

Data-derived weak universal consistency: the case of universal compression

Narayana Santhanam Venkat Anantharam Wojciech Szpankowski

Abstract

Many current applications in data science need rich model classes to adequately represent the statistics that may be driving the observations. But rich model classes may be too complex to admit estimators that converge to the truth with convergence rates that can be uniformly bounded over the entire collection of probability distributions comprising the model class, i.e. it may be impossible to guarantee uniform consistency of such estimators as the sample size increases. In such cases, it is conventional to settle for estimators with guarantees on convergence rate where the performance can be bounded in a model-dependent way, i.e. pointwise consistent estimators. But this viewpoint has the serious drawback that estimator performance is a function of the unknown model within the model class that is being estimated, and is therefore unknown. Even if an estimator is consistent, how well it is doing at any given time may not be clear, no matter what the sample size of the observations.

Departing from the classical uniform/pointwise consistency dichotomy that leads to this impasse, a new analysis framework is explored by studying rich model classes that may only admit pointwise consistency guarantees, yet all the information about the unknown model driving the observations that is needed to gauge estimator accuracy can be inferred from the sample at hand. We expect that this *data-derived* estimation framework will be broadly applicable to a wide range of estimation problems by providing a methodology to deal with much richer model classes. In this paper we analyze the lossless compression problem in detail in this novel data-derived framework.

I. INTRODUCTION AND MOTIVATION

Many of the most challenging problems in the data sciences stem from one or more of the following characteristics associated with data: extreme scale (typically requiring that the data reside on multiple storage nodes); high dimensionality and sparsity; patterns in the data that manifest at multiple scales; dynamic, temporal, and heterogeneous structure; complex dependencies between different parts of the data; and noise/ missing data. Tasks such as image recognition, classification, control and many others, which are built on such data sources, depend on estimating the relevant underlying structure in the data. Rich model classes, i.e. rich collections of probabilistic models, such as the collection of all probability distributions over a large or countably infinite support, or the set of long memory, slowly mixing Markov processes are often required to adequately model the complex characteristics of these data sources. A comprehensive approach to address these key challenges is critically needed.

N. Santhanam is with the Department of Electrical Engineering, University of Hawaii, Manoa. Email: nsanathan@hawaii.edu.

V. Anantharam is with the Department of Electrical Engineering and Computer Science, University of California, Berkeley. Email: ananth@eecs.berkeley.edu.

W. Szpankowski is with the Department of Computer Science, Purdue University. Email: szpan@purdue.edu.

33 Indeed, in bringing rigorous theory to bear on data science, the first question we face is related to
 34 model selection. There is often a tension between the need for rich model classes to better represent
 35 data and our ability to handle these collections from a mathematical point of view. Many applications,
 36 particularly in the big data regime, force us to consider model collections that are too complex to admit
 37 estimators with traditional model-agnostic *uniformly* consistent guarantees. These new collections often
 38 only admit *pointwise* convergent estimators [1] – i.e. convergence is only guaranteed individually for
 39 each model in the model class – which often are difficult to use predictively as their convergence cannot
 40 be verified. In this paper we depart from this dichotomy, and we propose a new analysis framework
 41 by characterizing rich model classes that may only admit pointwise guarantees, yet all the information
 42 about the unknown model needed to gauge estimator accuracy can be inferred from the observations.
 43 More precisely, we introduce here *data-derived* consistency, a new framework to analyze these rich model
 44 collections. To retain focus, in this document we concentrate on universal compression to bring out the
 45 salient features of this framework. We also make connections to a related prediction problem that was
 46 analyzed by us earlier in [2], and is now seen to fit into this broader framework.

47 The richness of a model class is often quantified by metrics such as its VC-dimension [3], Rademacher
 48 complexity [4], [5], [6], or – what is most relevant in the context of universal compression – its asymptotic
 49 per-symbol redundancy [7], [8], [9], [10], [11], [12]. Ideally, one would want an estimation algorithm with
 50 a model-agnostic guarantee on its performance, depending only on the sample size – this is the *uniform*
 51 consistency dogma that underlies most formulations of engineering applications today. But requiring
 52 such uniform consistency restricts the richness of the model classes we can deal with. Generally speaking,
 53 the more complex a model is, the less one could expect to be able to provide such uniform consistency
 54 guarantees.

55 When the model classes we are interested in are too complex to admit *uniformly consistent estimators*,
 56 the common belief is that the best we can do is to have estimators with convergence guarantees
 57 dependent on not just the sample size but on the underlying model in the model class that governs
 58 the statistics of the observations. These are called *pointwise consistent* estimators. It is well-understood
 59 that this viewpoint may not always be particularly useful, the problem now being that our gauge of
 60 the performance of the estimation scheme is dependent on the unknown underlying model – the very
 61 ambiguity we are addressing! Even if we have a pointwise consistent estimator, which is eventually
 62 almost surely accurate under the underlying model, for any fixed sample size we may never know how
 63 well the estimator is doing no matter how large the sample size is.

64 We illustrate this issue with a simple example below. Before doing so, we first introduce some of the
 65 notational conventions that will be used throughout this document. The symbol $:=$, and occasionally
 66 \equiv , is used to denote equality by definition. We write \log for logarithms to base 2 and \ln for logarithms
 67 to the natural base. The set of natural numbers, denoted \mathbb{N} , is the set $\{1, 2, \dots\}$, thought of as endowed
 68 with its usual σ -algebra comprised of all subsets of \mathbb{N} . For $n \geq 1$, we use \mathbb{N}^n to denote the set of strings
 69 of length n of natural numbers, with the product σ -algebra. The set of infinite sequences of natural
 70 numbers is denoted \mathbb{N}^∞ , and is thought of as endowed with the corresponding product σ -algebra. We
 71 will adopt the convention of thinking of a probability measure on \mathbb{N} as defined by a distribution, which
 72 assigns a probability to each natural number. A string of integers $(x_1, \dots, x_n) \in \mathbb{N}^n$ will be denoted
 73 by \mathbf{x} , or by x^n when it seems important to emphasize the specific length of the string. For a string
 74 of integers $\mathbf{x} := (x_1, \dots, x_n) \in \mathbb{N}^n$, its empirical distribution or *type* is the sequence of unnormalized

75 fractions on \mathbb{N} assigning the fraction $\frac{m}{n}$ to $x \in \mathbb{N}$ if x shows up m times in the string \mathbf{x} . It is conventional
 76 to think of the type as a probability distribution on \mathbb{N} and we will do so when convenient, but it is
 77 important at some places in the document to think of it as comprised of unnormalized fractions. \mathbb{N}^*
 78 denotes the set of strings of naturals of finite length, including the empty string. For the purposes of
 79 this paper it suffices to think of \mathbb{N}^* as a set with no additional structure.

80 **Example 1. (Hiding entropy)**

81 For $\epsilon > 0$ and $M \in \mathbb{N}$, let $p_{\epsilon, M}$ be the probability distribution that assigns probability $1 - \epsilon$ to the
 82 natural number 1 and assigns probability ϵ/M to the natural numbers 2 through $M + 1$. Denote the
 83 probability distribution that assigns probability 1 to the natural number 1 by p_0 . Let \mathcal{W} be the set
 84 comprised of the probability distributions $p_{\epsilon, M}$ for $\epsilon > 0$ and $M \in \mathbb{N}$, as well as p_0 .

85 Our task is to estimate the Shannon entropy of a probability distribution in \mathcal{W} using *i.i.d.* samples
 86 from it. However, we do not know which probability distribution in \mathcal{W} is governing the law of the
 87 observed samples. The natural *plug-in estimator* assigns to a sample X_1, \dots, X_n the entropy of its em-
 88 pirical distribution. Since every probability distribution in \mathcal{W} has finite support, the plug-in estimate is
 89 consistent almost surely, no matter which underlying distribution from \mathcal{W} is generating the observations.
 90 But at what point do we know that the plug-in estimate is close to the correct answer? Indeed, can we,
 91 at any point, get an upper bound for the true entropy using the plug-in estimate with, say, a confidence
 92 probability $3/4$, regardless of what the true probability distribution in \mathcal{W} is?

93 It turns out that it is *impossible* to provide such guarantees for \mathcal{W} . To see why, suppose we see a
 94 sequence of n successive 1s. This could have come from p_0 , or, with high probability, from any probability
 95 distribution $p_{\epsilon, M}$ with $0 < \epsilon \ll \frac{1}{n}$. What is worse, for any upper bound \hat{h} we may provide, however
 96 large, even if $0 < \epsilon \ll \frac{1}{n}$, the entropy of $p_{\epsilon, M}$ where $M \geq 2^{\hat{h}/\epsilon}$ is $h(\epsilon) + \epsilon \log M \geq \hat{h}$. Every such $p_{\epsilon, M}$
 97 gives the sample of n successive 1s a probability of at least $> 3/4$ if ϵ is sufficiently small, so our upper
 98 bound fails.

99 This argument applies whether we obtained \hat{h} from the plug-in estimator or *any* other estimator of
 100 the entropy. No upper bound that we propose on the entropy based on any finite sequence of 1s can
 101 hold with confidence probability $3/4$ under all probability distributions in \mathcal{W} . To make matters worse,
 102 the sequence of all 1s occurs with probability 1 when the underlying model in force is p_0 . Therefore,
 103 even when we could estimate the entropy consistently, we could never obtain even a trivial upper bound
 104 on the entropy with a confidence probability $\geq 3/4$. \square

105 We therefore challenge the dichotomy of *uniform* and *pointwise* consistency in the analysis of statis-
 106 tical estimators. This paper considers a new paradigm positioned in between these two extremes. We
 107 modify the definition of pointwise consistent estimators, keeping as far as possible the richness of the
 108 model class but ensuring that all the information needed about the unknown model to evaluate estimator
 109 accuracy can be gleaned from the observations. We call this modified notion of pointwise consistency
 110 *data-derived* consistency. The crux of the *data-derived* framework is to provide a mechanism that allows
 111 us to gauge from the observations how well we are doing.

112 We bring out the salient features of the data-derived framework in this document in the framework of
 113 universal compression. In the context of providing efficient compressed representations of samples from
 114 a data source, the goal of universal compression is to be able to work with a rich class of models for the
 115 source being compressed. Universal compression posits that we have a model class of source probability

116 measures, while we are required to come up with a universal probability measure that attempts to
 117 compress any source in the model class as well as possible without prior knowledge of the source. Since
 118 the universal probability measure is not exactly matched to any single source probability measure in
 119 the model class it incurs a redundancy, measured using the Kullback-Leibler (KL) divergence, against
 120 any source in the model class when compressing a sequence of observed samples whose statistics are
 121 governed by this source. The uniform consistency setup in this case corresponds to what is commonly
 122 known as the *strong* compression formulation, where we find universal probability measures whose
 123 per-symbol redundancy incurred against any source in the model class can be uniformly bounded
 124 over the entire model class and, in addition, diminishes to 0 as the sample size grows to infinity.
 125 The pointwise consistency setup in this case corresponds to what is commonly known as the *weak*
 126 compression formulation and is one where the universal probability measure incurs asymptotically zero
 127 per-symbol redundancy against each source in the model class, but the convergence to zero is not
 128 necessarily uniform over the entire model class.

129 We propose and study the *data-derived* weak compression formulation (*d.w.c.*) which identifies when,
 130 in the weak compression setup, we can also estimate the redundancy of the universal probability measure
 131 relative to the underlying source model generating the data. Broadly speaking, we aim to find a universal
 132 estimator/encoding with a given accuracy as well as a corresponding stopping rule that allows us to
 133 find out at what point the KL divergence from the true source becomes (and remains) small, for a
 134 predetermined sequence length. To characterize the classes of probability distributions on \mathbb{N} that are
 135 data-derived weakly compressible, we shall introduce the notion of what it means for a probability
 136 distribution in the class to be *deceptive* relative to the class. At a high level, a source probability
 137 distribution, viewed as a member of a collection of probability distributions, is *deceptive* if the asymptotic
 138 per-symbol redundancy of neighborhoods of the source within the model class is bounded away from 0,
 139 in the limit as the neighborhood shrinks to 0. Then, in our main finding, Theorem 17, we show that *a*
 140 *collection of probability measures is data-derived weakly compressible iff no source in the model class is*
 141 *deceptive.*

142 As we delve deeper into this formulation, we will see that data-derived consistency changes how we
 143 think of model classes. It shifts the focus away from the global complexity of the model class to some
 144 form of local complexity of each model within the model class, viewed as a member of the model class.

145 Our notion of data-derived consistency is closely related to other formulations in compression and
 146 learning theory – in particular hierarchical universal compression [13] and data-dependent structural risk
 147 minimization [14], as well as its subsequent development via the *luckiness* framework [15]. Fundamental
 148 to all these approaches is to balance the sample complexity of learning with the desire for richer model
 149 collections (or hypothesis collections as the case may be).

150 The paper is organized as follows. In the next section we develop our data-derived approach. Section III
 151 recalls some of the central prior results on universal compression that we build on in our work. Section IV
 152 discusses our main result (Theorem 17), which completely characterizes *d.w.c.* model classes of *i.i.d.*
 153 probability distributions on a countable set. We then illustrate several nuances in our formulation and
 154 results using several examples in Section V. Sections VI and VII are devoted to proving the main
 155 result. The main thread of the discussion is supported by several appendices. Appendix I reconciles
 156 the traditional definitions of strong and weak compressibility with those we work with in this paper.
 157 Appendix II gathers several basic results on entropy, redundancy and the Jensen-Shannon divergence

158 that we draw upon throughout the paper. Appendix III contains the details of the proof for the claims
 159 made regarding one of the examples in Section V. Appendix IV proves a lemma needed for the proof the
 160 sufficiency part of the main theorem. The last bit of the proof of the necessity part of the main theorem
 161 is in Appendix V and that of the sufficiency part in Appendix VI. Finally, Appendix VII corrects an
 162 erroneous claim made in passing in the concluding remarks in [2] (which does not in any way affect the
 163 rest of that paper).

164 II. FORMULATION OF THE PROBLEM

165 Let \mathcal{P} be a collection of probability distribution over \mathbb{N} . Given \mathcal{P} , we let \mathcal{P}^∞ denote the collection
 166 of probability measures on \mathbb{N}^∞ induced by *i.i.d.* sampling from the individual probability distributions
 167 in \mathcal{P} . We will use the term *source* to denote either $p \in \mathcal{P}$ or $p^\infty \in \mathcal{P}^\infty$ as appropriate. For notational
 168 simplicity and following the convention in literature, we will also often drop the superscript in p^∞
 169 and use p both for the probability distribution on \mathbb{N} and the corresponding *i.i.d.* probability measure
 170 induced on \mathbb{N}^∞ . Further, for $n \geq 1$ and a string of natural numbers $\mathbf{x} := (x_1, \dots, x_n) =: x^n \in \mathbb{N}^n$, we
 171 will write $p(\mathbf{x})$ or $p(x^n)$ for $\prod_{i=1}^n p(x_i)$. Here p can be thought of as a simplified notation for the product
 172 probability measure p^n on \mathbb{N}^n corresponding to the probability distribution p on \mathbb{N} .

173 We consider here the lossless compression problem for collections of large alphabet *i.i.d.* sources. The
 174 main contribution of this work is to propose and develop the data-derived framework for estimation prob-
 175 lems. The large alphabet *i.i.d.* compression problem is the vehicle we have used to do this, but the reader
 176 can no doubt easily come up with her or his favorite estimation problem where this framework might
 177 lead to interesting developments. In Example 8 we consider the problem of estimation of percentiles of
 178 the probability distribution defining the source – this has been studied in depth in [2], so here all we
 179 show is that this estimation task lies in the data-derived framework proposed in this document. Another
 180 example, which we have not studied in depth, but which seems to us to be particularly interesting, is
 181 that of entropy estimation, see Example 9.

182 Before embarking on the discussion, we introduce some additional notational conventions. For $1 \leq$
 183 $m \leq n$ and strings $\mathbf{y} \in \mathbb{N}^m$ and $\mathbf{x} \in \mathbb{N}^n$, we write $\mathbf{y} \preceq \mathbf{x}$ to denote that \mathbf{y} is a prefix of \mathbf{x} . We can also
 184 use this notation when $\mathbf{y} \in \mathbb{N}^m$ and $x \in \mathbb{N}^\infty$. The length of a finite string $\mathbf{x} \in \mathbb{N}^n$ is denoted by $|\mathbf{x}|$.

185 For a probability measure q on \mathbb{N}^∞ , given $n \geq 1$ and a string $\mathbf{x} \in \mathbb{N}^n$, we write $q(\mathbf{x})$ for the probability
 186 under q of the set of strings in \mathbb{N}^∞ whose prefix of length n is \mathbf{x} . In effect, we are treating \mathbf{x} as also
 187 denoting an event in \mathbb{N}^∞ . Note that, for $p \in \mathcal{P}$, $n \geq 1$, and $\mathbf{x} \in \mathbb{N}^n$, this notational convention is
 188 consistent with the earlier conventions of writing p for both $p^\infty \in \mathcal{P}^\infty$ and for the product probability
 189 measure on \mathbb{N}^n corresponding to p .

