Distribution).** Consider the random sample $\mathbf{X} = (X_1, \dots, X_n)$ with the data
 337 sample $\mathbf{x} = (x_1, \dots, x_n)$ from a Poisson random variable X . We consider three sufficient statistics for
 338 the parameter $\theta > 0$. These sufficient statistics are $T_1(\mathbf{X}) = \sum_{i=1}^n X_i$, the likelihood kernel $T_2(\mathbf{X}) =$
 339 $K(\mathbf{X}|\theta)$ for fixed but arbitrary θ and the likelihood function $T_3(\mathbf{X}) =$
 340 $L(\mathbf{X}|\theta)$ for fixed but arbitrary θ . In particular, we use $T_1(\mathbf{X})$ as a surrogate for $T'_1(\mathbf{X}) = \frac{\sum_{i=1}^n X_i}{n}$.
 341 Neither $T_1(\mathbf{X})$ or $T'_1(\mathbf{X})$ involves θ and can thus be used either to characterize \mathbf{X} or to estimate θ .
 342 Moreover, since there is an obvious one-to-one function relating $\frac{\sum_{i=1}^n X_i}{n}$ and $\sum_{i=1}^n X_i$, **Theorems 3.5**
 343 and **4.5** establish that $I_{\text{lost}}(\mathbf{x}|T'_1) = I_{\text{lost}}(\mathbf{x}|T_1)$ and $H_{\text{lost}}(\mathbf{X}, T'_1) = H_{\text{lost}}(\mathbf{X}, T_1)$, respectively. We
 344 consider $T_1(\mathbf{X})$ because it is also Poisson, whereas $T'_1(\mathbf{X})$ is not since $\frac{\sum_{i=1}^n X_i}{n}$ is not necessarily a
 345 nonnegative integer. In contrast to $T_1(\mathbf{X})$, both $T_2(\mathbf{X})$ and $T_3(\mathbf{X})$ contain θ and can only be used to
 346 characterize \mathbf{X} . For each of these three sufficient statistics we develop an expression for $I_{\text{lost}}(\mathbf{x}|T)$
 347 and describe how to obtain a numerical value. We then illustrate previous results with simple data and
 348 present computational results in Table 5.1.

349 **Case 1:** Let $T_1(\mathbf{X}) = \sum_{i=1}^n X_i$. Observe that $T_1(\mathbf{X})$ is a sufficient statistic for θ from **Result 2.2**
 350 since $f(\mathbf{x}|\theta) = P_\theta[\mathbf{X} = \mathbf{x}] = \frac{\theta^{\sum_{i=1}^n x_i} e^{-n\theta}}{\prod_{i=1}^n x_i!}$ can be factored in (2) into the functions $g[T_1(\mathbf{x})|\theta] = \theta^{\sum_{i=1}^n x_i} e^{-n\theta}$
 351 and $h(\mathbf{x}) = \frac{1}{\prod_{i=1}^n x_i!}$. Next recall that the statistic $\sum_{i=1}^n X_i$ has a Poisson distribution with parameter
 352 $n\theta$ [9]. Thus $P_\theta[\sum_{i=1}^n X_i = \sum_{i=1}^n x_i] = \frac{(n\theta)^{\sum_{i=1}^n x_i} e^{-n\theta}}{(\sum_{i=1}^n x_i)!}$, and so (8) becomes

$$P[\mathbf{X} = \mathbf{x} | \sum_{i=1}^n X_i = \sum_{i=1}^n x_i] = \frac{1}{n^{\sum_{i=1}^n x_i}} \binom{\sum_{i=1}^n x_i}{x_1, \dots, x_n}, \tag{50}$$

353 where the multinomial coefficient $\binom{\sum_{i=1}^n x_i}{x_1, \dots, x_n} = \frac{(\sum_{i=1}^n x_i)!}{\prod_{i=1}^n x_i!}$. It follows from (50) and (10) that

$$I_{\text{lost}}(\mathbf{x}|T_1) = -\log \binom{\sum_{i=1}^n x_i}{x_1, \dots, x_n} + (\log n) \sum_{i=1}^n x_i, \tag{51}$$

354 which is also $I_{\text{lost}}(\mathbf{x}|T'_1)$.

355 For a data sample (x_1, \dots, x_n) , the evaluation of $I_{\text{lost}}(\mathbf{x}|T_1)$ in (51) involves computing
 356 factorials. For realistic data, the principal limitation in calculating them by direct multiplication is
 357 their magnitude. See [11] for a discussion. However, (51) can be approximated using either the well-
 358 known Stirling formula or the more accurate Ramanujan approximation [12]. The online multinomial
 359 coefficient calculator [13] can evaluate multinomial coefficients for both x_i and n less than
 360 approximately 50 if any $x_i = 0$ is removed from $\binom{\sum_{i=1}^n x_i}{x_1, \dots, x_n}$. Such deletions do not affect the calculation
 361 since $0! = 1$.

362 As a numerical example, consider a data sample \mathbf{x} of size $n = 34$ from a Poisson random
 363 variable X with $\theta = 3$. On the average, $T_1(\mathbf{X}) = \sum_{i=1}^n X_i = n\theta = 102$, so we take $\sum_{i=1}^n x_i = 102$ for
 364 the data sample $\mathbf{x} = (4, 7, 1, 3, 4, 2, 5, 0, 1, 2, 3, 6, 8, 0, 1, 2, 4, 9, 0, 2, 3, 1, 4, 2, 0, 1, 5, 6, 2, 7, 0, 1, 4, 2)$. Then
 365 the calculator at [13] gives that $\binom{\sum_{i=1}^n x_i}{x_1, \dots, x_n} \approx 1.574 \times 10^{123}$ in (50). Moreover, $(\log n) \sum_{i=1}^n x_i = 518.915$.
 366 Hence from (51), $I_{\text{lost}}(\mathbf{x}|T_1) = I_{\text{lost}}(\mathbf{x}|T'_1) \approx 109.667$ bits of Shannon information. This value
 367 corresponds to 13.708 bytes at 8 bits per byte or to 0.013 kilobytes (KB) at 1024 bytes per kilobyte [14]. It
 368 thus follows from the discussion at the beginning of this example that

$$I_{\text{lost}}(\mathbf{x}|T_1) = I_{\text{lost}}(\mathbf{x}|T'_1) \approx 0.013 \text{ KB}. \tag{52}$$

369 **Case 2:** Let $T_2(\mathbf{X}) = K(\mathbf{X}|\theta)$ for fixed but arbitrary $\theta > 0$. For a data sample (x_1, \dots, x_n) write

$$L(\mathbf{x}|\theta) = f(\mathbf{x}|\theta) = \frac{\theta^{\sum_{i=1}^n x_i} e^{-n\theta}}{\prod_{i=1}^n x_i!}, \tag{53}$$

370 from which

$$K(\mathbf{x}|\theta) = \theta^{\sum_{i=1}^n x_i} e^{-n\theta} \tag{54}$$

371 and $R(\mathbf{x}) = \frac{1}{\prod_{i=1}^n x_i!}$ in (4). Note that for all fixed $\theta > 0$ except $\theta = 1$, there is an obvious one-to-one
 372 function between $T_1(\mathbf{x}) = \sum_{i=1}^n x_i$ and (54). Hence in the numerical example of **Case 1**,
 373 $I_{\text{lost}}(\mathbf{x}|K(\mathbf{x}|\theta)) = I_{\text{lost}}(\mathbf{x}|T_1) \approx 0.013$ KB from **Theorem 3.5** for all $\theta > 0$ except $\theta = 1$. For $\theta = 1$,
 374 $K(\mathbf{x}|\theta) = e^{-n}$ and is constant with respect to any data sample \mathbf{x} . Thus $I_{\text{comp}}(\mathbf{x}|1, K) = 0$ and
 375 $I_{\text{lost}}(\mathbf{x}|K(\mathbf{x}|1)) = I_{\text{total}}(\mathbf{x}|1, K)$. It follows that $K(\mathbf{x}|1)$ provides no information about \mathbf{X} .