190 It is a standard fact that a probability measure q on \mathbb{N}^∞ is completely specified by $q(\mathbf{x})$ for all $\mathbf{x} \in \mathbb{N}^n$
 191 for all $n \geq 1$, subject to the consistency conditions $q(\mathbf{x}) = \sum_{\mathbf{y} \in \mathbb{N}^m : \mathbf{x} \preceq \mathbf{y}} q(\mathbf{y})$ for all $1 \leq n \leq m$ and
 192 $\mathbf{x} \in \mathbb{N}^n$.

193 $\{0, 1\}^*$ denotes the set of binary strings of finite length. The notation $\{0, 1\}^* \setminus \emptyset$ is used for the set of
 194 binary strings of finite length, excluding the empty string. For $\mathbf{b} \in \{0, 1\}^* \setminus \emptyset$, the length of \mathbf{b} is denoted
 195 by $l(\mathbf{b})$.

196 We write $\mathbb{1}(A)$ to denote the indicator of an event A .

197 It is convenient to state some of the supporting results in this document at a level of generality where
 198 the underlying set is a countable set, in which case we denote such a set by \mathcal{X} . Also, we will state some

199 results that apply to arbitrary collections of probability measures on \mathbb{N}^∞ , i.e. not necessarily of the form
 200 \mathcal{P}^∞ for some collection of probability distributions \mathcal{P} on \mathbb{N} . In such cases, we denote such a collection
 201 of probability measures on \mathbb{N}^∞ by Λ .

202 If q and r are arbitrary probability measures on \mathbb{N}^∞ , then

$$D_n(q||r) := E_q \log \frac{q(X^n)}{r(X^n)},$$

203 denotes the KL divergence over length n strings of q with respect to r . If p and \tilde{p} are probability
 204 distributions on \mathbb{N} , then $D(p||\tilde{p})$ denotes the KL divergence of p with respect to \tilde{p} , which is $E_p \log \frac{p(X)}{\tilde{p}(X)}$.
 205 Note that, with our conventions, the expression $D_n(p||\tilde{p})$ is also well-defined, and can be viewed as a
 206 shorthand notation for $D_n(p^\infty||\tilde{p}^\infty)$. We thus have $D_n(p||\tilde{p}) = nD(p||\tilde{p})$ for all $n \in \mathbb{N}$, since p^∞ and \tilde{p}^∞
 207 are *i.i.d.* probability measures on \mathbb{N}^∞ . KL divergence is also called *relative entropy*.

208 For probability distributions p and \tilde{p} on \mathbb{N} , their ℓ_1 distance is

$$|p - \tilde{p}|_1 := \sum_{i \in \mathbb{N}} |p(i) - \tilde{p}(i)|.$$

209 II-A. Strong compressibility and weak compressibility

210 In the lossless data compression problem for the collection of probability measures \mathcal{P}^∞ on \mathbb{N}^∞
 211 corresponding to a collection of probability distributions \mathcal{P} on \mathbb{N} , our estimator is a probability measure
 212 q on \mathbb{N}^∞ .¹ The problem formulation can be understood by thinking of the *loss* $L(p, q, \mathbf{x})$ incurred by the
 213 estimator q against a source p , given the length n observation $\mathbf{x} \in \mathbb{N}^n$, as being the *excess codelength*,

$$L(p, q, \mathbf{x}) := \log \frac{p(\mathbf{x})}{q(\mathbf{x})}.$$

214 The terminology is justified by thinking of $\log \frac{1}{p(\mathbf{x})}$ as an indication of the length of the binary string
 215 one would want to use to represent \mathbf{x} in an ideal prefix-free scheme for compressing strings of length n
 216 from the source p if one knew what p was, and thinking of $\log \frac{1}{q(\mathbf{x})}$ as the length of the binary string one
 217 would be led to use for representing \mathbf{x} in the prefix-free compression scheme suggested by the estimator
 218 q . For more on this, see the discussion in Appendix I on how strong and weak compressibility is typically
 219 defined in the literature.

220 With this loss function in mind, we now make the following definitions.

221 **Definition 2.** Let \mathcal{P} be a collection of probability distributions on \mathbb{N} , and \mathcal{P}^∞ the corresponding
 222 collection of probability measures on \mathbb{N}^∞ induced by *i.i.d.* sampling from the individual probability
 223 distributions in \mathcal{P} . Then \mathcal{P}^∞ , or equivalently \mathcal{P} , is called *strongly compressible* if there is a probability
 224 measure q on \mathbb{N}^∞ satisfying

$$\limsup_{n \rightarrow \infty} \sup_{p \in \mathcal{P}^\infty} \frac{1}{n} E_p \log \frac{p(X^n)}{q(X^n)} = 0. \quad (1)$$

225 □

226 The preceding definition may seem unusual relative to the definition of strong compressibility that is
 227 traditionally encountered in the literature on data compression [8], [1]. In Appendix I we establish that
 228 it is identical to the traditional definition.

¹It is not required that the probability measure q be generated by *i.i.d.* sampling.

Discussions of data compression in the literature are often framed in the language of *redundancy*. We formalize this notion in the following definition.

Definition 3. Let Λ be any collection of probability measures on \mathbb{N}^∞ . The *length- n redundancy* of Λ is defined to be

$$R_n(\Lambda) := \inf_q \sup_{r \in \Lambda} E_r \log \frac{r(X^n)}{q(X^n)}, \quad (2)$$

where the outer infimum is taken over all probability measures on \mathbb{N}^∞ , or equivalently over all probability measures on \mathbb{N}^n . The redundancy in the special case $n = 1$ is called the *single letter redundancy* of Λ , and $R_n(\Lambda)/n$ is called the *per-symbol length- n redundancy* of Λ . The *asymptotic per-symbol redundancy* of Λ is $\limsup_{n \rightarrow \infty} R_n(\Lambda)/n$.

More generally, given a probability measure \hat{q}_n on \mathbb{N}^n one can define the length- n redundancy of Λ with respect to \hat{q}_n to be $\sup_{r \in \Lambda} E_r \log \frac{r(X^n)}{\hat{q}_n(X^n)}$ and similarly for the per-symbol length- n redundancy of Λ with respect to \hat{q}_n . Given a probability measure q on \mathbb{N}^∞ , one can define the asymptotic-per-symbol redundancy of Λ with respect to q to be $\limsup_{n \rightarrow \infty} \frac{1}{n} \sup_{r \in \Lambda} E_r \log \frac{r(X^n)}{q(X^n)}$.

Even more generally, given a probability measure \hat{q}_n on \mathbb{N}^n one can define the length- n redundancy of $r \in \Lambda$ with respect to \hat{q}_n to be $E_r \log \frac{r(X^n)}{\hat{q}_n(X^n)}$ and define the per-symbol length- n redundancy of $r \in \Lambda$ with respect to \hat{q}_n similarly. Given a probability measure q on \mathbb{N}^∞ , one can define the asymptotic-per-symbol redundancy of $r \in \Lambda$ with respect to q to be $\limsup_{n \rightarrow \infty} \frac{1}{n} E_r \log \frac{r(X^n)}{q(X^n)}$.

When \mathcal{P} is a collection of probability distributions on \mathbb{N} , and \mathcal{P}^∞ the corresponding collection of probability measures on \mathbb{N}^∞ induced by *i.i.d.* sampling from the individual probability distributions in \mathcal{P} , we will talk about each of the redundancy quantities as properties of \mathcal{P} when in fact they are defined for \mathcal{P}^∞ . Similarly, given a probability measure \hat{q}_n on \mathbb{N}^n or a probability measure q on \mathbb{N}^∞ we will talk about each of the redundancy quantities for a given $p \in \mathcal{P}$ with respect to \hat{q}_n or q (as appropriate) when we mean the corresponding quantities for the $p^\infty \in \mathcal{P}^\infty$ corresponding to p . \square

It is worth noting that a collection of probability distributions on \mathbb{N} is strongly compressible iff its asymptotic per-symbol redundancy is zero. For completeness, we give a proof of this claim in Lemma 31 in Appendix I. We also observe that the asymptotic per-symbol redundancy of a collection of probability measures Λ on \mathbb{N}^∞ can also be written as

$$\limsup_{n \rightarrow \infty} R_n(\Lambda)/n = \limsup_{n \rightarrow \infty} \frac{1}{n} \inf_q \sup_{r \in \Lambda} E_r \log \frac{r(X^n)}{q(X^n)} = \inf_q \limsup_{n \rightarrow \infty} \frac{1}{n} \sup_{r \in \Lambda} E_r \log \frac{r(X^n)}{q(X^n)},$$

where the infimum on both sides of the equality is over probability measures q on \mathbb{N}^∞ . Namely, the $\limsup_{n \rightarrow \infty}$ can be interchanged with the \inf_q . A proof of this is given in Lemma 36 in Appendix II.

We can allow for much richer collections of probability distributions if we work with a weaker notion of compressibility.

Definition 4. Let \mathcal{P} be a collection of probability distributions on \mathbb{N} , and \mathcal{P}^∞ the collection of probability measures on \mathbb{N}^∞ induced by *i.i.d.* sampling from the individual probability distributions in \mathcal{P} . Then \mathcal{P}^∞ , or equivalently \mathcal{P} , is called *weakly compressible* if there exists a probability measure q over \mathbb{N}^∞ such that, for all $p \in \mathcal{P}^\infty$ with finite entropy rate, we have

$$\limsup_{n \rightarrow \infty} \frac{1}{n} E_p \log \frac{p(X^n)}{q(X^n)} = 0. \quad (3)$$

\square

264 One artifact of the above definition is that any collection of probability distributions on \mathbb{N} where
 265 every source has infinite entropy is vacuously weakly compressible. In Appendix I we establish that
 266 this definition of weak compressibility is identical to the definition of weak compressibility commonly
 267 encountered in the literature on data compression, see e.g. Kieffer [16]. Also, in Lemma 32 of Appendix I
 268 we formally establish the essentially tautological fact that a collection of probability distributions \mathcal{P} on
 269 \mathbb{N} is weakly compressible iff there exists a probability measure q on \mathbb{N}^∞ such that every $p \in \mathcal{P}$ with
 270 finite entropy has vanishing asymptotic per-symbol redundancy with respect to q .

271 II-B. Compression in the data-derived sense

272 Working with collections of probability distributions on \mathbb{N} that are weakly compressible gives us a
 273 richer class of models than working with those that are strongly compressible. Weak compressibility of a
 274 collection \mathcal{P} of probability distributions on \mathbb{N} ensures that there is a probability measure q on \mathbb{N}^∞ such
 275 that q is essentially as good an encoder as the underlying p for long enough strings of natural numbers
 276 drawn *i.i.d.* from p , where goodness is measured in terms of the number of bits used per symbol encoded.
 277 This is what it means to say that the asymptotic per-symbol redundancy of every $p^\infty \in \mathcal{P}^\infty$ with respect
 278 to q is 0,

279 But observe that what one means by “long enough” depends on the unknown p , since convergence
 280 to the limit in (3) need not be uniform over $p \in \mathcal{P}$. The main contribution of our work is to come to
 281 grips with this issue without having to back off all the way to being able to deal only with strongly
 282 compressible collections of probability distributions.

283 Our ideas are built around the notion of a *universal stopping rule*, which we introduce next. Recall
 284 that a stopping rule is a function of observed strings where the decision to *stop* or not at any given time
 285 is based only on what has been observed thus far. We formalize a stopping rule by a function τ from
 286 \mathbb{N}^* , the set of all finite strings of naturals, to the set $\{0, 1\}$,

$$\tau : \mathbb{N}^* \rightarrow \{0, 1\}.$$

287 When τ assigns value 0 on a finite string x^n , possibly the empty string, it indicates that the stopping
 288 rule is still waiting after having observed x^n . A string x^n , possibly the empty string, is assigned 1 if
 289 the stopping rule has stopped on any prefix of x^n . From a notational point of view, since τ quantifies a
 290 stopping rule, we will have for all strings x^n with prefix x^m that $\tau(x^n) \geq \tau(x^m)$.

291 The stopping rule τ is required to be universal for \mathcal{P} . In other words, the stopping rule cannot
 292 change depending on the unknown probabilistic model $p \in \mathcal{P}$ that is generating the observations. In the
 293 formulation that we will develop in this paper, for a given a threshold $\delta > 0$, a stopping rule, call it τ
 294 for now, will be based on some fixed probability measure q on \mathbb{N}^∞ , and will signify when the length is
 295 “long enough” that the normalized KL divergence between the underlying source distribution and the
 296 probability measure q has fallen below δ and will remain below δ henceforth. We will insist that τ stops
 297 at a finite time for all $p \in \mathcal{P}$, *i.e.*,

$$p(\lim_{n \rightarrow \infty} \tau(X^n) = 1) = 1, \text{ for all } p \in \mathcal{P}. \quad (4)$$

298 We will include the condition in (4) in the concept of what we mean by a universal stopping rule.

299 To understand this requirement better, fix a probability measure q on \mathbb{N}^∞ , and for $p \in \mathcal{P}$ let

$$\mathcal{N}_{p,\delta;q} := \{n : \frac{1}{n} E_p \log \frac{p(X^n)}{q(X^n)} > \delta\}.$$

300 Thus $\mathcal{N}_{p,\delta;q}$ is the set of all lengths $n \geq 1$ such that the length- n KL divergence of the *i.i.d.* probability
 301 measure p^∞ corresponding to p with respect to the probability measure q is worse than the accuracy
 302 required. Now consider the set

$$\mathcal{N}_{\delta;q} := \cup_{p \in \mathcal{P}} \mathcal{N}_{p,\delta;q}.$$

303 In the trivial case where $\mathcal{N}_{\delta;q}$ is a finite set, let N denote the largest element in $\mathcal{N}_{\delta;q}$. Then, for all
 304 $n \geq N$, we have

$$\sup_{p \in \mathcal{P}} \frac{1}{n} E_p \log \frac{p(X^n)}{q(X^n)} \leq \delta.$$

305 Clearly we can choose the stopping rule to be 0 for all sequences with length $n \leq N$ and 1 for all
 306 sequences with length $> N$, and this is universal.

307 It is more interesting when $\mathcal{N}_{\delta;q}$ defined above is not a finite set. Even in this case, the stopping
 308 rule τ has to stop at a finite time almost surely no matter which source is governing the observations.
 309 Naturally, no matter when τ stops waiting, the sequence length may not be long enough for some sources
 310 in \mathcal{P} , so τ fails on such sequences. More formally, for $\delta > 0$, τ fails with respect to q or is δ -premature
 311 *with respect to q* for a source $p \in \mathcal{P}$ and at time i if there is some string x^i such that

$$\tau(x_1^i) = 1 \text{ and } \frac{1}{i} E_p \log \frac{p(X^i)}{q(X^i)} > \delta. \quad (5)$$

312 For $p \in \mathcal{P}$, consider the subset of \mathbb{N}^∞ defined as

$$\left\{ x \in \mathbb{N}^\infty : \exists i \text{ such that } \tau(x^i) = 1 \text{ and } \frac{1}{i} \sum_{y^i \in \mathbb{N}^i} p(y^i) \log \frac{p(y^i)}{q(y^i)} > \delta \right\}. \quad (6)$$

313 For $p \in \mathcal{P}$, the above set is the set of strings on which τ is δ -premature with respect to q . While this
 314 set depends on which $p \in \mathcal{P}$ is driving the observations, this set is an event in the product σ -algebra
 315 on \mathbb{N}^∞ whatever the underlying $p \in \mathcal{P}$. To see this, note that it is a countable union of sets of the form
 316 $\{x \in \mathbb{N}^\infty : \tau(x^i) = 1\}$, $i \geq 1$ (which of the components sets lie in the union is determined, for the fixed
 317 probability measure q on \mathbb{N}^∞ , by the underlying source probability distribution p).

318 While the set in (6) may not be an empty set, we can at least try to ensure that its probability under
 319 p is small. This thought process leads to what we mean by a collection of probability distributions on
 320 \mathbb{N} being weakly compressible in the data-derived sense, formalized below. This is the central concept
 321 investigated in this paper.