376 **Case 3:** Let $T_3(\mathbf{X}) = L(\mathbf{X}|\theta)$ for fixed but arbitrary $\theta > 0$. We attempt to obtain $I_{\text{lost}}(\mathbf{x}|L(\mathbf{x}|\theta))$
 377 for a data sample $\mathbf{x} = (x_1, \dots, x_n)$ by determining $|A_{L(\mathbf{x}|\theta)}|$ and using (20). From (53), note that for all
 378 fixed $\theta > 0$ except $\theta = 1$, $\mathbf{y} \in A_{L(\mathbf{x}|\theta)}$ if and only if

$$\frac{\theta^{\sum_{i=1}^n y_i}}{\prod_{i=1}^n y_i!} = \frac{\theta^{\sum_{i=1}^n x_i}}{\prod_{i=1}^n x_i!}. \tag{55}$$

379 Thus for any fixed θ satisfying $\theta > 0$ and $\theta \neq 1$, $\mathbf{y} \in A_{L(\mathbf{x}|\theta)}$ if both $\sum_{i=1}^n y_i = \sum_{i=1}^n x_i$ and
 380 $\prod_{i=1}^n y_i! = \prod_{i=1}^n x_i!$. However, for some $\theta > 0$ and $\theta \neq 1$, it is possible that $\mathbf{y} \in A_{L(\mathbf{x}|\theta)}$ when neither
 381 $\sum_{i=1}^n y_i = \sum_{i=1}^n x_i$ nor $\prod_{i=1}^n y_i! = \prod_{i=1}^n x_i!$. For example, let $\theta = 2$, $\mathbf{x} = (4, 1, 1, 0)$, and $\mathbf{y} = (3, 2, 0, 0)$.
 382 Then $\sum_{i=1}^n x_i = 6$, $\sum_{i=1}^n y_i = 5$, $\prod_{i=1}^n x_i! = 24$, and $\prod_{i=1}^n y_i! = 12$. However, (55) is satisfied.

383 Such complications suggest that an efficient implicit enumeration of the \mathbf{y} satisfying (55) would
 384 be required to obtain $|A_{L(\mathbf{x}|\theta)}|$ for calculating $I_{\text{lost}}(\mathbf{x}|L(\mathbf{x}|\theta))$ from (20). Using such an algorithm, a
 385 conventional computer could probably compute $I_{\text{lost}}(\mathbf{x}|L(\mathbf{x}|\theta))$ for the numerical data and value of
 386 θ in **Case 1** since there is now a 250 petabyte, 200 petaflop conventional computer [15]. Substantially
 387 larger problems, if not already tractable, will likely be so in the foreseeable future on quantum
 388 computers. Recently the milestone of quantum supremacy was achieved where the various possible

389 combinations of a certain randomly generated output were obtained in 110 seconds, whereas this task
 390 would have taken the above conventional supercomputer 10,000 years [16]. Regardless, for the data
 391 of **Case 1**, we have the upper bound $I_{\text{lost}}(\mathbf{x}|L(\mathbf{x}|\theta)) \leq 0.013$ KB from (32).

392 We present some simple further simple computational results for the Poisson example distribution
 393 to illustrate relationships between among T_1, T_2, T_3 . Table 5.1 below summarizes the results for sample
 394 data (x_1, x_2, x_3) with $\sum_{i=1}^3 x_i \leq 2$. In particular, a complete enumeration of $A_{L(\mathbf{x}|\theta)}$ gives
 395 $I_{\text{lost}}(\mathbf{x}|L(\mathbf{x}|\theta))$ from (20).