322 **Definition 5.** Let \mathcal{P} be a collection of probability distributions on \mathbb{N} and \mathcal{P}^∞ the associated collection
 323 of probability measures on \mathbb{N}^∞ got by *i.i.d.* sampling from the individual distributions in \mathcal{P} . We say
 324 that \mathcal{P}^∞ , or equivalently \mathcal{P} , is *weakly compressible in the data-derived sense* or *data-derived weakly*
 325 *compressible* (d.w.c.) if there is a probability measure q on \mathbb{N}^∞ such that, for any accuracy $\delta > 0$ and
 326 confidence probability $0 < 1 - \eta < 1$, there is a universal stopping rule $\tau_{\delta,\eta}$ with the property that, no
 327 matter what $p^\infty \in \mathcal{P}^\infty$ is in force, we have

$$\begin{aligned} p(\tau_{\delta,\eta} \text{ is } \delta\text{-premature with respect to } q \text{ for } p) & \quad (7) \\ := p(\exists i \text{ such that } \tau_{\delta,\eta}(X^i) = 1 \text{ and } \frac{1}{i} \sum_{y^i \in \mathbb{N}^i} p(y^i) \log \frac{p(y^i)}{q(y^i)} > \delta) & < \eta. \end{aligned}$$

328 □

329 **Claim 6. (Strongly compressible implies *d.w.c.*)** Suppose \mathcal{P} is a collection of probability
 330 distributions on \mathbb{N} that is strongly compressible, namely there exists a probability measure q on \mathbb{N}^∞
 331 that satisfies (1). It follows then that, for all $\delta > 0$, the sets

$$N_{\delta;q} := \{n : \sup_{p \in \mathcal{P}^\infty} \frac{1}{n} E_p \log \frac{p(X^n)}{q(X^n)} > \delta\}$$

332 are finite. For any $\eta > 0$, suppose we set $\tau_{\delta,\eta}(x^i) = 1$ if $i > \max N_{\delta;q}$ and 0 else, we obtain for all
 333 $p \in \mathcal{P}^\infty$ that $p(\tau_{\delta,\eta}$ is δ -premature with respect to q) = 0. Thus every strongly compressible collection
 334 of probability distributions on \mathbb{N} is *d.w.c.*. \square

335 **Claim 7. (*d.w.c.* implies weakly compressible)** Suppose \mathcal{P} is a collection of probability distri-
 336 butions on \mathbb{N} that is *d.w.c.*, as in Definition 5. Let q be a probability measure on \mathbb{N}^∞ such that, for
 337 every accuracy $\delta > 0$ and confidence probability $0 < 1 - \eta < 1$ there is a universal stopping time $\tau_{\delta,\eta}$
 338 satisfying (7) for every $p \in \mathcal{P}$. Fix $p \in \mathcal{P}$. From (7) we conclude that, for all $i \geq 1$, we have

$$p(\tau_{\delta,\eta}(X^i) = 1) \mathbb{1} \left(\frac{1}{i} \sum_{y^i \in \mathbb{N}^i} p(y^i) \log \frac{p(y^i)}{q(y^i)} > \delta \right) < \eta.$$

339 However, since the stopping time $\tau_{\delta,\eta}$ is universal, it must satisfy (4), i.e. it stops eventually. Hence we
 340 have

$$\lim_{i \rightarrow \infty} p(\tau_{\delta,\eta}(X^i) = 1) = 1.$$

341 From this, it follows that

$$\limsup_{i \rightarrow \infty} \frac{1}{i} \sum_{y^i \in \mathbb{N}^i} p(y^i) \log \frac{p(y^i)}{q(y^i)} \leq \delta,$$

342 (in fact, for this to hold, it suffices to have the condition in (7) hold for some $0 < 1 - \eta < 1$ and not
 343 necessarily for all $\eta > 0$, for the given $\delta > 0$). Letting $\delta \rightarrow 0$, we see that the condition in (3) holds,
 344 for the given probability measure q on \mathbb{N}^∞ , for all $p \in \mathcal{P}$. This means, by definition, that \mathcal{P} is weakly
 345 compressible. \square

346 Claims 6 and 7 imply that

$$\text{Strongly compressible} \subseteq \textit{d.w.c.} \subseteq \text{weakly compressible}.$$

347 In Section V-A we will see examples of model classes demonstrating that each of these inclusions is
 348 strict.

349 As can be seen from the preceding discussion, our formulation of *d.w.c.* model classes is aimed at
 350 addressing the most interesting case from a statistical modeling viewpoint, which is the case where
 351 \mathcal{P}^∞ is weakly compressible, but not strongly compressible. Typically, we need global constraints on the
 352 collection of sources that comprise a model class to render the model class strongly compressible – for
 353 example, that the square root of the Fisher information be integrable over the model class for a class to
 354 be strongly compressible [10]. By contrast, as we will see, data-derived weak compressibility does not
 355 depend on controlling the entire class \mathcal{P}^∞ , but requires only that local neighborhoods of each $p \in \mathcal{P}$,
 356 viewed as a member of \mathcal{P} , be simple. Indeed, one of the main contributions of this paper is to obtain a
 357 condition that is both necessary and sufficient for an *i.i.d.* collection \mathcal{P}^∞ to be *d.w.c.*.

358 The operational interpretation for our formulation of *d.w.c.* model classes comprised of *i.i.d.* sources
 359 can be articulated as follows. Given such a model class, let q be any measure on \mathbb{N}^∞ that verifies the
 360 definition, i.e. such that for every $\delta > 0$ and $\eta > 0$ there is some universal stopping rule $\tau_{\delta,\eta} : \mathbb{N}^* \mapsto \{0, 1\}$
 361 for which the probability under every p in the model class that $\tau_{\delta,\eta}$ is δ -premature with respect to q for
 362 p is less than η .

363 As we observe the realization of the *i.i.d.* data samples from the (unknown) source p in the model class,
 364 we will eventually see a string of some (random) length $n = n(\delta, \eta, p)$ (say x^n) such that $\tau_{\delta,\eta}(x^n) = 1$.
 365 Now, even though we do not know p , we get the guarantee (with confidence probability $\geq 1 - \eta$) that
 366 using q to compress any subsequent length- n or longer sequence of symbols in the usual way (i.e.,
 367 $-\log q(x^k)$ bits for a sequence x^k) incurs an expected per-symbol redundancy $\leq \delta$.

368 II-C. Other examples of data-derived problem formulations

369 To clarify that the ideas in our framework have the potential to apply much more broadly to estimation
 370 problems other than the lossless compression problem that we have focused on in this document, we
 371 highlight in this section data-derived formulations for two other estimation problems. The first is a
 372 prediction task from [2], which we call the *insurance* problem, while the second is an entropy estimation
 373 task. In later sections, we will also make some comparisons between the insurance problem and the
 374 universal lossless compression problem studied here.

375 **Example 8. (Insurability)** Suppose we have a collection \mathcal{P}^∞ of *i.i.d.* measures over \mathbb{N}^∞ . Given a
 376 finite sample (X_1, \dots, X_n) with *i.i.d.* marginals from an unknown $p \in \mathcal{P}$ we want to estimate a finite
 377 upper bound on the next symbol X_{n+1} in a data-derived sense. If there are $p \in \mathcal{P}$ with unbounded
 378 support then for any finite upper bound we propose there is a probability under such p that it may
 379 not be valid. In our data-derived formulation, we therefore want to provide an estimated upper bound
 380 $\Phi(X_1^n)$, and a universal stopping rule τ that tells us from what point we should believe that our estimates
 381 $\Phi(X_1^n)$ are at least as big as X_{n+1} , while allowing for some probability of being wrong.

382 Formally, given a confidence probability $0 < 1 - \eta < 1$, we seek to come up with a mapping $\Phi : \mathbb{N}^* \rightarrow \mathbb{R}$
 383 and a stopping rule τ such that, for all $p \in \mathcal{P}$, we have

$$p(\exists i \in \mathbb{N} \text{ such that } \Phi(X^i) < X_{i+1} \text{ and } \tau(X^i) = 1) < \eta.$$

384 If this is possible, we say that the model class \mathcal{P}^∞ is *insurable*. In prior work, in [2], the collections \mathcal{P}^∞
 385 that are insurable were completely characterized. See Corollary 19 and Corollary 20 for more details
 386 and connections with the results developed in this document. \square

387 **Example 9. (Entropy estimation)** Let \mathcal{P} be a collection of probability distributions on \mathbb{N} . Given
 388 a finite sample (X_1, \dots, X_n) sampled *i.i.d.* from an unknown $p \in \mathcal{P}$, we want to provide a data-derived
 389 finite upper bound \hat{H} on the entropy of p . Formally, given a confidence probability $0 < 1 - \eta < 1$, we
 390 would like to come up with a mapping $\hat{H} : \mathbb{N}^* \rightarrow \mathbb{R}$ and a universal stopping rule τ such that, for all
 391 $p \in \mathcal{P}$, we have

$$p(\exists i \in \mathbb{N} \text{ such that } \hat{H} < H(p) \text{ and } \tau(X^i) = 1) < \eta. \quad \square$$

392 We do not know the answer to this question, in the sense that we do not know a simple intuitive
 393 necessary and sufficient condition that will characterize which collections \mathcal{P} of probability distributions
 394 on \mathbb{N} admit data-derived estimates of entropy and which do not.

III. BACKGROUND

395
396 This section highlights some interesting prior results on universal compression that will be used in
397 this paper. Readers can skip the proofs in this section if they are willing to take the results here at face
398 value when they are referred to. We have collected in this section the more interesting prior results we
399 use. Other, more basic, prior results that we also use are collected in Appendix II.

III-A. Weak compression

400
401 Let \mathcal{P} be a collection of probability distribution on \mathbb{N} and \mathcal{P}^∞ the collection of probability measures
402 on \mathbb{N}^∞ induced by *i.i.d.* sampling from the individual probability distributions in \mathcal{P} . In Appendix I we
403 have demonstrated that the notion of weak compressibility of \mathcal{P}^∞ in the sense of Kieffer [16] is identical
404 to the definition of weak compressibility of \mathcal{P}^∞ that we have made in Definition 4.

405 The following lemma gives a useful characterization of weak compressibility.

406 **Lemma 10.** Let \mathcal{P} be a collection of probability distributions on \mathbb{N} and \mathcal{P}^∞ the associated set of
407 *i.i.d.* probability measures on \mathbb{N}^∞ . Then \mathcal{P}^∞ is weakly compressible iff there exists a distribution q on
408 \mathbb{N} such that for all $p \in \mathcal{P}$ with finite entropy we have

$$\sum_{x \in \mathbb{N}} p(x) \log \frac{1}{q(x)} < \infty. \quad (8)$$

409 **Proof** From [16, Theorem 1] we know that \mathcal{P}^∞ is weakly compressible iff there is a countable set
410 $\mathcal{Q} := \{q_1, q_2, \dots\}$ of probability distributions on \mathbb{N} such that for all $p \in \mathcal{P}$ with finite entropy there is
411 some $q_i \in \mathcal{Q}$ satisfying

$$\sum_{x \in \mathbb{N}} p(x) \log \frac{1}{q_i(x)} < \infty.$$

412 Therefore, if there is a probability distribution q on \mathbb{N} satisfying (8) for all $p \in \mathcal{P}$, we can immediately
413 conclude that \mathcal{P}^∞ is weakly compressible. It remains to show the converse.

414 To do this, suppose that \mathcal{P}^∞ is weakly compressible and let \mathcal{Q} be a choice of the countable set of
415 probability distributions on \mathbb{N} guaranteed by [16, Theorem 1]. Fix some enumeration of \mathcal{Q} as $\mathcal{Q} =$
416 $\{q_1, q_2, \dots\}$.

417 Consider the probability distribution q on \mathbb{N} given by

$$q(n) := \frac{\sum_{i=1}^{|\mathcal{Q}|} \frac{q_i(n)}{i(i+1)}}{\sum_{j=1}^{|\mathcal{Q}|} \frac{1}{j(j+1)}}, \quad n \in \mathbb{N},$$

418 where the upper limit of the summation is understood to be ∞ if \mathcal{Q} is countably infinite. Observe that,
419 for all i and for all n , we have

$$q(n) \geq \frac{q_i(n)}{i(i+1)}.$$

420 Therefore, for all $p \in \mathcal{P}$ with finite entropy and all $q_i \in \mathcal{Q}$, we have

$$\sum_{x \in \mathbb{N}} p(x) \log \frac{1}{q(x)} \leq \sum_{x \in \mathbb{N}} p(x) \log \frac{i(i+1)}{q_i(x)}.$$

421 Since the right hand side of the preceding equation is finite for at least one $q_i \in \mathcal{Q}$, this completes the
422 proof. \square

III-B. Finite redundancy implies tightness

Let us recall the definition of *tightness* of a collection of probability distributions on \mathbb{N} .

Definition 11. A collection \mathcal{P} of probability distributions on \mathbb{N} is said to be tight if for every $\gamma > 0$ there is a natural number M_γ such that

$$\sup_{p \in \mathcal{P}} p(X > M_\gamma) < \gamma.$$

□

We now show that tightness of a collection of probability distributions on \mathbb{N} is implied by finiteness of the single letter redundancy of the collection. The result we present is a well-known folk theorem, see for example [17, Lemma 4]. Here we give an elementary proof of this result.

Lemma 12. Let \mathcal{P} be a collection of probability distributions on \mathbb{N} . If the single letter redundancy of \mathcal{P} is finite, then \mathcal{P} is tight.

Proof We reproduce the proof from [18, Lemma 1]. Since \mathcal{P} has finite single letter redundancy, there is a probability distribution q on \mathbb{N} such that

$$\dot{R} := \sup_{p \in \mathcal{P}} D(p||q) < \infty.$$

Proposition 33 in Appendix II implies that for all $p \in \mathcal{P}$ we have

$$E_p \left| \log \frac{p(X)}{q(X)} \right| \leq \dot{R} + 2(\log e)/e.$$

Hence, for all $p \in \mathcal{P}$ and all integers $m \geq 1$, we have

$$p \left(\left| \log \frac{p(X)}{q(X)} \right| > m \right) \leq (\dot{R} + (2 \log e)/e)/m. \quad (9)$$

To complete the argument, we need to define the *linearly interpolated cumulative distribution function* of a probability distribution on \mathbb{N} .

Definition 13. For a probability distribution q on \mathbb{N} , the linearly interpolated cumulative distribution $\dot{F}_q(n)$ for $n \in \mathbb{N} \cup \{0\}$ follows the standard definition of the cumulative distribution function, i.e.

$$\dot{F}_q(n) := F_q(n) = \mathbb{P}(X \leq n) \quad (10)$$

where X is a random variable distributed according to q . For $n \in \mathbb{N} \cup \{0\}$ and a real number $n \leq x \leq n+1$, however, we define

$$\dot{F}_q(x) := (n+1-x)\dot{F}_q(n) + (x-n)\dot{F}_q(n+1).$$

Note that \dot{F}_q is a nondecreasing function with domain the nonnegative real numbers and range either $[0, 1]$ or $[0, 1)$. For $t \in [0, 1)$, we define $\dot{F}_q^{-1}(t)$ to be the right continuous inverse of \dot{F}_q , i.e.

$$\dot{F}_q^{-1}(t) := \sup\{x \geq 0 : \dot{F}_q(x) \leq t\}.$$

□

446 Given $\gamma > 0$, pick m so large that $(\dot{R} + (\log e)/e)/m < \gamma/2$. For all $p \in \mathcal{P}$, we then have

$$\begin{aligned}
p(X > \dot{F}_q^{-1}(1 - \gamma/2^{m+1})) &= p(\log \frac{p(X)}{q(X)} > m, X > \dot{F}_q^{-1}(1 - \gamma/2^{m+1})) + p(\log \frac{p(X)}{q(X)} \leq m, X > \dot{F}_q^{-1}(1 - \gamma/2^{m+1})) \\
&\leq p(|\log \frac{p(X)}{q(X)}| > m) + 2^m q(X > \dot{F}_q^{-1}(1 - \gamma/2^{m+1})) \\
&< (\dot{R} + (\log e)/e)/m + \gamma/2 \\
&< \gamma.
\end{aligned}$$

447 This establishes that \mathcal{P} is tight. □

448

III-C. Bounds on redundancy

449 The following technical lemma is used in Example 23 and in Example 27. Its roots go back to [13].

450 **Lemma 14.** Let \mathcal{X} be a countable set, and \mathcal{P} be a collection of probability distributions on \mathcal{X} . For
451 i ranging over the finite set of indices $\{1, \dots, M\}$ or over all indices $i \geq 1$, let $S_i \subset \mathcal{X}$ be a subset of \mathcal{X} ,
452 and assume that these sets are pairwise disjoint. Suppose that for each i there exists $p_i \in \mathcal{P}$ such that

$$p_i(S_i) \geq \delta.$$

453 Then, for all probability distributions q on \mathcal{X} , we have

$$\sup_{p \in \mathcal{P}} D(p||q) \geq \delta \log(M) - 1,$$

454 if the number of subsets in the collection is finite, equal to M , and

$$\sup_{p \in \mathcal{P}} D(p||q) = \infty,$$

455 if the number of subsets in the collection is infinite.

456 **Proof** This is a simplified formulation of the *distinguishability* concept in [13]. To prove the claim,
457 note that for any m at most equal to the number of subsets in the collection, we must have $q(S_i) \leq 1/m$
458 for some i . For such a choice of i we can write

$$\begin{aligned}
D(p_i||q) &= \sum_{x \in S_i} p_i(x) \log \frac{p_i(x)}{q(x)} + \sum_{x \in S_i^c} p_i(x) \log \frac{p_i(x)}{q(x)} \\
&\stackrel{(a)}{\geq} p_i(S_i) \log \frac{p_i(S_i)}{q(S_i)} + p_i(S_i^c) \log \frac{p_i(S_i^c)}{q(S_i^c)} \\
&\geq p_i(S_i) \log \frac{1}{q(S_i)} + p_i(S_i^c) \log \frac{1}{q(S_i^c)} - 1 \\
&\geq \delta \log m - 1,
\end{aligned}$$

459 where step (a) is from the log sum inequality. This completes the proof. □

IV. CHARACTERIZATION OF *d.w.c.* MODEL CLASSES

460
 461 In this section we state our primary result, which is a necessary and sufficient condition for a model
 462 class comprised of a collection of probability distributions \mathcal{P} on \mathbb{N} to be data-derived weak compressible.
 463 We will see that what decides whether a model class \mathcal{P} is *d.w.c.* or not is a *local* property of
 464 the probability distributions in \mathcal{P} , viewed as members of \mathcal{P} . Namely, the characterization of data-
 465 derived weak compressibility is based on considering a property of local neighborhoods, as defined in
 466 Section IV-A, of the individual probability distributions in the model class. Distributions having bad
 467 local neighborhoods are what we call *deceptive* distributions, defined and studied in detail in Section IV-
 468 B. The notion of deceptive distributions lies at the heart of our characterization, in Theorem 17, of which
 469 model classes are *d.w.c.*.