397 Table 5.1. Poisson Example
 398

$\mathbf{x} = (x_1, x_2, x_3)$	$T_1(\mathbf{x})$	$I_{\text{lost}}(\mathbf{x} T_1)$	$T_2(\mathbf{x})$	$I_{\text{lost}}(\mathbf{x} T_2)$	$T_3(\mathbf{x})$	$I_{\text{lost}}(\mathbf{x} T_3)$
(0,0,0)	0	0	$e^{-3\theta}$	0	$e^{-3\theta}$	0
(0,0,1)	1	$\log 3$	$\theta e^{-3\theta}$	$\log 3$	$\theta e^{-3\theta}$	$\log 3$
(0,1,0)						
(1,0,0)						
(1,1,0)	2	$\log \frac{9}{2}$	$\theta^2 e^{-3\theta}$	$\log \frac{9}{2}$	$\theta^2 e^{-3\theta}$	$\log 3$
(1,0,1)						
(0,1,1)						
(2,0,0)	2	$\log 9$	$\theta^2 e^{-3\theta}$	$\log 9$	$\frac{\theta^2 e^{-3\theta}}{2}$	$\log 3$
(0,2,0)						
(0,0,2)						

399 **Example 5.2 (Binomial Distribution).** Consider a random sample $\mathbf{X} = (X_1, \dots, X_n)$ from a
 400 binomial random variable X with parameters m and θ , where θ is the probability of success on
 401 any of the m Bernoulli trials associated with the $X_i, i = 1, \dots, n$. Let m be fixed, so the only
 402 parameter is θ . Moreover, the sample space of the underlying random variable X is now finite.

403 **Case 1:** $T_1(\mathbf{X}) = \sum_{i=1}^n X_i$. Again $\sum_{i=1}^n X_i$ is an SS for θ . From [9], $\sum_{i=1}^n X_i$ has a binomial
 404 distribution with parameter θ for fixed nm . Hence

$$P_\theta \left[\sum_{i=1}^n X_i = \sum_{i=1}^n x_i \right] = \theta^{\sum_{i=1}^n x_i} \theta^{mn - \sum_{i=1}^n x_i} \binom{mn}{\sum_{i=1}^n x_i} \tag{56}$$

406 and

$$P_\theta[\mathbf{X} = \mathbf{x}] = \theta^{\sum_{i=1}^n x_i} \theta^{mn - \sum_{i=1}^n x_i} \prod_{i=1}^n \binom{m}{x_i}. \tag{57}$$

407 From (1), dividing (57) by (56) gives

$$P[\mathbf{X} = \mathbf{x} | \sum_{i=1}^n X_i = t] = \frac{\prod_{i=1}^n \binom{m}{x_i}}{\binom{mn}{t}}. \tag{58}$$

408 By taking the $-\log$ of (58) gives the lost information as

$$I_{\text{lost}}(\mathbf{x}|T_1) = -\log \frac{\prod_{i=1}^n \binom{m}{x_i}}{\binom{mn}{t}} = -\sum_{i=1}^n \log \binom{m}{x_i} + \log \binom{mn}{t}. \tag{59}$$

409 **Case 2:** $T_2(\mathbf{X}) = K(\mathbf{X}|\theta)$. In this case we use (16) as in **Example 5.1**. Write

$$L(\mathbf{x}|\theta) = f(\mathbf{x}|\theta) = \theta^{\sum_{i=1}^n x_i} (1 - \theta)^{mn - \sum_{i=1}^n x_i} \prod_{i=1}^n \binom{m}{x_i}, \tag{60}$$

410 from which $K(\mathbf{x}|\theta) = \theta^{\sum_{i=1}^n x_i} (1 - \theta)^{mn - \sum_{i=1}^n x_i}$ and $R(\mathbf{x}) = \prod_{i=1}^n \binom{m}{x_i}$ in (4). To factor the right side of
 411 (60) as in (2), let g be the identity function and $h(\mathbf{x}) = \prod_{i=1}^n \binom{m}{x_i}$. Hence,

$$I_{\text{lost}}(\mathbf{x}|T_2) = -\log \frac{\prod_{i=1}^n \binom{m}{x_i}}{\sum_{\mathbf{y} \in A_{K(\mathbf{x}|\theta)}} \prod_{i=1}^n \binom{m}{y_i}}, \tag{61}$$

412 and (61) yields

$$I_{\text{lost}}(\mathbf{x}|T_2) = -\sum_{i=1}^n \log \binom{m}{x_i} + \log \sum_{\mathbf{y} \in A_{K(\mathbf{x}|\theta)}} \prod_{i=1}^n \binom{m}{y_i}, \tag{62}$$