IV-A. Local neighborhoods

470
 471 We will see in this section that what makes the local neighborhoods of a probability distribution $p \in \mathcal{P}$
 472 bad and kills *d.w.c.* is that when a stopping rule is forced by $p^\infty \in \mathcal{P}^\infty$ into certifying the accuracy
 473 of the estimate at some time (which will have to be the case, since the stopping rule has to stop with
 474 probability 1 under p), it will nevertheless be the case that there are other probability distributions in
 475 \mathcal{P} , potentially arbitrarily close to p , which induce inadequate performance on the estimator. We now
 476 proceed to make this vague description of the underlying ideas precise.

477 For probability distributions p and \tilde{p} on \mathbb{N} , we define

$$\mathcal{J}(p, \tilde{p}) := D\left(p \parallel \frac{p + \tilde{p}}{2}\right) + D\left(\tilde{p} \parallel \frac{p + \tilde{p}}{2}\right), \quad (11)$$

478 which, up to a scaling factor, is a Jensen-Shannon divergence between p and \tilde{p} [19]. The primary reason
 479 we use the Jensen-Shannon divergence in place of the KL divergence is that \mathcal{J} satisfies a pseudo-
 480 triangle inequality, as shown in Lemma 37 in Appendix II, while still retaining much of the statistical
 481 interpretation that a KL divergence has. Lemma 37, reproduced from [2], also shows that \mathcal{J} is intimately
 482 connected with the ℓ_1 -distance between probability distributions on \mathbb{N} . Indeed, \mathcal{J} generates the same
 483 topology on the set of probability distributions on \mathbb{N} that the ℓ_1 -distance does.

484 More generally, for probability measures q and r on \mathbb{N}^∞ , we use the notation

$$\mathcal{J}(q, r) := D_1\left(q \parallel \frac{q + r}{2}\right) + D_1\left(r \parallel \frac{q + r}{2}\right),$$

485 where, in the above, the KL divergences are taken between the single letter marginals of q and r on the
 486 first sample. Note that in this case it is no longer necessary that q should equal r when $\mathcal{J}(q, r) = 0$.
 487 Also note that this notation is consistent with our convention of using the notation p to represent both
 488 a probability distribution on \mathbb{N} and the corresponding $p^\infty \in \mathcal{P}^\infty$.

489 **Definition 15.** An ϵ -neighborhood of $p \in \mathcal{P}$ is the set $B(p, \epsilon; \mathcal{P})$ of all $p' \in \mathcal{P}$ such that $\mathcal{J}(p, p') < \epsilon$.
 490 □

491 For technical reasons, we will also make use of ℓ_1 -neighborhoods in the paper in addition to the
 492 ϵ -neighborhoods defined above (which we will refer to simply as *neighborhoods*). The ℓ_1 -neighborhood
 493 of radius $\epsilon > 0$ around $p \in \mathcal{P}$ is comprised of all $p' \in \mathcal{P}$ such that $|p - p'|_1 < \epsilon$.

IV-B. Deceptive distributions

494
 495 **Definition 16.** $p^\infty \in \mathcal{P}^\infty$ is said to be *deceptive* if the asymptotic per-symbol redundancy of
 496 neighborhoods of p is bounded away from 0 in the limit as the neighborhood shrinks to 0. More precisely,
 497 we define $p^\infty \in \mathcal{P}^\infty$, or equivalently $p \in \mathcal{P}$, to be deceptive if

$$\liminf_{\epsilon \rightarrow 0} \limsup_q \sup_{p' \in B(p, \epsilon; \mathcal{P})} \frac{1}{n} D_n(p' || q) > 0. \quad (12)$$

498 In the above, the infimum is over all q that are probability measures on \mathbb{N}^∞ (not necessarily obtained
 499 by *i.i.d.* sampling). The verbal description of this condition in terms of the asymptotic per-symbol
 500 redundancy of the neighborhoods of p is justified by Lemma 36, which is proved in Appendix II. \square

501 Our main result is the following Theorem 17. The necessity part of this theorem is proved in Section VI
 502 and the sufficiency part in Section VII.

503 **Theorem 17.** Let \mathcal{P} be a collection of probability distributions on \mathbb{N} and \mathcal{P}^∞ the associated collection
 504 of probability measures on \mathbb{N}^∞ got by *i.i.d.* sampling. Then \mathcal{P}^∞ is *d.w.c.* iff no $p \in \mathcal{P}$ is deceptive. \square

505 In the rest of this section we explore the concept of deceptive distributions to flesh out a few properties
 506 of such distributions and their neighborhoods. This will help to better understand Definition (12) and
 507 will set the stage for understanding the proof of Theorem 17.

508 *IV-B.1) A simpler characterization of deceptive distributions*

509 In determining whether a source $p \in \mathcal{P}$ is deceptive, (12) allows us to choose q depending on ϵ . We
 510 now show that this degree of freedom is unnecessary.

511 **Lemma 18.** If $p \in \mathcal{P}$ is not deceptive, then there is a single probability measure q^* on \mathbb{N}^∞ such that

$$\lim_{\epsilon \rightarrow 0} \limsup_{n \rightarrow \infty} \sup_{p' \in B(p, \epsilon; \mathcal{P})} \frac{1}{n} D_n(p' || q^*) = 0.$$

512 Hence, we have that p is deceptive iff for all probability measures q on \mathbb{N}^∞ we have

$$\lim_{\epsilon \rightarrow 0} \limsup_{n \rightarrow \infty} \sup_{p' \in B(p, \epsilon; \mathcal{P})} \frac{1}{n} D_n(p' || q) > 0.$$

513 **Proof** Because p is not deceptive, there exists a sequence $(\delta_m > 0, m \geq 1)$, with $\lim_{m \rightarrow \infty} \delta_m \rightarrow 0$,
 514 and a sequence of probability measures $(q_m, m \geq 1)$ on \mathbb{N}^∞ such that, for all sufficiently large $m \geq 1$,
 515 we have

$$\limsup_{n \rightarrow \infty} \sup_{p' \in B(p, 1/m; \mathcal{P})} \frac{1}{n} D_n(p' || q_m) \leq \delta_m.$$

516 Define the probability measure q^* on \mathbb{N}^∞ that, for each $n \geq 1$ and $\mathbf{x} \in \mathbb{N}^n$, assigns to the string \mathbf{x} the
 517 probability

$$q^*(\mathbf{x}) := \sum_{m \geq 1} \frac{q_m(\mathbf{x})}{m(m+1)}.$$

518 For all $m \geq 1$, $n \geq 1$ and $p' \in B(p, 1/m; \mathcal{P})$, we have

$$\frac{1}{n} D_n(p' || q^*) \leq \frac{1}{n} D_n(p' || q_m) + \frac{\log(m(m+1))}{n}.$$

519 This implies that

$$\limsup_{n \rightarrow \infty} \sup_{p' \in B(p, 1/m; \mathcal{P})} \frac{1}{n} D_n(p' || q^*) \leq \delta_m + \lim_{n \rightarrow \infty} \frac{\log(m(m+1))}{n} = \delta_m,$$

520 and so

$$\lim_{\epsilon \rightarrow 0} \limsup_{n \rightarrow \infty} \sup_{p' \in B(p, \epsilon; \mathcal{P})} \frac{1}{n} D_n(p' || q^*) = \lim_{m \rightarrow \infty} \limsup_{n \rightarrow \infty} \sup_{p' \in B(p, 1/m; \mathcal{P})} \frac{1}{n} D_n(p' || q^*) \leq \lim_{m \rightarrow \infty} \delta_m = 0.$$

521 This concludes the proof. \square

522 *IV-B.2) Neighborhoods of non-deceptive distributions are tight*

523 Recall the definition of *tightness* of a collection of probability distributions on \mathbb{N} from Definition 11.

524 The following corollary is immediate.

525 **Corollary 19.** If $p \in \mathcal{P}$ is not deceptive, then some neighborhood of p is tight.

526 **Proof** If $p \in \mathcal{P}$ is not deceptive then, for some $\epsilon > 0$, there exists $n \geq 1$ and a probability measure q
527 on \mathbb{N}^∞ such that

$$\sup_{p' \in B(p, \epsilon)} D_n(p' || q) < \infty.$$

528 From Proposition 34 in Appendix II, it follows that the single letter redundancy of the neighborhood
529 $B(p, \epsilon)$ is finite, which implies that $B(p, \epsilon)$ is tight, from Lemma 12. \square

530 The above corollary helps to make a connection between two data-derived formulations – *d.w.c.*, which
531 is considered in this document, and insurability, from Example 8. We showed in [2] that a collection of
532 *i.i.d.* probability measures \mathcal{P}^∞ on \mathbb{N}^∞ is insurable iff some neighborhood, exactly as defined here, of
533 every $p \in \mathcal{P}$ is tight. We therefore obtain

534 **Corollary 20.** Let \mathcal{P} be a collection of probability distributions on \mathbb{N} and let \mathcal{P}^∞ denote the
535 associated collection of *i.i.d.* probability measures on \mathbb{N}^∞ . If \mathcal{P}^∞ is *d.w.c.*, then \mathcal{P}^∞ is insurable. \square

536 In both cases, note that the condition relies on some neighborhood within the model class of every
537 model being simple. We expect this kind of locality to appear as a feature of the characterization of
538 which model classes admit data-derived estimators in most data-derived formulations.

539 V. EXAMPLES

540 We now discuss a series of examples that highlight various aspects of our formulation. These examples
541 also help flesh out the notion of what it means for a probability distribution to be deceptive.

542 *V-A. Strongly compressible \subsetneq d.w.c. \subsetneq weakly compressible*

543 We first give examples showing that weakly compressible collections of probability distribution on \mathbb{N}
544 are a strictly richer class of models than *d.w.c.* collections. We also show that there are collections of
545 probability distributions on \mathbb{N} that are *d.w.c.* but are not strongly compressible.

546 **Weakly compressible but not *d.w.c.***

547 We consider two examples in this category.

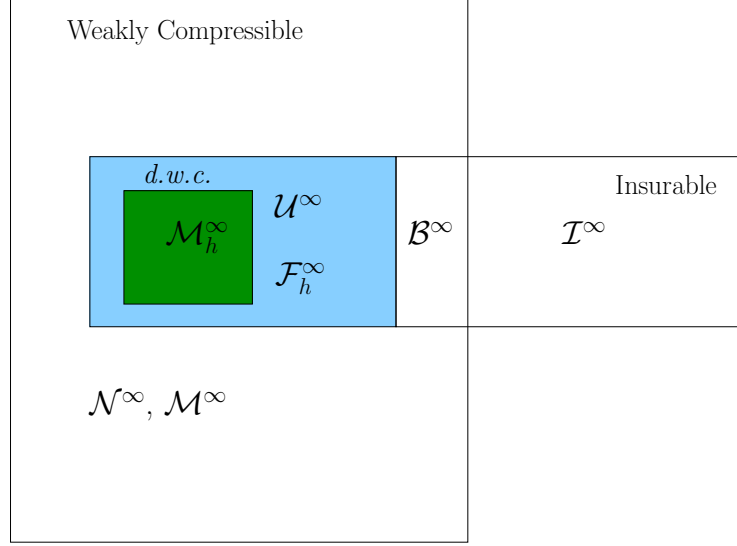


Fig. 1. Summary of examples: \mathcal{M}_h^∞ is strongly compressible (hence *d.w.c.*, insurable and weakly compressible), \mathcal{U}^∞ and \mathcal{F}_h^∞ are *d.w.c.* (hence insurable and weakly compressible), \mathcal{B}^∞ is weakly compressible and insurable but not *d.w.c.*, \mathcal{N}^∞ and \mathcal{M}^∞ are weakly compressible, but not insurable nor *d.w.c.*, while \mathcal{I}^∞ is insurable but not weakly compressible. Note that Corollary 20 shows that all *d.w.c.* collections are insurable, while Claim 6 and Claim 7 show that strong compressibility implies *d.w.c.* and that *d.w.c.* implies weak compressibility respectively.

548 A monotone probability distribution p on \mathbb{N} is one that satisfies $p(y) \geq p(y+1)$ for all $y \in \mathbb{N}$. Let
 549 \mathcal{M} denote the collection of all monotone probability distributions on \mathbb{N} and \mathcal{M}^∞ be the corresponding
 550 collection of *i.i.d.* probability measures on \mathbb{N}^∞ .

551 **Example 21.** (\mathcal{M}^∞ is weakly compressible but not *d.w.c.*)

552 To see that \mathcal{M}^∞ is weakly compressible [20] note that, for all $p \in \mathcal{M}$ and all $n \in \mathbb{N}$, we have

$$p(n) \leq \frac{1}{n}.$$

553 It follows that every $p \in \mathcal{M}$ with finite entropy must satisfy

$$\sum_{n \geq 1} p(n) \log n \leq \sum_{n \geq 1} p(n) \log \frac{1}{p(n)} < \infty. \quad (13)$$

554 Now consider the probability distribution q on \mathbb{N} assigning probability $q(n) = \frac{6}{\pi^2 n^2}$ to $n \in \mathbb{N}$. From (13)
 555 we see that, for all $p \in \mathcal{M}$ with finite entropy, we have

$$\sum_{n \geq 1} p(n) \log \frac{1}{q(n)} < \infty.$$

556 From Lemma 10 we conclude that \mathcal{M}^∞ is weakly compressible.

557 It turns out that all the probability distributions $p \in \mathcal{M}$ are deceptive. To conclude this, we show
 558 that no neighborhood around any $p \in \mathcal{M}$ is tight and then appeal to Corollary 19. This would then
 559 imply, by Theorem 17, that \mathcal{M}^∞ is not *d.w.c.*. In fact, it would have been enough to show that there
 560 exists some $p \in \mathcal{M}$ such that that no neighborhood of p is tight.

561 Let \mathcal{U} denote the collection of all uniform distributions over finite supports of form $\{m, m + 1, \dots, M\}$
 562 where m and M are positive integers with $m \leq M$. For $p \in \mathcal{M}$ and $\epsilon > 0$, consider the collection

$$\mathcal{M}(p, \epsilon) := \{p' : p' = (1 - \alpha)p + \alpha q \text{ for } q \in \mathcal{U} \cap \mathcal{M} \text{ and } 0 \leq \alpha < \epsilon\}. \quad (14)$$

563 In (14) q can be any monotone uniform distribution, namely a uniform distribution with support
 564 $\{1, \dots, M\}$ for some $M > 0$. Clearly $\mathcal{M}(p, \epsilon) \subset \mathcal{M}$. Note also that $\mathcal{M}(p, \epsilon)$ is a subset of an ℓ_1 -neighborhood
 565 of p corresponding to ℓ_1 -distance 2ϵ . We will show that $\mathcal{M}(p, \epsilon)$ is not tight for all p and all $\epsilon > 0$. By
 566 Lemma 37 and the definition of neighborhoods in Definition 15, it follows that no neighborhood of any
 567 $p \in \mathcal{M}$ is tight.

568 For $0 < \alpha < \epsilon$, let $0 < \delta < \alpha$ and $n \geq 1$. Observe that if the support $\{1, \dots, M\}$ of a uniform
 569 distribution $q' \in \mathcal{U} \cap \mathcal{M}$ satisfies $M \geq \frac{n}{1-\delta}$, then we have

$$q'\{j : j > n\} = 1 - \frac{n}{M} \geq \frac{\delta}{\alpha}.$$

570 Thus, given any $p \in \mathcal{M}$, we have a distribution $p' = (1 - \alpha)p + \alpha q' \in \mathcal{M}(p, \epsilon)$ that satisfies $p'\{j : j > n\} \geq$
 571 δ . Therefore, $\mathcal{M}(p, \epsilon)$ is not tight. This completes the argument. \square

572 For our second example, we consider the set \mathcal{N}_1^∞ of all *i.i.d.* probability measures on \mathbb{N}^∞ corresponding
 573 to the set of all probability distributions p on \mathbb{N} such that $E_p X < \infty$, denoted \mathcal{N}_1 .

574 **Example 22.** (\mathcal{N}_1^∞ is weakly compressible but not *d.w.c.*)

575 Note that every $p \in \mathcal{N}_1$ has finite entropy. Also, by definition, all $p \in \mathcal{N}_1$ satisfy $\sum_{i \geq 1} ip_i < \infty$.
 576 Therefore the simplified version of Kieffer's condition for weak compressibility, as stated in Lemma 10,
 577 is satisfied by the distribution $q(i) := 1/2^i$ ($i \geq 1$). Thus we conclude that \mathcal{N}_1 is weakly compressible.

578 We can show that every $p \in \mathcal{N}_1$ is deceptive by showing that no neighborhood of any $p \in \mathcal{N}_1$ is
 579 tight. The approach is similar to that in Example 21. Given $\epsilon > 0$, consider distributions of the form
 580 $p' = (1 - \alpha)p + \alpha q$, where $q \in \mathcal{U}$ is a uniform distribution over a support of the form $\{m, m + 1, \dots, M\}$,
 581 and $0 < \alpha < \epsilon$. Since q has finite support, we have $p' \in \mathcal{N}_1$.

582 As in Example 21 we observe that (i) the ℓ_1 distance between p' and q is strictly less than 2ϵ ; (ii)
 583 for all $0 < \delta < \alpha$ and $n \geq 1$, we can pick $q' \in \mathcal{U}$, with \mathcal{U} defined as in Example 21, whose support
 584 satisfies $M \geq \frac{n}{1-\delta}$, which then implies that the $(1 - \delta)$ -percentile of $p' := (1 - \alpha)p + \alpha q'$ can be made
 585 to lie above n . Since the above construction works for arbitrary $n \geq 1$ and in view of Lemma 37 and
 586 the way in which neighborhoods are defined in Definition 15, no neighborhood of any $p \in \mathcal{N}_1$ is tight,
 587 which shows that every $p \in \mathcal{N}_1$ is deceptive and hence, by Theorem 17, that \mathcal{N}_1 cannot be *d.w.c.* As in
 588 Example 21, to apply Theorem 17 it would have been enough to show that there is at least one $p \in \mathcal{N}_1$
 589 which is deceptive. \square

590 ***d.w.c.* but not strongly compressible**

591 The example we consider in this category is \mathcal{U} , which is defined in Example 21. Let \mathcal{U}^∞ denote the
 592 collection of all *i.i.d.* probability measures on \mathbb{N}^∞ corresponding to \mathcal{U} .

593 **Example 23.** (\mathcal{U}^∞ is not strongly compressible but is *d.w.c.*)

594 We first show that \mathcal{U} has infinite single letter redundancy. To see this, we partition \mathbb{N} into disjoint
 595 subsets $(T_i, i \geq 0)$, where $T_i := \{2^i, \dots, 2^{i+1} - 1\}$. For each T_i there is an associated distribution $p_i \in \mathcal{U}$

596 such that $p_i(T_i) = 1$. Since the number of these disjoint sets T_i is infinite, we conclude from Lemma 14
 597 that the single redundancy of \mathcal{U} is ∞ .

598 From the second part of Proposition 34 we can now conclude that the length- n redundancy of \mathcal{U} is
 599 ∞ for all $n \geq 1$, so its asymptotic per-symbol redundancy is also ∞ , which means, by Lemma 31, that
 600 \mathcal{U} is not strongly compressible.

601 To see that \mathcal{U} is *d.w.c.*, note that around each probability distribution $p \in \mathcal{U}$ there is an ℓ_1 -
 602 neighborhood that contains no other probability distribution in \mathcal{U} . Such a neighborhood has length- n
 603 redundancy equal to 0 for all n because the only possible distribution in the neighborhood is p . Hence
 604 the asymptotic per-symbol redundancy of all sufficient small neighborhoods of each $p \in \mathcal{U}$ is zero, which
 605 means, by definition, that each $p \in \mathcal{U}$ is not deceptive, see Definition 16. \square

606 Strongly compressible and *d.w.c.*

607 For completeness we next give an example of a collection of probability distributions on \mathbb{N} which is
 608 strongly compressible, hence automatically *d.w.c.*.

609 For $h > 0$, we consider the set $\mathcal{M}_h \subset \mathcal{M}$ of all monotone probability distributions on \mathbb{N} where the
 610 second moment of the self information satisfies the bound

$$E_p \left(\log \frac{1}{p(X)} \right)^2 \leq h.$$

611 Let \mathcal{M}_h^∞ denote the set of all *i.i.d.* probability measures on \mathbb{N}^∞ corresponding to \mathcal{M}_h .

612 Example 24. (\mathcal{M}_h^∞ is strongly compressible, hence *d.w.c.*)

613 Note that for any monotone probability distribution p on \mathbb{N} and all $i \geq 1$ we have $p(i) \leq 1/i$. Therefore
 614 for any $p \in \mathcal{M}_h$, if X is a random variable taking values in \mathbb{N} with the probability distribution p , we
 615 have

$$E_p \log^2(X) \leq E_p \log^2 \frac{1}{p(X)} \leq h.$$

616 Therefore, for all $p \in \mathcal{M}_h$, we have by the Cauchy-Schwarz inequality that $E_p \log X \leq \sqrt{h}$. Now, for
 617 the probability distribution q on \mathbb{N} given by $q(i) = \frac{1}{i(i+1)}$, $i \geq 1$, we have

$$\sup_{p \in \mathcal{M}_h} E_p \left(\lceil \log \frac{1}{q(X)} \rceil \right)^2 \leq \sup_{p \in \mathcal{M}_h} E_p (\log(X^2 + X) + 1)^2 \leq \sup_{p \in \mathcal{M}_h} E_p (2 \log X + 2)^2 \leq 4(\sqrt{h} + 1)^2,$$

618 where the last inequality follows because, for all $p \in \mathcal{M}_h$, we have

$$E_p (2 \log(X) + 2)^2 = 4E_p (\log^2(X) + 2 \log X + 1) \leq 4(h + 2\sqrt{h} + 1) = 4(\sqrt{h} + 1)^2.$$

619 Therefore (see Appendix III for a proof), we can construct a probability measure q^* on \mathbb{N}^∞ such that

$$\sup_{p \in \mathcal{M}_h^\infty} \frac{1}{n} D_n(p || q^*) \leq \frac{2h^{\frac{1}{4}}(\sqrt{h} + 1)}{\sqrt{\ln n}} + \pi \sqrt{\frac{2}{3n}} \log e.$$

620 From this it follows that the collection \mathcal{M}_h^∞ is strongly compressible, and therefore *d.w.c.* trivially from
 621 Claim 6. \square

622 Comparing Examples 21 and 24, we observe, that countable unions of *d.w.c.* model classes need not
 623 be *d.w.c.*. In fact, as we will see in Example 27, even finite unions of *d.w.c.* model classes need not be
 624 *d.w.c.*.

V-B. *d.w.c. collections*

625
626 Thus far, we have seen two *d.w.c.* class – \mathcal{U}^∞ and \mathcal{M}_h^∞ . But neither is completely satisfying. In the
627 collection \mathcal{U} above, there was a neighborhood around each probability measure $p \in \mathcal{U}$ with no other
628 element of \mathcal{U} . Thus \mathcal{U} trivially satisfied the local condition characterizing *d.w.c.* in Theorem 17. The \mathcal{M}_h
629 case falls into another extreme – the entire model collection \mathcal{M}_h is strongly compressible, and therefore
630 the condition characterizing *d.w.c.* in Theorem 17 was again satisfied in a trivial way.

631 We now therefore construct two additional examples of *d.w.c.* model classes that are much more
632 interesting. Our first example is of *d.w.c.* model classes \mathcal{F}_h , where neither of the two extreme situations
633 mentioned above holds. Our second example is of a *d.w.c.* model class \mathcal{H} with a source none of whose
634 neighborhoods are strongly compressible, but where the asymptotic per-symbol redundancy diminishes
635 to 0 as the neighborhood shrinks to the defining probability distribution.

636 **More interesting *d.w.c.* model classes**

637 For a probability distribution p on \mathbb{N} and a number $M > 0$, define the probability measure

$$p^{(M)}(n) := \begin{cases} p(n - M) & n \geq M + 1 \\ 0 & \text{else.} \end{cases}$$

638 Namely, $p^{(M)}$ shifts p to the right by M . Furthermore, let the *span* of any probability distribution p
639 on \mathbb{N} having finite support be defined to be the largest natural number which has non-zero probability
640 under p .

641 For $h > 0$, we consider the model classes

$$\mathcal{F}_h := \left\{ (1 - \epsilon)p_1 + \epsilon p_2^{(\text{span}(p_1)+1)} : p_1 \in \mathcal{U}, p_2 \in \mathcal{M}_h \text{ and } 0 < \epsilon < 1 \right\}.$$

642 As usual, let \mathcal{F}_h^∞ denote the set of *i.i.d.* probability measures on \mathbb{N}^∞ associated to \mathcal{F}_h . Note that the
643 initial uniform component of any $p \in \mathcal{F}_h$ is uniquely determined.

644 **Example 25. (\mathcal{F}_h^∞ is *d.w.c.*)**

645 **Proof** Let the *base* of any probability distribution over the naturals be the smallest natural number
646 which has non-zero probability. Consider any probability distribution $p = (1 - \epsilon)p_1 + \epsilon p_2^{(\text{span}(p_1)+1)} \in \mathcal{F}_h$
647 with $p_1 \in \mathcal{U}$, $p_2 \in \mathcal{M}_h$, and $0 < \epsilon < 1$. Let m denote $\text{base}(p)$ (which clearly equals $\text{base}(p_1)$), and let
648 $m + M - 1$ denote the $\text{span}(p_1)$, where $M \geq 1$. Thus $|\text{support}(p_1)| = M$.

649 Consider any probability distribution $u \in \mathcal{F}_h$, written as $u = (1 - \epsilon')u_1 + \epsilon' u_2^{(\text{span}(q_1)+1)}$, where $u_1 \in \mathcal{U}$,
650 $u_2 \in \mathcal{M}_h$, and $0 < \epsilon' < 1$. Suppose that u is within ℓ_1 distance $\frac{(1-\epsilon')^2}{M(M+1)}$ from p . We show that

$$|\text{span}(u_1)| \leq m + \left\lceil \frac{M}{1 - \epsilon} \right\rceil.$$

651 To see this, suppose to the contrary that we have

$$|\text{span}(u_1)| \geq m + \left\lceil \frac{M}{1 - \epsilon} \right\rceil + 1.$$

652 If $\text{base}(u_1) \leq m$, all elements in the support of p_1 are assigned probability $\leq \frac{1}{\frac{M}{1-\epsilon}+1}$ from u . If $\text{base}(u_1) >$
653 m , then $u(\text{base}(p_1)) = 0$. Thus, in either case, we have $u(\text{base}(p_1)) \leq \frac{1}{\frac{M}{1-\epsilon}+1}$.

654 We can now lower bound the ℓ_1 distance between p and u by

$$\frac{(1-\epsilon)}{M} - \frac{1}{\frac{M}{1-\epsilon} + 1} = \frac{(1-\epsilon)^2}{M(M+1-\epsilon)} > \frac{(1-\epsilon)^2}{M(M+1)}.$$

655 This contradiction proves the claim.

656 Now, for fixed numbers m' and M' , consider the collection $\mathcal{P}_{m',M'} \subseteq \mathcal{F}_h$ of all probability distributions
657 with base m' , and whose support of the initial uniform component is M' . Recall that \mathcal{M}_h was shown
658 to be strongly compressible in Example 24. Observe that the redundancy of $\mathcal{P}_{m',M'}$ will be at most the
659 redundancy of \mathcal{M}_h plus 1. Therefore we must also have that $\mathcal{P}_{m',M'}$ is strongly compressible.

660 The set of all probability distributions in the ℓ_1 -neighborhood of $p \in \mathcal{F}_h$ with radius $\frac{(1-\epsilon)^2}{M(M+1)}$ can be
661 decomposed into the finite union

$$\bigcup_{\substack{m',M' \\ m'+M' \leq \lceil m + \frac{M}{1-\epsilon} \rceil}} \mathcal{P}_{m',M'}.$$

662 Each component of the finite union is strongly compressible. Therefore it follows that this neighborhood
663 of $p \in \mathcal{F}_h$ is strongly compressible. Thus no $p \in \mathcal{F}_h$ is deceptive and the collection is *d.w.c.* \square

664 We construct a *d.w.c.* collection \mathcal{H} where one of the probability distributions in \mathcal{H} has no non-zero
665 neighborhood that is also strongly compressible.

666 We again partition \mathbb{N} into $(T_i, i \geq 0)$ as before, where $T_i = \{2^i, \dots, 2^{i+1} - 1\}$ for $i \geq 0$. Let \mathcal{H} contain
667 the probability distribution p_0 that assigns probability $\frac{1}{(i+1)(i+2)}$ to 2^i for all $i \geq 0$. We will construct
668 \mathcal{H} in such a way that while p_0 is not going to be deceptive in \mathcal{H} , no neighborhood of p_0 in \mathcal{H} will be
669 strongly compressible.

670 We construct \mathcal{H} in several steps. We first fix a sequence $(\epsilon_m, m \geq 2)$ such that $0 < \epsilon_m < \frac{1}{2}$ and

$$\lim_{m \rightarrow \infty} \epsilon_m = 0.$$

671 Next, for $m \geq 2$, $k \geq m$, and $j \in \{2^k + 1, \dots, 2^k + 2^{\lceil k\epsilon_m \rceil}\}$, we define the probability distribution

$$p_{m,k,j}(r) := \begin{cases} p_0(r), & \text{if } 1 \leq r \leq 2^{m-1} - 1, \\ \frac{1}{m} - \frac{1}{k+1}, & \text{if } r = 2^{m-1} + 1, \\ \frac{1}{k+1}, & \text{if } r = j, \\ 0, & \text{else.} \end{cases}$$

672 Now, for $m \geq 2$ and $k \geq m$, let

$$\mathcal{H}_{m,k} := \left\{ p_{m,k,j} : 2^k + 1 \leq j \leq 2^k + 2^{\lceil k\epsilon_m \rceil} \right\},$$

673 let

$$\mathcal{H}_m := \bigcup_{k \geq m} \mathcal{H}_{m,k},$$

674 and, finally, let

$$\mathcal{H} := \{p_0\} \cup \left(\bigcup_{m \geq 2} \mathcal{H}_m \right).$$

675 A few observations about our construction. For all $m \geq 2$, all the probability distributions in \mathcal{H}_m
 676 assign probabilities exactly as p_0 does to every element in $\cup_{i=0}^{m-2} T_i$, and the rest of their support is
 677 disjoint from that of p_0 . It follows that, for all $m \geq 2$. for all $p \in \mathcal{H}_m$, we have

$$\|p - p_0\|_1 = \frac{2}{m}.$$

678 Hence, for all $m \geq 2$, the set of probability distributions in \mathcal{H} within ℓ_1 distance $\leq \frac{2}{m}$ from p_0 is
 679 precisely $\{p_0\} \cup (\cup_{r \geq m} \mathcal{H}_r)$. Around any probability distribution in \mathcal{H} other than p_0 , there is a non-zero
 680 neighborhood containing no other probability distribution that belongs to \mathcal{H} . Therefore, none of the
 681 probability distributions in \mathcal{H} other than p_0 can possibly be deceptive. Hence, to show that \mathcal{H} is *d.w.c.*,
 682 we have to prove that p_0 is not deceptive.

683 **Example 26.** (None of the neighborhoods of $p_0 \in \mathcal{H}$ is strongly compressible.)

684 We show that for all $m \geq 2$ the collection of probability distributions \mathcal{H}_m is not strongly compressible,
 685 *i.e.*, its asymptotic per-symbol redundancy is bounded away from zero.

686 To see this, for $2^k + 1 \leq j \leq 2^k + 2^{\lceil k\epsilon_m \rceil}$, let $S_j \subset \mathbb{N}^{k+1}$ be the set of all length- $(k+1)$ sequences all
 687 of whose symbols but one are from $\cup_{i=0}^{m-1} T_i$, and there is exactly one occurrence of the number j in the
 688 sequence. Clearly, for distinct j , S_j are disjoint. Observe that

$$p_{m,k,j}(S_j) = \left(1 - \frac{1}{k+1}\right)^k \geq \frac{1}{e}.$$

689 Therefore, from Lemma 14, we have that the length- $(k+1)$ redundancy of $\mathcal{H}_{m,k}$, which we denote by
 690 $R_{k+1}(\mathcal{H}_{m,k})$, satisfies

$$\frac{R_{k+1}(\mathcal{H}_{m,k})}{k+1} \geq \frac{1}{k+1} \left(\frac{\log |\mathcal{H}_{k,m}|}{e} - 1 \right) = \frac{1}{k+1} \left(\frac{\lceil k\epsilon_m \rceil}{e} - 1 \right).$$

691 Since for all $k \geq m \geq 2$ we have $\mathcal{H}_{m,k} \subset \mathcal{H}_m$, it follows that for $m \geq 2$ the length- n redundancy of \mathcal{H}_m ,
 692 for $n \geq m+1$, which we denote by $R_n(\mathcal{H}_m)$, satisfies

$$\frac{R_n(\mathcal{H}_m)}{n} \geq \frac{R_n(\mathcal{H}_{m,n-1})}{n} \geq \frac{1}{n} \left(\frac{\lceil (n-1)\epsilon_m \rceil}{e} - 1 \right).$$

693 Hence, the asymptotic per-symbol redundancy of \mathcal{H}_m satisfies

$$\limsup_{n \rightarrow \infty} \frac{R_n(\mathcal{H}_m)}{n} \geq \frac{\epsilon_m}{e}. \quad (15)$$

694 Thus \mathcal{H}_m is not strongly compressible and, in particular, neither is any ℓ_1 neighborhood of p_0 .

695 Nevertheless, we can show that p_0 is not deceptive. We will verify that, as $m \rightarrow \infty$, the asymptotic
 696 per-symbol redundancy of an ℓ_1 neighborhood of radius $\frac{2(m+1)}{m^2}$ around p_0 goes to 0.²

697 To do so, observe from Proposition 35 that the asymptotic per-symbol redundancy of any collection
 698 of probability distributions on \mathbb{N} is upper bounded by the single-letter redundancy of the collection.
 699 Recall that for $m \geq 2$ the ℓ_1 neighborhood of radius $\frac{2(m+1)}{m^2}$ around p_0 is the collection $\{p_0\} \cup (\cup_{l \geq m} \mathcal{H}_l)$.
 700 We will verify that the single-letter redundancy of $\{p_0\} \cup (\cup_{l \geq m} \mathcal{H}_l)$ diminishes to 0 as $m \rightarrow \infty$, which
 701 will then imply that p_0 is not deceptive, using Proposition 35.