413 where

$$A_{K(\mathbf{x}|\theta)} = \{ \mathbf{y} \in S^n | \theta^{\sum_{i=1}^n y_i} (1-\theta)^{mn-\sum_{i=1}^n y_i} = \theta^{\sum_{i=1}^n x_i} (1-\theta)^{mn-\sum_{i=1}^n x_i} \}. \tag{63}$$

414 From (63), for any fixed θ satisfying $0 < \theta < 1$ and $\theta \neq 1/2$, it can easily be shown that $\mathbf{y} \in$
 415 $A_{K(\mathbf{x}|\theta)}$ if and only if $\sum_{i=1}^n y_i = \sum_{i=1}^n x_i$. Thus in general, for a given \mathbf{x} and fixed θ , determining
 416 $A_{K(\mathbf{x}|\theta)}$ in **Case 2** would require an enumeration of the \mathbf{y} satisfying (63) to compute (62). We perform
 417 such an enumeration below for a simple example.

418 **Case 3:** $T_3(\mathbf{X}) = L(\mathbf{X}|\theta)$. For a data sample $\mathbf{x} = (x_1, \dots, x_n)$ we now have

$$L(\mathbf{x}|\theta) = \left(\frac{\theta}{1-\theta}\right)^{\sum_{i=1}^n x_i} (1-\theta)^{mn} \prod_{i=1}^n \binom{m}{x_i} \tag{64}$$

419 with g be the identity function and $h(\mathbf{x}) = 1$ in (2). For fixed θ satisfying $0 < \theta < 1$ and $\theta \neq 1/2$,
 420 we obtain that $\mathbf{y} \in A_{L(\mathbf{x}|\theta)}$ if and only if

$$\left(\frac{\theta}{1-\theta}\right)^{\sum_{i=1}^n y_i} \prod_{i=1}^n \binom{m}{y_i} = \left(\frac{\theta}{1-\theta}\right)^{\sum_{i=1}^n x_i} \prod_{i=1}^n \binom{m}{x_i}. \tag{65}$$

421 As in **Case 3** of **Example 5.1**, developing an algorithm to use (65) and determine $|A_{L(\mathbf{x}|\theta)}|$ for
 422 calculating $I_{\text{lost}}(\mathbf{x}|L(\mathbf{x}|\theta))$ from (20) is beyond the scope of this paper.

423 As a simple example, consider the experiment of flipping a possibly biased coin twice ($m = 2$). The
 424 total number of heads follows a binomial distribution with the parameter θ , which is the probability
 425 of getting a head on any flip. By doing this experiment three times we generate the random variables
 426 X_1, X_2, X_3 with possible values 0, 1, 2. Table 5.2 shows all the possibilities and the lost information for
 427 the statistics. The small size of this example allows the computation of I_{lost} in **Cases 2** and **3** via total
 428 enumeration.

429
 430

Table 5.2. Binomial Example

$\mathbf{x} = (x_1, x_2, x_3)$	$T_1(\mathbf{x})$	$I_{\text{lost}}(\mathbf{x} T_1)$	$T_2(\mathbf{x})$	$I_{\text{lost}}(\mathbf{x} T_2)$	$T_3(\mathbf{x})$	$I_{\text{lost}}(\mathbf{x} T_3)$
(0,0,0)	0	0	$(1-\theta)^6$	0	$(1-\theta)^6$	0
(0,0,1)	1	$\log 3$	$(1-\theta)^5\theta^1$	$\log 3$	$2(1-\theta)^5\theta^1$	$\log 3$
(0,1,0)						
(1,0,0)						
(1,1,0)	2	$\log \frac{15}{4}$	$(1-\theta)^4\theta^2$	$\log \frac{15}{4}$	$4(1-\theta)^4\theta^2$	$\log 3$
(1,0,1)						
(0,1,1)						
(2,0,0)	2	$\log 15$	$(1-\theta)^4\theta^2$	$\log 15$	$(1-\theta)^4\theta^2$	$\log 3$
(0,2,0)						
(0,0,2)						
(1,1,1)	3	$\log \frac{5}{2}$	$(1-\theta)^3\theta^3$	$\log \frac{5}{2}$	$8(1-\theta)^3\theta^3$	0
(2,1,0)						