²The choice of radius $\frac{2(m+1)}{m^2}$ is made since it satisfies $\frac{2}{m} < \frac{2(m+1)}{m^2} < \frac{2}{m-1}$ for $m \geq 2$, and we defined ℓ_1 neighborhoods to be open sets.

702 For $m \geq 2$, let q_m be the probability distribution on \mathbb{N} defined by

$$q_m(r) := \begin{cases} p_0(r), & \text{if } 1 \leq r \leq 2^{m-1} - 1, \\ \frac{1}{m} - \frac{1}{m+1}, & \text{if } r = 2^{m-1} + 1, \\ \frac{1}{(k+1)(k+2) 2^{\lceil k\epsilon_m \rceil}}, & \text{if } r \in \{2^k + 1, \dots, 2^{\lceil k\epsilon_m \rceil}\}, k \geq m, \\ 0, & \text{else.} \end{cases}$$

Let $l \geq m \geq 2$. Then, for every $k \geq l$ and $j \in \{2^k + 1, \dots, 2^k + 2^{\lceil k\epsilon_l \rceil}\}$, note that $p_{l,k,j} \in \mathcal{H}_{l,k}$ and q_l assign the same probabilities as those assigned by p_0 to every number $\leq 2^{l-1} - 1$. It follows that

$$\begin{aligned} D(p_{l,k,j} \| q_l) &= p_{l,k,j}(2^{l-1} + 1) \log \frac{p_{l,k,j}(2^{l-1} + 1)}{q_l(2^{l-1} + 1)} + p_{l,k,j}(j) \log \frac{p_{l,k,j}(j)}{q_l(j)} \\ &\leq \frac{1}{l} \log(l+1) + \frac{1}{k+1} \log(k+2) + \frac{1}{k+1} \log 2^{\lceil k\epsilon_l \rceil} \\ &\leq \epsilon_l + \frac{2}{l} \log(l+1) + \frac{1}{l+1}. \end{aligned} \tag{16}$$

703 Now, for $m \geq 2$, consider the mixture probability distribution \bar{q}_m on \mathbb{N} given by

$$\bar{q}_m(r) := \sum_{l \geq m} \frac{m}{l(l+1)} q_l(r).$$

704 Fix $m \geq 2$. We have seen that any probability distribution in \mathcal{H} in the ℓ_1 neighborhood of radius $\frac{2(m+1)}{m^2}$
 705 around p_0 must belong to $\{p_0\} \cup (\cup_{l \geq m} \mathcal{H}_l)$. For every $k \geq l \geq m$, and $j \in \{2^k + 1, \dots, 2^k + 2^{\lceil k\epsilon_l \rceil}\}$, we
 706 observe that $p_{l,k,j} \in \mathcal{H}_{l,k}$ and \bar{q}_m assign the same probabilities as those assigned by p_0 to every number
 707 $\leq 2^{m-1} - 1$. Also, p_0 and \bar{q}_m assign the same probabilities as those assigned by p_0 to every number
 708 $\leq 2^{m-1} - 1$. We will now use this observation to find upper bounds for $D(p_{m,k,j} \| \bar{q}_m)$ for $k \geq m$ and
 709 $j \in \{2^k + 1, \dots, 2^k + 2^{\lceil k\epsilon_m \rceil}\}$, then for $D(p_{l,k,j} \| \bar{q}_m)$ for $k \geq l \geq m+1$ and $j \in \{2^k + 1, \dots, 2^k + 2^{\lceil k\epsilon_l \rceil}\}$,
 710 and finally for $D(p_0 \| \bar{q}_m)$.

For $k \geq m$ and $j \in \{2^k + 1, \dots, 2^k + 2^{\lceil k\epsilon_m \rceil}\}$, we write

$$\begin{aligned} D(p_{m,k,j} \| \bar{q}_m) &= p_{m,k,j}(2^{m-1} + 1) \log \frac{p_{m,k,j}(2^{m-1} + 1)}{\bar{q}_m(2^{m-1} + 1)} + p_{m,k,j}(j) \log \frac{p_{m,k,j}(j)}{\bar{q}_m(j)} \\ &\leq p_{m,k,j}(2^{m-1} + 1) \log \frac{(m+1)p_{m,k,j}(2^{m-1} + 1)}{q_m(2^{m-1} + 1)} + p_{m,k,j}(j) \log \frac{(m+1)p_{m,k,j}(j)}{q_m(j)} \\ &\leq \epsilon_m + \frac{4}{m} \log(m+1) + \frac{1}{m+1}, \end{aligned} \tag{17}$$

711 where the last step uses (16) for the choice $l = m$.

For $k \geq l \geq m + 1$ and $j \in \{2^k + 1, \dots, 2^k + 2^{\lceil k\epsilon_l \rceil}\}$, we write

$$\begin{aligned}
D(p_{l,k,j} || \bar{q}_m) &= \sum_{n=m-1}^{l-2} p_{l,k,j}(2^n) \log \frac{p_{l,k,j}(2^n)}{\bar{q}_m(2^n)} + \sum_{r=2^{l-1}}^{\infty} p_{l,k,j}(r) \log \frac{p_{l,k,j}(r)}{\bar{q}_m(r)} \\
&\leq \sum_{n=m-1}^{l-2} p_{l,k,j}(2^n) \log \frac{p_{l,k,j}(2^n)}{\frac{mq_{n+2}(2^n)}{(n+2)(n+3)}} + \sum_{r=2^{l-1}}^{\infty} p_{l,k,j}(r) \log \frac{p_{l,k,j}(r)l(l+1)}{mq_l(r)} \\
&\leq \sum_{n=m-1}^{l-2} p_{l,k,j}(2^n) \log \frac{p_{l,k,j}(2^n)}{\frac{q_{n+2}(2^n)}{(n+2)(n+3)}} + \sum_{r=2^{l-1}}^{\infty} p_{l,k,j}(r) \log \frac{p_{l,k,j}(r)}{q_l(r)} + \frac{\log(\frac{l(l+1)}{m})}{l} \\
&= \sum_{n=m-1}^{l-2} \frac{\log((n+2)(n+3))}{(n+1)(n+2)} + \sum_{r=2^{l-1}}^{\infty} p_{l,k,j}(r) \log \frac{p_{l,k,j}(r)}{q_l(r)} + \frac{\log(\frac{l(l+1)}{m})}{l} \\
&\stackrel{(a)}{\leq} \sum_{n=m-1}^{l-2} \frac{\log((n+2)(n+3))}{(n+1)(n+2)} + \epsilon_l + 4 \frac{\log(l+1)}{l} + \frac{1}{l+1} \\
&\leq \sum_{n=m-1}^{\infty} \frac{\log((n+2)(n+3))}{(n+1)(n+2)} + \epsilon_m + \frac{4 \log(m+1)}{m} + \frac{1}{m+1}, \tag{18}
\end{aligned}$$

712 where (a) uses the bound $\log(l(l+1)/m) \leq 2 \log(l+1)$, observes that $q_{n+2}(2^n) = p_0(2^n) = p_{l,k,j}(2^n)$,
713 and uses (16).

To bound $D(p_0 || \bar{q}_m)$ from above, note that $\bar{q}_m(2^n) = \frac{m}{n+2} p_0(2^n)$ for $n \geq m-1$. Therefore we have

$$\begin{aligned}
D(p_0 || \bar{q}_m) &= \sum_{n=m-1}^{\infty} p_0(2^n) \log \frac{p_0(2^n)}{\bar{q}_m(2^n)} \\
&\leq \sum_{n=m-1}^{\infty} \frac{\log(n+1)}{(n+1)(n+2)}. \tag{19}
\end{aligned}$$

714 From (17), (18), and (19), the single letter redundancy of all sources around p_0 within ℓ_1 distance $\frac{2(m+1)}{m^2}$
715 of p_0 satisfies the upper bound

$$\sup_{p \in \{p_0\} \cup (\cup_{l \geq m} \mathcal{H}_l)} D(p || \bar{q}_m) \leq \sum_{n=m-1}^{\infty} \frac{\log((n+2)(n+3))}{(n+1)(n+2)} + \epsilon_m + \frac{4 \log(m+2)}{m+1} + \frac{1}{m+1}. \tag{20}$$

716 Note that

$$\sum_{n=1}^{\infty} \frac{\log((n+2)(n+3))}{(n+1)(n+2)} < \infty.$$

717 Hence, as $m \rightarrow \infty$, each of the terms on the right side of (20) converges to 0. Since the single letter
718 redundancy of $\{p_0\} \cup (\cup_{l \geq m} \mathcal{H}_l)$ diminishes to 0 as $m \rightarrow \infty$, from Proposition 35, the asymptotic per-
719 symbol redundancy of $\{p_0\} \cup (\cup_{l \geq m} \mathcal{H}_l)$ also diminishes to zero as $m \rightarrow \infty$. Therefore p_0 is not deceptive.

720 In conclusion, none of the neighborhoods of p_0 is strongly compressible, from (15), since the asymptotic
721 per-symbol redundancy of a $\frac{2(m+1)}{m^2}$ size ℓ_1 neighborhood of p_0 is lower bounded by $\epsilon_m/e > 0$. Yet, as we
722 showed above, p_0 is not deceptive. As noted above, no other probability distribution in \mathcal{H} can possibly
723 be deceptive since it has a neighborhood of nonzero radius around it containing no other probability
724 distribution from \mathcal{H} . Therefore, \mathcal{H} is *d.w.c.*

725 □

726 *V-C. Non-d.w.c. collections*

727 We now construct two examples of non-*d.w.c.* model classes to illustrate some additional points.

728 In Example 27 we define a model class \mathcal{B} where exactly one source in the model class is deceptive. This
 729 would mean that \mathcal{B} is not *d.w.c.*. However, even though \mathcal{B} is not *d.w.c.*, removing the single deceptive
 730 source renders the rest of the model class *d.w.c.*. Put another way, adding a single source to a *d.w.c.*
 731 model class may make the resulting bigger model class not *d.w.c.*. Since a model class with one source
 732 is trivially *d.w.c.*, it follows that even finite unions of *d.w.c.* classes may not be *d.w.c.*.

733 The second example we give here is of an insurable model class \mathcal{I} that is not *d.w.c.*. See Example 8
 734 for the definition of insurability of a model class.

735 Partition \mathbb{N} into $(T_i, i \geq 0)$, where $T_i := \{2^i, \dots, 2^{i+1} - 1\}$, $i \geq 0$. For $0 < \epsilon < 1$, let $n_\epsilon = \lceil \frac{1}{\epsilon} \rceil$. Note
 736 that ϵ lies in the range $[\frac{1}{n_\epsilon}, \frac{1}{n_\epsilon - 1})$. For $1 \leq j \leq 2^{n_\epsilon}$, let $p_{\epsilon,j}$ be the probability distribution on \mathbb{N} that
 737 assigns probability $1 - \epsilon$ to the natural number 1 (or equivalently, to the set T_0), and ϵ to the natural
 738 number $2^{n_\epsilon} + j - 1$. Finally, let p_0 be a singleton probability distribution assigning probability 1 to the
 739 natural number 1.

740 Now, let \mathcal{B} (mnemonic for binary, since every probability distribution in \mathcal{B} has support of cardinality
 741 at most 2) be the collection of probability distributions on \mathbb{N} defined by

$$\mathcal{B} := \{p_{\epsilon,j} : 0 < \epsilon < 1, 1 \leq j \leq 2^{n_\epsilon}\} \cup \{p_0\}.$$

742 As usual, \mathcal{B}^∞ denotes the set of *i.i.d.* probability measures on \mathbb{N}^∞ corresponding to \mathcal{B} .

743 **Example 27.** (p_0 is the unique probability distribution in \mathcal{B} that is deceptive.)

744 An ℓ_1 neighborhood of radius δ around p_0 is comprised of p_0 and the $p_{\epsilon,j}$ for all $0 < \epsilon < \delta/2$, and all
 745 $1 \leq j \leq 2^{n_\epsilon}$. For all $n \geq 1$ and $j \in \mathcal{T}_n$, let $S_{n,j}$ denote the set of all length n strings of natural numbers
 746 with exactly one appearance of j and the remaining $n - 1$ elements of the string being 1. Then, we have

$$p_{\frac{1}{n},j}(S_{n,j}) = \left(1 - \frac{1}{n}\right)^{n-1} \geq \frac{1}{e}.$$

747 For each $n \geq 1$, the sets $S_{n,j}$ are disjoint as j ranges over \mathcal{T}_n . Further, they are subsets of \mathbb{N}^n . Therefore,
 748 Lemma 14 implies that the length- n redundancy of the collection $\{p_{\frac{1}{n},j} : j \in \mathcal{T}_n\}$ is lower bounded by

$$\frac{n}{e} - 1.$$

749 Therefore, for all $n > \frac{2}{\delta}$, the length- n redundancy of the ℓ_1 neighborhood of radius δ is bounded below
 750 by $\frac{n}{e} - 1$. This implies that the asymptotic per-symbol redundancy of the ℓ_1 neighborhood of size δ is
 751 bounded below by $\frac{1}{e}$. From the second part of Lemma 18, we conclude that p_0 is deceptive.

752 On the other hand, for $0 < \epsilon < 1$, around every other probability distribution $p_{\epsilon,j} \in \mathcal{B}$, there is an
 753 ℓ_1 -neighborhood of radius $\frac{1}{n_\epsilon}$ that contains only probability distributions in \mathcal{B} that have support equal
 754 to $\{1, 2^{n_\epsilon} + j - 1\}$. For $n \geq 1$, let \hat{r}_n denote the probability measure on \mathbb{N}^n giving probability $\frac{1}{(n+1)\binom{n}{k}}$ to
 755 each of the strings in \mathbb{N}^n comprised of k occurrences of $2^{n_\epsilon} + j - 1$ and $n - k$ occurrences of 1, $0 \leq k \leq n$.

756 Let r_n be the probability measure corresponding to \hat{r}_n , as in Lemma 29. Then, for all $p \in \mathcal{B}$ in this
 757 ℓ_1 -neighborhood of $p_{\epsilon,j} \in \mathcal{B}$, we have for all n

$$D_n(p||r_n) \leq \log(n+1).$$

758 Noting that the measure r on \mathbb{N}^∞ that assigns probability

$$r(\mathbf{x}) = \sum_{m \geq 1} \frac{r_m(\mathbf{x})}{m(m+1)}$$

759 satisfies

$$\limsup_{n \rightarrow \infty} \sup_{p: |p - p_{\epsilon,j}| < \frac{1}{n_\epsilon}} \frac{1}{n} D_n(p||q) \leq \lim_{n \rightarrow \infty} \frac{\log n}{n} = 0,$$

760 we conclude that for every $p_{\epsilon,j} \in \mathcal{B}$ there is an ℓ_1 -neighborhood of $p_{\epsilon,j}$ that has zero asymptotic per-
 761 symbol redundancy. Hence, by Lemma 37, there is a neighborhood of $p_{\epsilon,j}$ that has zero asymptotic
 762 per-symbol redundancy. We conclude that, while p_0 is deceptive, no other probability distribution in \mathcal{B}
 763 is deceptive.

764 Indeed, this is quite intuitive when we think about what is involved operationally in compressing
 765 strings of integers whose statistics are *i.i.d.* and governed by a probability distribution in \mathcal{B} . If at any
 766 point we see two distinct symbols in such a string, there is no ambiguity about what the underlying
 767 distribution is from that point on, and very little ambiguity in the probabilities of the two distinct
 768 symbols seen, of which one must be the symbol 1. But if we see a string of all 1s we can never be sure
 769 (no matter what the length of the string) what the underlying source is. One possibility is that the
 770 source is p_0 .

771 But having seen a string of 1s of length m , there is also a reasonable chance that the underlying
 772 source could be $p_{\epsilon,j}$ for some $\epsilon \ll \frac{1}{m}$ and any $j \in T_{n_\epsilon}$. There are 2^{n_ϵ} such possible values j can take in
 773 T_{n_ϵ} , so any description of j requires an additional n_ϵ bits or $\gg m$ bits.

774 However, if we remove p_0 from the collection, we have no such trouble. We have no obligation to stop
 775 on any finite length string of all 1s, no matter how long it is, since the sequence of all 1s has probability
 776 0 under every source in \mathcal{B} other than p_0 . \square

777 The last example is a collection \mathcal{I} of probability measures over \mathbb{N} that is insurable but not *d.w.c.*. In
 778 fact \mathcal{I} is not even weakly compressible.

779 Partition \mathbb{N} into the sets $(T_i, i \geq 0)$ as before, where $T_i := \{2^i, \dots, 2^{i+1} - 1\}$. For each $i \geq 1$, pick
 780 exactly one element of T_i and assign it probability $1/(i(i+1))$. We define \mathcal{I} to be the collection of all
 781 probability distributions on \mathbb{N} that can be formed in this way. \mathcal{I}^∞ denotes the set of *i.i.d.* probability
 782 measures on \mathbb{N}^∞ corresponding to \mathcal{I} .

783 **Example 28.** (\mathcal{I} is insurable but not weakly compressible, hence not *d.w.c.*)

784 For all $p \in \mathcal{I}$ and all $k \geq 1$, we have

$$\sum_{n \geq 2^k} p(n) = \frac{1}{k}.$$

785 This means that the entire set \mathcal{I} is tight. By [2, Theorem 1], we can therefore conclude that \mathcal{I} is
 786 insurable.

787 On the other hand, for every probability distribution q on \mathbb{N} , for all $i \geq 1$ there is $x_i \in T_i$ such that

$$q(x_i) \leq \frac{1}{2^i}.$$

788 By the definition of \mathcal{I} , there is a probability distribution $p \in \mathcal{I}$ that has support $\{x_i : i \geq 1\}$. Note that
 789 $D(p||q) = \infty$. Since every probability distribution in \mathcal{I} has finite entropy (in fact they all have the same
 790 entropy), from Lemma 10 we conclude that \mathcal{I} is not weakly compressible. In particular, \mathcal{I} is not *d.w.c.*
 791 □

792 VI. NECESSITY PART OF THEOREM 17

793 In this section we prove the necessity part of Theorem 17. Namely, we prove that the existence of
 794 deceptive distributions kills *d.w.c.*. More precisely, we prove that if \mathcal{P} is a collection of probability
 795 distributions on \mathbb{N} and \mathcal{P}^∞ the associated collection of *i.i.d.* probability measures on \mathbb{N}^∞ , then \mathcal{P}^∞ is
 796 *d.w.c.* only if no $p \in \mathcal{P}^\infty$ is deceptive.

797 To prove this, suppose $p \in \mathcal{P}$ is deceptive. Then, by the second part of Lemma 18, for every probability
 798 measure q on \mathbb{N}^∞ we can find $\delta > 0$ such that

$$\lim_{\epsilon' \rightarrow 0} \limsup_{n \rightarrow \infty} \sup_{p' \in B(p, \epsilon'; \mathcal{P})} \frac{1}{n} D_n(p'||q) > \delta.$$

799 Pick any $0 < \eta < 1$, and let τ be a stopping rule. We will demonstrate that there is some $\tilde{p} \in \mathcal{P}$ such
 800 that

$$\tilde{p}(\tau \text{ is } \delta\text{-premature with respect to } q \text{ for } \tilde{p}) > \eta,$$

801 where we refer to the discussion around (5) to recall what it means for a stopping time to be δ -premature
 802 for the probability distribution $\tilde{p} \in \mathcal{P}$, with respect to the probability measure q on \mathbb{N}^∞ .

803 In order to do this, for all $n \geq 1$ let

$$A_n := \{x^n \in \mathbb{N}^n : \tau(x^n) = 1\}$$

804 denote the set of sequences of length n on which τ has entered. Note that $p(A_n)$ is increasing with n
 805 and $\lim_{n \rightarrow \infty} p(A_n) = 1$. We can therefore pick $n \geq 4/(1 - \eta)$ large enough such that $p(A_n) \geq (1 + \eta)/2$.

806 Let $\epsilon := \frac{\log \epsilon}{16n^8}$. Applying Lemma 39 in Appendix II to *i.i.d.* probability measures over length- n
 807 strings, we see that for all $\tilde{p} \in \mathcal{P}$ such that $\mathcal{J}(p, \tilde{p}) \leq \epsilon$, we have

$$\tilde{p}(A_n) > (1 + \eta)/2 - \frac{2}{n} \geq \eta.$$

808 Since $\limsup_{n \rightarrow \infty} \sup_{p' \in B(p, \epsilon'; \mathcal{P})} \frac{1}{n} D_n(p'||q)$ is nondecreasing in ϵ' , we can choose $\tilde{p} \in B(p, \epsilon; \mathcal{P})$ such
 809 that for some $n \geq 1$ we have

$$\tilde{p}(A_n) > \eta \text{ and } \frac{1}{n} D_n(\tilde{p}||q) > \delta.$$

810 This in turn means, for the choice of η and δ above, that $\tilde{p}(\tau \text{ is } \delta\text{-premature with respect to } q \text{ for } \tilde{p}) >$
 811 η . This completes the proof of the necessity part of Theorem 17.

812 As a caveat regarding the structure of this proof, we remark that the presence of a deceptive distri-
 813 bution $p \in \mathcal{P}$ does not automatically imply that any other probability distribution in any neighborhood
 814 of the deceptive distribution p is also deceptive. For example, the class \mathcal{B} in Example 27 has only p_0
 815 deceptive, while no other distribution in its neighborhood is.

³Please note that in the interest of simplicity, we have not attempted to provide the best scaling for ϵ or the tightest possible bounds.

VII. SUFFICIENCY PART OF THEOREM 17

816
817 In this section we prove the sufficiency part of Theorem 17. Namely, we prove that if a collection \mathcal{P}
818 of probability distributions on \mathbb{N} does not contain any deceptive distributions, then \mathcal{P} is *d.w.c.*. We do
819 this by explicitly constructing a probability measure q^* on \mathbb{N}^∞ such that, given any desired confidence
820 probability $0 < 1 - \eta < 1$ and accuracy $\delta > 0$, there is a stopping rule τ such that, for every $p \in \mathcal{P}$,
821 under p , τ is δ -premature with respect to q^* for p , as defined in (5), with probability at most η . Note
822 that it suffices to prove this for all δ of the form $\frac{1}{m}$ for $m \geq 1$, so will restrict attention to this case, and
823 denote the corresponding stopping rule we construct by $\tau_{\eta,m}$.

824 Suppose $p \in \mathcal{P}$ is not deceptive. From Lemma 18, there is a probability measure q_p on \mathbb{N}^∞ such that
825 for all $m \geq 1$ we can pick $\epsilon_{p,m} > 0$ satisfying

$$\limsup_{n \rightarrow \infty} \sup_{p' \in B(p, \epsilon_{p,m}; \mathcal{P})} \frac{1}{n} D_n(p' || q_p) < \frac{1}{m}. \quad (21)$$

826 We fix such an $\epsilon_{p,m} > 0$ for each $p \in \mathcal{P}$ and $m \geq 1$, satisfying the additional technical requirement that
827 $\epsilon_{p,m} < 16 \log e$.

828 For $\delta \geq 1$, let $m = 1$ and for $0 < \delta < 1$ let $m = \lceil 1/\delta \rceil$. Therefore m is the natural number such that
829 $\frac{1}{m} \leq \delta < \frac{1}{m-1}$. For any $\delta > 0$, we call $\epsilon_{p, \lceil 1/\delta \rceil}$ the δ -reach of p . In particular, $\epsilon_{p,m} > 0$ is the $\frac{1}{m}$ -reach
830 of p .

831 The intuitive meaning of the reach $\epsilon_{p,\delta}$ of a probability distribution $p \in \mathcal{P}$ is that, even if the statistics
832 of the observations are being determined by some probability distribution in \mathcal{P} within the reach of p that
833 is not necessarily p , we have control, by waiting long enough, over the amount of harm, as determined
834 by δ , that will be done if we decide instead that the statistics of the observations are being determined
835 by p . This rough heuristic will be made more precise in what follows. Note that, for any $m \geq 1$, we do
836 not require any regularity over $p \in \mathcal{P}$ of $\epsilon_{p,m}$. The reason this does not matter will also soon become
837 apparent, and is basically because, for each $m \geq 1$, it will suffice to focus on only a countable collection
838 of $p \in \mathcal{P}$.

839 Given $m \geq 1$, the zone $Q_{p,m}$ of a probability distribution $p \in \mathcal{P}$ is defined to be the set of probability
840 distributions u on \mathbb{N} given by

$$Q_{p,m} \stackrel{\text{def}}{=} \left\{ u : |p - u|_1 < \frac{\epsilon_{p,m}^2 (\ln 2)^2}{16} \right\}, \quad (22)$$

841 where, $\epsilon_{p,m}$ is the $\frac{1}{m}$ -reach of p . Note that the probability distributions in $Q_{p,m}$ are not necessarily in
842 \mathcal{P} .

843 Note that for all $p \in \mathcal{P}$, because we have assumed that $\epsilon_{p,m} < 16 \log e$, Lemma 37 in Appendix II
844 implies that the zone $Q_{p,m}$ satisfies $Q_{p,m} \cap \mathcal{P} \subseteq B(p, \epsilon_{p,m}; \mathcal{P})$. Trivially $p \in Q_{p,m} \cap \mathcal{P}$. Therefore we
845 have we have

$$\mathcal{P} = \cup_{p \in \mathcal{P}} (Q_{p,m} \cap \mathcal{P}).$$

846 Further, since $Q_{p,m}$ is open in the ℓ_1 topology, each of the intersections $Q_{p,m} \cap \mathcal{P}$ is relatively open in
847 the ℓ_1 topology on \mathcal{P} . Since \mathcal{P} is Lindelöf under the ℓ_1 topology (see [2, Sec. 6.1] for a proof), there is a
848 countable set $\tilde{\mathcal{P}}_m \subseteq \mathcal{P}$, such that \mathcal{P} is covered by the collection of relatively open sets $(Q_{\tilde{p},m} \cap \mathcal{P}, \tilde{p} \in \tilde{\mathcal{P}}_m)$,
849 i.e. we have

$$\mathcal{P} = \cup_{\tilde{p} \in \tilde{\mathcal{P}}_m} (Q_{\tilde{p},m} \cap \mathcal{P}). \quad (23)$$

850 For any fixed $m \geq 1$, we will make a choice of such a $\tilde{\mathcal{P}}_m$ and refer to it as the *quantization* of \mathcal{P} and
 851 to elements of $\tilde{\mathcal{P}}_m$ as the *centroids* of the quantization, borrowing from commonly used literature in
 852 classification. We index the countable set of centroids, $\tilde{\mathcal{P}}_m$ by $\iota_m : \tilde{\mathcal{P}}_m \rightarrow \mathbb{N}$.

853 We now construct a probability measure q^* on \mathbb{N}^∞ and, for each $0 < \eta < 1$ and $m \geq 1$, a stopping rule
 854 $\tau_{\eta,m}$, such that the pair q^* and $\tau_{\eta,m}$ will together satisfy the required guarantee that for every $p \in \mathcal{P}$,
 855 the probability that the stopping time $\tau_{\eta,m}$ is $\frac{1}{m}$ -premature with respect to q^* for p is at most η .

856 *a) Construction of the probability measure q^* on \mathbb{N}^∞ :* For each $\tilde{p} \in \tilde{\mathcal{P}}_m$ there is a probability
 857 measure $q_{\tilde{p}}$ on \mathbb{N}^∞ satisfying (21) for \tilde{p} , with $\epsilon_{\tilde{p},m}$ denoting the $\frac{1}{m}$ -reach of \tilde{p} . Let

$$\tilde{Q}_m := \{q_{\tilde{p}} : \tilde{p} \in \tilde{\mathcal{P}}_m\}$$

858 denote the collection of these probability measures as \tilde{p} ranges over $\tilde{\mathcal{P}}_m$. Note that \tilde{Q}_m is countable and
 859 is a collection of not necessarily *i.i.d.* probability measures on \mathbb{N}^∞ . For $\tilde{q} \in \tilde{Q}_m$, set the index $\iota_m(\tilde{q})$ to
 860 be equal the index assigned to the corresponding centroid \tilde{p} in the enumeration of $\tilde{\mathcal{P}}_m$. Then define a
 861 probability measure q_m on \mathbb{N}^∞ by setting, for each $n \geq 1$ and each $\mathbf{x} \in \mathbb{N}^n$, the probability

$$q_m(\mathbf{x}) := \sum_{\tilde{q} \in \tilde{Q}_m} \frac{\tilde{q}(\mathbf{x})}{\iota_m(\tilde{q})(\iota_m(\tilde{q}) + 1)}.$$

862 Finally, let q^* be the probability measure on \mathbb{N}^∞ defined by letting

$$q^*(\mathbf{x}) := \sum_{m \geq 1} \frac{q_m(\mathbf{x})}{m(m+1)},$$

863 for each $n \geq 1$ and $\mathbf{x} \in \mathbb{N}^n$,

Now, for all $\tilde{p} \in \tilde{\mathcal{P}}_m$, we have

$$\begin{aligned} \limsup_{n \rightarrow \infty} \sup_{p' \in B(\tilde{p}, \epsilon_{\tilde{p}, \frac{1}{m}}; \mathcal{P})} \frac{1}{n} D_n(p' || q^*) &= \limsup_{n \rightarrow \infty} \sup_{p' \in B(\tilde{p}, \epsilon_{\tilde{p}, \frac{1}{m}}; \mathcal{P})} \frac{1}{n} D_n(p' || q_m) \\ &= \limsup_{n \rightarrow \infty} \sup_{p' \in B(\tilde{p}, \epsilon_{\tilde{p}, \frac{1}{m}}; \mathcal{P})} \frac{1}{n} D_n(p' || q_{\tilde{p}}) \\ &< \frac{1}{m}. \end{aligned} \tag{24}$$

864 We turn next to construct a stopping rule $\tau_{\eta,m}$ having the property that, for all $p \in \mathcal{P}$, we have

$$p(\tau_{\eta,m} \text{ is } \frac{1}{m}\text{-premature with respect to } q^* \text{ for } p) < \eta.$$

865 *b) Description of the stopping rule $\tau_{\eta,m}$:* Fix $0 < \eta < 1$ and $m \geq 1$. Let $p \in \mathcal{P}$ be the probability
 866 distribution in force, which is unknown. The idea is that we want sequences generated by the (unknown)
 867 $p \in \mathcal{P}$ to be captured by one of the centroids of the quantization $\tilde{\mathcal{P}}_m$ that have p in their $\frac{1}{m}$ -reach.

868 Consider a length- n sequence x^n on which we have not yet decided the value of $\tau_{\eta,m}(x^n)$ for any
 869 $1 \leq m \leq n$. Let x^n have type (*i.e.*, empirical distribution) t , which we now insist on thinking of as a
 870 sequence of unnormalized fractions on \mathbb{N} , in order to ensure that t determines the length of the sequence
 871 x^n that defines it. The set of centroids in $\tilde{\mathcal{P}}_m$ that can potentially *capture* t is defined to be

$$\tilde{\mathcal{P}}_{m,t} := \{\tilde{p} \in \tilde{\mathcal{P}}_m : t \in Q_{\tilde{p},m}\}.$$

872 Since $\cup_{\tilde{p} \in \tilde{\mathcal{P}}_m} (Q_{\tilde{p},m} \cap \mathcal{P})$ is an open set containing \mathcal{P} , the probability under p of the set of all sequences
 873 of length n whose type is captured by some centroid in $\tilde{\mathcal{P}}_m$ approaches 1 as $n \rightarrow \infty$.

874 Not every centroid in $\tilde{\mathcal{P}}_{m,t}$ is necessarily benign, since some of these centroids may not have the
 875 generating probability measure p within their $\frac{1}{m}$ -reach. Given that the number of centroids is countably
 876 infinite, there is no easy union bound based approach that could be invoked to resolve the issue.
 877 Therefore, when $\tilde{\mathcal{P}}_{m,t} \neq \emptyset$, we refine $\tilde{\mathcal{P}}_{m,t}$ further to $\hat{\mathcal{P}}_{m,t} \subset \tilde{\mathcal{P}}_{m,t}$ in a way that will allow us to use
 878 Lemma 40 to bound the probability of wrong capture.

879 To do so, for every $\tilde{p} \in \tilde{\mathcal{P}}_m$, with $\frac{1}{m}$ -reach $\epsilon_{\tilde{p},m}$, let

$$D_{\tilde{p},m} := \frac{\epsilon_{\tilde{p},m}^4 (\ln 2)^4}{256}.$$

880 The quantity above plays the role of γ when using Lemma 40.

881 To understand the core of our sufficiency proof, consider what happens when the underlying p happens
 882 to be outside the $\frac{1}{m}$ -reach of $p' \in \tilde{\mathcal{P}}_{m,t}$. Since p is far from p' (out of its $\frac{1}{m}$ -reach), but p' is close to
 883 the empirical distribution, t , of the observed sequence, our pseudo-triangle inequality from Lemma 37
 884 will use the quantity $D_{p',m}$ to lower bound the distance of t from the underlying p , which allows us to
 885 conclude that sequences with type t have a small probability under p .

886 The centroids in $\tilde{\mathcal{P}}_{m,t}$ that get placed into $\hat{\mathcal{P}}_{m,t}$ are those that satisfy (26) and (27) below. In what
 887 follows, the quantity $\log C(p', m)$ of a centroid $p' \in \tilde{\mathcal{P}}_{m,t}$ plays the role of the “effective size” of the
 888 support size of p' , corresponding to the number k of Lemma 40. Given $\tilde{p} \in \tilde{\mathcal{P}}_m$, we define $C(\tilde{p}, m)$ via

$$C(\tilde{p}, m) := 2^3 \left(\sup_{r \in B(\tilde{p}, \epsilon_{\tilde{p},m}; \mathcal{P})} \dot{F}_r^{-1}(1 - \sqrt{D_{\tilde{p},m}/6}) \right), \quad (25)$$

889 and we note that $C(\tilde{p}, m)$ is finite from the tightness result in Lemma 12. This is because we have

$$\limsup_{n \rightarrow \infty} \sup_{r \in B(\tilde{p}, \epsilon_{\tilde{p},m}; \mathcal{P})} \frac{1}{n} D_n(r || q^*) < \frac{1}{m},$$

890 from (24), which implies that for sufficiently large n the single letter redundancy of the family of n -fold
 891 product measures on \mathbb{N}^n corresponding to the probability distributions in $B(\tilde{p}, \epsilon_{\tilde{p},m}; \mathcal{P})$ is finite, which,
 892 by Lemma 12, implies that this family of n -fold product measures on \mathbb{N}^n is tight, which implies that
 893 the family of product distributions $B(\tilde{p}, \epsilon_{\tilde{p},m}; \mathcal{P})$ is tight.