(2,0,1)	3	log 10	$(1 - \theta)^3 \theta^3$	log 10	$2(1 - \theta)^3 \theta^3$	log 6
(1,0,2)						
(1,2,0)						
(0,1,2)						
(0,2,1)						
(2,1,1)	4	$\log \frac{15}{4}$	$(1 - \theta)^2 \theta^4$	$\log \frac{15}{4}$	$4(1 - \theta)^2 \theta^4$	log 3
(1,2,1)						
(1,1,2)						
(2,2,0)	4	log 15	$(1 - \theta)^2 \theta^4$	log 15	$(1 - \theta)^2 \theta^4$	log 3
(2,0,2)						
(0,2,2)						
(2,2,1)	5	log 3	$(1 - \theta)^1 \theta^5$	log 3	$2(1 - \theta)^1 \theta^5$	log 3
(2,1,2)						
(1,2,2)						
(2,2,2)	6	0	θ^6	0	θ^6	0

431
 432 Now using (40), we give in Table 5.3 the entropic losses of **Example 5.2** for T_1, T_2, T_3 . Note that
 433 $H_{\text{lost}}(\mathbf{X}, T)$ is the same for the sum T_1 and the likelihood kernel T_2 , which are related by a one-to-
 434 one function. Hence **Result 4.5** is corroborated. Also observe that $H_{\text{lost}}(\mathbf{X}, T)$ is smallest for the
 435 likelihood function T_3 .

437 Table 5.3. Entropic loss over different statistics for a binomial distribution

$H_{\text{lost}}(\mathbf{X}, T_1)$	$H_{\text{lost}}(\mathbf{X}, T_2)$	$H_{\text{lost}}(\mathbf{X}, T_3)$
1.4722	1.4722	1.2095

438
 439
 440
 441
 442
 443 **Example 5.3 (Geometric Distribution).** Consider a random sample $\mathbf{X} = (X_1, \dots, X_n)$ with sample
 444 data $\mathbf{x} = (x_1, \dots, x_n)$ from a geometric random variable X , where the parameter θ is the probability
 445 of success on any of the series of independent Bernoulli trials for which X is the trial number on
 446 which the first success is obtained. It readily follows from [5] that

$$P[\mathbf{X} = \mathbf{x}] = \theta^n (1 - \theta)^{\sum_{i=1}^n x_i - n}. \tag{66}$$

447 **Case 1:** $T_1(\mathbf{X}) = \sum_{i=1}^n X_i$. For fixed n , $\sum_{i=1}^n X_i$ has a negative binomial distribution with
 448 parameter θ . Hence,
 449

$$P\left[\sum_{i=1}^n X_i = \sum_{i=1}^n x_i\right] = \binom{\sum_{i=1}^n x_i - 1}{n - 1} \theta^n (1 - \theta)^{\sum_{i=1}^n x_i - n}. \tag{67}$$

450 Thus $T_1(\mathbf{X}) = \sum_{i=1}^n X_i$ is an SS for θ since it satisfies (2) with $g[T_1(\mathbf{x})|\theta] = \theta^n (1 - \theta)^{T_1(\mathbf{x}) - n}$ and
 451 $h(x_1, \dots, x_n) = \binom{\sum_{i=1}^n x_i - 1}{n - 1}$. Moreover, substitution of (66) and (67) into (8) gives

$$P\left[\mathbf{X} = \mathbf{x} \mid \sum_{i=1}^n X_i = \sum_{i=1}^n x_i\right] = \frac{1}{\binom{\sum_{i=1}^n x_i - 1}{n - 1}}. \tag{68}$$

452 Then from (14) and (68) we obtain that

$$I_{\text{lost}}(\mathbf{x}|T_1) = \log\left(\binom{\sum_{i=1}^n x_i - 1}{n - 1}\right). \tag{69}$$

453 **Case 2:** $T_2(\mathbf{X}) = K(\mathbf{X}|\theta)$. From (66), for all $\mathbf{x} \in S^n$, $R(\mathbf{x}) = 1$ and

$$K(\mathbf{x}|\theta) = L(\mathbf{x}|\theta) = \left(\frac{\theta}{1 - \theta}\right)^n (1 - \theta)^{\sum_{i=1}^n x_i}. \tag{70}$$

454 Thus for $0 < \theta < 1$, there is an obvious one-to-one function between $T_1(\mathbf{x}) = \sum_{i=1}^n x_i$ and $T_2(\mathbf{x}) =$
 455 $K(\mathbf{x}|\theta)$ in (70). Thus from **Theorem 3.5**, $I_{\text{lost}}(\mathbf{x}|T_2(\mathbf{x})) = I_{\text{lost}}(\mathbf{x}|T_1)$ as given in (69).

456 **Case 3:** $T_3(\mathbf{X}) = L(\mathbf{X}|\theta)$. Since $K(\mathbf{X}|\theta) = L(\mathbf{X}|\theta)$ from (70), then

$$I_{\text{lost}}(\mathbf{x}|T_3) = \log \left(\frac{\sum_{i=1}^n x_i - 1}{n - 1} \right) \tag{71}$$

457 from (69). However, there is an alternate derivation of (71). For $0 < \theta < 1$ it follows from (70) that then
 458 $\mathbf{y} \in A_{L(\mathbf{x}|\theta)}$ if and only if

$$\sum_{i=1}^n y_i = \sum_{i=1}^n x_i. \tag{72}$$

459 But for fixed positive integers x_1, \dots, x_n we have from [17] that the number of solutions $|A_{L(\mathbf{x}|\theta)}|$ to
 460 (72) in positive integers y_1, \dots, y_n is

$$\binom{\sum_{i=1}^n x_i - 1}{n - 1}. \tag{73}$$

461 Thus (71) follows for $L(\mathbf{X}|\theta)$ from (73) and (20), so $I_{\text{lost}}(\mathbf{x}|T_1) = I_{\text{lost}}(\mathbf{x}|T_2) = I_{\text{lost}}(\mathbf{x}|T_3)$ from
 462 **Theorem 3.5**.

463 As a numerical illustration, let the random variable X denote the number of flips of a possibly
 464 biased coin until a head is obtained. Then X has a geometric distribution with the parameter θ as the
 465 probability of getting a head on any flip. Suppose this experiment is performed three times yielding the
 466 possible sample data $\mathbf{x} = (x_1, x_2, x_3)$ shown in Table 5.4. $I_{\text{lost}}(\mathbf{x}|T)$ is then calculated for each of the
 467 sufficient statistics for θ of **Example 5.3**. Observe that the individual statistics depend on θ while the
 468 lost information does not. Moreover, $I_{\text{lost}}(\mathbf{x}|T_1) = I_{\text{lost}}(\mathbf{x}|T_2) = I_{\text{lost}}(\mathbf{x}|T_3)$ for all the sample data as
 469 established analytically above.

Table 5.4. Geometric Example

$\mathbf{x} = (x_1, x_2, x_3)$	$T_1(\mathbf{x})$	$I_{\text{lost}}(\mathbf{x} T_1)$	$T_2(\mathbf{x})$	$I_{\text{lost}}(\mathbf{x} T_2)$	$T_3(\mathbf{x})$	$I_{\text{lost}}(\mathbf{x} T_3)$
(1,1,1)	3	0	θ^3	0	θ^3	0
(2,1,1)	4	$\log 3$	$\theta^3(1 - \theta)$	$\log 3$	$\theta^3(1 - \theta)$	$\log 3$
(1,2,1)						
(1,1,2)						
(2,2,1)	5	$\log 6$	$\theta^3(1 - \theta)^2$	$\log 6$	$\theta^3(1 - \theta)^2$	$\log 6$
(2,1,2)						
(1,2,2)						
(2,2,2)	6	$\log 10$	$\theta^3(1 - \theta)^3$	$\log 10$	$\theta^3(1 - \theta)^3$	$\log 10$

472 **6. Conclusion**

473 In this paper, the Shannon information obtained from a random sample \mathbf{X} for a discrete random
 474 variable X with a single parameter θ was decomposed into two components: (i) the compressed
 475 information obtained by the value of a real-valued statistic $T(\mathbf{X})$ for the sample data \mathbf{x} and (ii) the
 476 information lost by using this statistic to characterize \mathbf{X} . We focused on this lost information caused
 477 by multiple data sets having the same value of the statistic. This possibility is typical of data analysis,
 478 where the data uniquely determines the value of the statistic, but a value of the statistic does not
 479 uniquely determine the data yielding it. In other words, we answered the question: how much
 480 Shannon information is lost about a data sample when only the value of a sufficient statistic is known
 481 but not the original data. We also defined the entropic loss associated with a sufficient statistic
 482 T under consideration as the expected lost information over all possible samples to give a metric

483 dependent only on T . Our approach is applicable to any T , but we focused on sufficient statistics for
484 θ for simplicity. Applications of our results were computationally intensive.
485

486 References

- 487 1. Landauer, R. Irreversibility and heat generation in the computing process. *IBM Journal of*
488 *Research and Development*. 1961, 5, 183–191.
- 489 2. Hodge, S.E.; Vieland, V.J. Information loss in binomial data due to data compression. *Entropy*
490 2017, 19, 75–81.
- 491 3. Casella, G.; Berger, R.L. *Statistical Inference*, 2nd ed.; Cengage Learning, Delhi, India, 2002.
- 492 4. Pawitan, Y. In *All Likelihood: Statistical Modeling and Inference Using Likelihood*, 1st ed.; The
493 Clarendon Press: Oxford, UK, 2013.
- 494 5. Rohatgi, V.K.; Ehsanes Saleh, A.K. *An Introduction to Probability and Statistics*, 2nd ed.; John
495 Wiley & Sons, Inc., NY, USA, 2001.
- 496 6. Shannon, C.E.; Weaver, W. *The Mathematical Theory of Communication*, 1st ed.; The University
497 of Illinois Press, Urbana, Illinois, 1964.
- 498 7. Shannon, C. A mathematical theory of communication. *Bell. Syst. Tech. J.* 1948, 27, 379–423.
- 499 8. Kapur, J.N.; Kesavan, H.K. *Entropy Optimization Principles with Applications*, 1st ed.;
500 Academic Press, Inc., San Diego, CA, USA, 1992.
- 501 9. Johnson, J.L. *Probability and Statistics for Computer Science*, 1st ed.; John Wiley & Sons, Inc., NJ,
502 USA, 2003.
- 503 10. Beeler, R.A. *How to Count: An Introduction to Combinatorics and Its Applications*, 1st ed.;
504 Springer, Switzerland, 2015.
- 505 11. <https://en.wikipedia.org/wiki/Factorial#Computation>, accessed 7/1/2019.
- 506 12. Mortici, C. Ramanujan formula for the generalized Stirling approximation. *Applied*
507 *Mathematics and Computation* 2010, 19, 2579–2585.
- 508 13. <https://mathcracker.com/multinomial-coefficient-calculator.php>, accessed 7/27/2019.
- 509 14. <https://en.wikipedia.org/wiki/Kilobyte>, accessed 7/27/2019.
- 510 15. [https://en.wikipedia.org/wiki/Summit_\(supercomputer\)](https://en.wikipedia.org/wiki/Summit_(supercomputer)), accessed 7/1/2019.
- 511 16. Arute, F.; Arya, K.; Martinis, J.M. Quantum supremacy using a programmable
512 superconducting processor. *Nature* 2019, 574, 505–510.
- 513 17. Mahmoudvand, R.; Hassani, H.; Farzaneh, A.; Howell, G. The exact number of nonnegative
514 integer solutions for a linear Diophantine inequality. *IAENG International Journal of Applied*
515 *Mathematics* 2010, 40, 5 pages.



© 2019 by the authors. Submitted for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).