894 With $C(p', m)$ for $p' \in \tilde{\mathcal{P}}_{m,t}$ defined as in (25), the conditions we require on $p' \in \tilde{\mathcal{P}}_{m,t}$ in order to
 895 place it in $\hat{\mathcal{P}}_{m,t}$ are

$$\exp(-n D_{p',m}/18) \leq \frac{\eta}{2C(p', m) \iota(p')^2 n(n+1)}, \quad (26)$$

896 and

$$2\dot{F}_t^{-1}(1 - \sqrt{D_{p',m}/6}) \leq \log C(p', m). \quad (27)$$

897 Note that given $\tilde{p} \in \tilde{\mathcal{P}}_m$ and a type t (which we recall determines the length n of the sequence defining
 898 it), one could ask if the conditions analogous to (26) and 27 hold or not for the pair (\tilde{p}, t) ; this observation
 899 will become important in Appendix V. It is also worth remarking that the proof of sufficiency of the
 900 necessary and sufficient condition for the insurability of a model class in [2, Thm. 1] also uses a similar
 901 criterion to bound the probability of wrong capture.

902 We are now in a position to specify the stopping rule $\tau_{\eta,m}$. Consider a sequence of natural numbers,
 903 x^n , having type t , which we recall determines the length n of the sequence, and assume that we have
 904 not yet specified $\tau_{\eta,m}$ for any prefix x^l of the sequence x^n for $1 \leq l \leq n$.

905 If $\hat{\mathcal{P}}_{m,t} = \emptyset$ there is no way to assign any element of $\hat{\mathcal{P}}_{m,t}$ to this sequence and its suffixes and so we
 906 move on to all the possible single letter extensions of the sequence x^n , without for the moment deciding
 907 what $\tau_{\eta,m}(x^n)$ is, although we it will eventually turn out to be 0.

908 If $\hat{\mathcal{P}}_{m,t} \neq \emptyset$, let \hat{p} denote the probability distribution in $\hat{\mathcal{P}}_{m,t}$ with the smallest index. All suffixes of
 909 x^n are then said to be *trapped* by \hat{p} , which means that they are assigned to $\hat{p} \in \hat{\mathcal{P}}_{m,t}$. From (24), we
 910 have

$$\limsup_{n \rightarrow \infty} \sup_{r \in B(\hat{p}, \epsilon_{\hat{p}, m}; \mathcal{P})} \frac{1}{n} D_n(r || q^*) < \frac{1}{m}.$$

911 This means that the set

$$N_{\hat{p}} := \{n : \sup_{r \in B(\hat{p}, \epsilon_{\hat{p}, m}; \mathcal{P})} \frac{1}{n} D_n(r || q^*) \geq \frac{1}{m}\} \quad (28)$$

912 is finite. For any suffix x^N of x^n , when $N > \max N_{p^*}$, we set $\tau_{\eta,m}(x^N) = 1$, 0 else.

913 Finally for each finite string x^n for which the value of $\tau_{\eta,m}(x^n)$ has not yet been decided, we set this
 914 value to be 0. It can be checked that $\tau_{\eta,m}$ so defined is a stopping time. This is because if $\tau_{\eta,m}(x^n) = 0$
 915 for any sequence $x^n \in \mathbb{N}^n$, then we also have $\tau_{\eta,m}(x^m) = 0$ for $1 \leq m \leq n$, i.e. for all its prefixes.

916 c) $\tau_{\eta,m}$ enters with probability 1: This is proved in Appendix V, using an argument similar to that
 917 used in the sufficiency proof in [2].

918 d) *Probability under any $p \in \mathcal{P}$ that $\tau_{\eta,m}$ is $\frac{1}{m}$ -premature with respect to q^* for p is strictly less*
 919 *than η :* Consider any $p \in \mathcal{P}$. Among sequences of natural numbers on which $\tau_{\eta,m}$ has entered, we will
 920 distinguish between those that are in *good* traps and those in *bad* traps. If a sequence x^n is trapped
 921 by $\hat{p} \in \tilde{\mathcal{P}}_m$ such that $p \in B(\hat{p}, \epsilon_{\hat{p}, m}; \mathcal{P})$, we call \hat{p} is a good trap for that sequence. Conversely, if
 922 $p \notin B(\hat{p}, \epsilon_{\hat{p}, m}; \mathcal{P})$, \hat{p} is called a bad trap for that sequence.

923 (*Good traps*) Suppose a length- n sequence x^n is in a good trap. Namely, it is trapped by a probability
 924 distribution $\hat{p} \in \tilde{\mathcal{P}}_m$ such that $p \in B(\hat{p}, \epsilon_{\hat{p}, m}; \mathcal{P})$. Then, if $\tau_{\eta,m}(x^n) = 1$ it must be the case that
 925 $\frac{1}{n} D(p || q^*) < \frac{1}{m}$. Thus such sequences cannot contribute to the probability under p of $\tau_{\eta,m}$ being
 926 $\frac{1}{m}$ -premature with respect to q^* for p .

927 (*Bad traps*) We can show that the probability with which sequences generated by p fall into bad
 928 traps is strictly less than η using an argument, which is essentially identical to the one used in [2], based
 929 on the pseudo-triangle inequality from Lemma 37. This argument is reproduced in Appendix VI for the
 930 sake of completeness. Pessimistically, we assume that $\tau_{\eta,m}$ is $\frac{1}{m}$ -premature with respect to q^* for p on
 931 every sequence that falls into a bad trap.

932 This completes the proof of the sufficiency part of Theorem 17.

933

ACKNOWLEDGMENTS

934 This work was in part supported by the NSF Science & Technology Center for Science of Information
 935 Grant number CCF-0939370. In addition, Santhanam was also supported by NSF Grants CCF-1065632
 936 and CCF-1619452; Anantharam was also supported by the ARO MURI grant W911NF- 08-1-0233,
 937 “Tools for the Analysis and Design of Complex Multi-Scale Networks”, Marvell Semiconductor Inc.,

938 the U.C. Discovery program, the William and Flora Hewlett Foundation supported Center for Long
 939 Term Cybersecurity at Berkeley, and the NSF grants CNS-0910702, ECCS-1343998, CNS-1527846,
 940 CCF-1618145 and CCF-1901004; Szpankowski was also supported by NSF Grants CCF-1524312, CCF-
 941 2006440, and CCF-2007238.

942 REFERENCES

- 943 [1] L.D. Davisson. Universal noiseless coding. *IEEE Transactions on Information Theory*, 19(6):783–795, November
 944 1973.
- 945 [2] N. Santhanam and V. Anantharam. Agnostic insurance of model classes. *Journal of Machine Learning Research*,
 946 pages 2329–2355, 2015. Full version available from arXiv doc id: 1212:3866.
- 947 [3] C.M. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2006.
- 948 [4] V. Koltchinskii. Rademacher penalties and structural risk minimization. *IEEE Transactions on Information Theory*,
 949 47:1902–1914, 2001.
- 950 [5] P.L. Bartlett, S. Boucheron, and G. Lugosi. Model selection and error estimation. *Machine Learning*, 48:85–113, 2002.
- 951 [6] P.L. Bartlett and S. Mendelson. Rademacher and Gaussian complexities: Risk bounds and structural results. *Journal*
 952 *of Machine Learning Research*, 3:463–482, 2002.
- 953 [7] Y.M. Shtarkov. Universal sequential coding of single messages. *Problems of Information Transmission*, 23(3):3–17,
 954 1987.
- 955 [8] B. Fittingoff. Universal methods of coding for the case of unknown statistics. In *Proceedings of the 5th Symposium*
 956 *on Information Theory*, pages 129–135. Moscow-Gorky, 1972.
- 957 [9] R.E. Krichevsky and V.K. Trofimov. The performance of universal coding. *IEEE Transactions on Information Theory*,
 958 27(2):199–207, March 1981.
- 959 [10] J. Rissanen. Universal coding, information, prediction, and estimation. *IEEE Transactions on Information Theory*,
 960 30(4):629–636, July 1984.
- 961 [11] A. Barron, J. Rissanen, and B. Yu. The minimum description length principle in coding and modeling. *IEEE*
 962 *Transactions on Information Theory*, 44(6):2743–2760, October 1998.
- 963 [12] M. Drmota and W. Szpankowski. Precise minimax redundancy and regrets. *IEEE Trans. Information Theory*,
 964 50:2686–2707, 2004.
- 965 [13] N. Merhav and M. Feder. Universal prediction. *IEEE Transactions on Information Theory*, 44(6):2124–2147, October
 966 1998.
- 967 [14] S. Ben-David and S. Shalev-Schwartz. *Understanding Machine Learning*. Cambridge University Press, 2012.
- 968 [15] P. Grunwald. *The Minimum Description Length Principle*. MIT Press, 2007.
- 969 [16] J.C. Kieffer. A unified approach to weak universal source coding. *IEEE Transactions on Information Theory*,
 970 24(6):674–682, November 1978.
- 971 [17] David Haussler. A general minimax result for relative entropy. *IEEE Transactions on Information Theory*, 43(4):1276–
 972 1280, 1997.
- 973 [18] M. Hosseini and N. Santhanam. Single letter characterization of average-case strong redundancy of compressing
 974 memoryless sequences. In *Allerton conference on communication, control and computing*, 2015.
- 975 [19] Jianhua Lin. Divergence measures based on the Shannon entropy. *IEEE Transactions on Information Theory*,
 976 37(1):145–151, 1991.
- 977 [20] P. Elias. Universal codeword sets and representations of integers. *IEEE Transactions on Information Theory*,
 978 21(2):194–203, March 1975.
- 979 [21] T. Cover and J. Thomas. *Elements of Information Theory*. John Wiley and sons., 1991.
- 980 [22] A. Orłitsky, N.P. Santhanam, and J. Zhang. Universal compression of memoryless sources over unknown alphabets.
 981 *IEEE Transactions on Information Theory*, 50(7):1469–1481, July 2004.
- 982 [23] A. Orłitsky, N.P. Santhanam, K. Viswanathan, and J. Zhang. On modeling profiles instead of values. In *Uncertainty*
 983 *in Artificial Intelligence*, 2004.
- 984 [24] S. Ho and R. Yeung. On information divergence measures and joint typicality. *IEEE Transactions on Information*
 985 *Theory*, 56(12):5893–5905, 2010.
- 986 [25] T. Weissman, E. Ordentlich, G. Seroussi, S. Verdu, and M. Weinberger. Universal discrete denoising: known channel.
 987 *IEEE Transactions on Information Theory*, 51(1):5–28, 2005. See also HP Labs Tech Report HPL-2003-29, Feb 2003.

988 [26] K.L. Chung. A note on the ergodic theorem of information theory. *Annals of Mathematical Statistics*, 32:612—614,
989 1961.

APPENDIX I

ALTERNATE DEFINITIONS OF STRONG AND WEAK COMPRESSIBILITY

992 We first establish the following elementary result.

993 **Lemma 29.** For $n \geq 1$, let \hat{q}_n be a probability measure on \mathbb{N}^n . Then there is a probability measure
994 q_n on \mathbb{N}^∞ such that, for all $\mathbf{x} \in \mathbb{N}^n$, we have $q_n(\mathbf{x}) = \hat{q}_n(\mathbf{x})$.

995 **Proof** We define q_n by specifying $q_n(\mathbf{y})$ for all $\mathbf{y} \in \mathbb{N}^m$ for all $m \geq 1$. If $1 \leq m \leq n$ and $\mathbf{y} \in \mathbb{N}^m$, let

$$q_n(\mathbf{y}) := \sum_{\mathbf{x}' \in \mathbb{N}^n: \mathbf{y} \preceq \mathbf{x}'} \hat{q}_n(\mathbf{x}').$$

996 For $m \geq n$ and $\mathbf{y} \in \mathbb{N}^m$, if \mathbf{y} is \mathbf{x} followed by a string of 1s, for some $\mathbf{x} \in \mathbb{N}^n$, let

$$q_n(\mathbf{y}) := \hat{q}_n(\mathbf{x}),$$

997 else let $q_n(\mathbf{y}) := 0$. It can be checked that q_n , defined in this way, satisfies the consistency conditions
998 $q_n(\mathbf{z}) = \sum_{\mathbf{y} \in \mathbb{N}^m: \mathbf{z} \preceq \mathbf{y}} q_n(\mathbf{y})$ for all $1 \leq l \leq m$ and $\mathbf{z} \in \mathbb{N}^l$. Hence q_n defines a probability measure on
999 \mathbb{N}^∞ . It can also be checked that q_n satisfies the requirement in the statement of the lemma. \square

1000 Using Lemma 29, we now get the following result, which will help establish the equivalence of our
1001 definitions of strong and weak compressibility with those common in literature.

1002 **Lemma 30.** Let Λ be any collection of probability measures on \mathbb{N}^∞ (not necessarily *i.i.d.*). Suppose
1003 there exists a sequence of probability measures \hat{q}_n on \mathbb{N}^n such that

$$\limsup_{n \rightarrow \infty} \sup_{r \in \Lambda} \frac{1}{n} E_r \log \frac{r(X^n)}{\hat{q}_n(X^n)} = 0.$$

1004 Then there is a probability measure q on \mathbb{N}^∞ such that

$$\limsup_{n \rightarrow \infty} \sup_{r \in \Lambda} \frac{1}{n} E_r \log \frac{r(X^n)}{q(X^n)} = 0.$$

1005 **Proof** For each $n \geq 1$, let the probability measure q_n on \mathbb{N}^∞ be constructed to match the probability
1006 measure \hat{q}_n on \mathbb{N}^n , as in Lemma 29. Define the probability measure q on \mathbb{N}^∞ that, for each $n \geq 1$ and
1007 $\mathbf{x} \in \mathbb{N}^n$, assigns to \mathbf{x} the probability

$$q(\mathbf{x}) := \sum_{i=1}^{\infty} \frac{q_i(\mathbf{x})}{i(i+1)}.$$

1008 For all $n \geq 1$ we therefore have

$$\begin{aligned} \sup_{r \in \Lambda} \frac{1}{n} E_r \log \frac{r(X^n)}{q(X^n)} &\leq \sup_{r \in \Lambda} \frac{1}{n} E_r \log \frac{r(X^n)}{q_n(X^n)} + \frac{\log(n(n+1))}{n} \\ &= \sup_{r \in \Lambda} \frac{1}{n} E_r \log \frac{r(X^n)}{\hat{q}_n(X^n)} + \frac{\log(n(n+1))}{n}. \end{aligned}$$

1009 Hence

$$\limsup_{n \rightarrow \infty} \sup_{r \in \Lambda} \frac{1}{n} E_r \log \frac{r(X^n)}{q(X^n)} = 0.$$

1010 □

1011 Let \mathcal{P} be a collection of probability distributions on \mathbb{N} and \mathcal{P}^∞ the collection of probability measures
 1012 on \mathbb{N}^∞ induced by *i.i.d.* sampling from the individual probability distributions in \mathcal{P} . In most prior
 1013 work [8], [1], [16] the collection \mathcal{P} is called strongly compressible if there is a sequence of probability
 1014 measures \hat{q}_n on \mathbb{N}^n such that

$$\limsup_{n \rightarrow \infty} \sup_{p \in \mathcal{P}^\infty} \frac{1}{n} E_p \log \frac{p(X^n)}{\hat{q}_n(X^n)} = 0.$$

1015 Lemma 30 immediately establishes that this definition is equivalent to the definition of strong com-
 1016 pressibility that we have made in Definition 2.

1017 The most commonly used definition of weak compressibility in prior work is due to Kieffer [16],
 1018 and is framed in the language of length functions of compression schemes. Let Λ be any collection
 1019 of stationary ergodic probability measures on \mathbb{N}^∞ (not necessarily *i.i.d.*). A compression scheme is a
 1020 sequence of mappings $\phi_n : \mathbb{N}^n \rightarrow \{0, 1\}^* \setminus \emptyset$ whose image satisfies the prefix condition, i.e. for any two
 1021 distinct elements in the domain the image of the first is not a prefix of the image of the second. The
 1022 collection Λ is called weakly compressible if there is a compression scheme $(\phi_n, n \geq 1)$ such that, for all
 1023 $r \in \Lambda$, we have

$$\lim_{n \rightarrow \infty} \frac{1}{n} E_r l(\phi_n(X^n)) = H(r),$$

1024 where $H(r)$ denotes the entropy rate of r .

1025 Let \mathcal{P} be a collection of probability distributions on \mathbb{N} and \mathcal{P}^∞ the corresponding collection of *i.i.d.*
 1026 probability measures on \mathbb{N}^∞ . Note that \mathcal{P}^∞ is a collection of stationary ergodic probability measures.
 1027 We now show that the definition of weak compressibility of \mathcal{P}^∞ in the sense of Kieffer [16] is identical
 1028 to the definition of weak compressibility of \mathcal{P}^∞ that we have made in Definition 4.

1029 Suppose first that \mathcal{P}^∞ is weakly compressible in the sense of Definition 4. If every probability
 1030 distribution in \mathcal{P} has infinite entropy, consider an arbitrary compression scheme $(\phi_n, n \geq 1)$, for instance
 1031 by defining $\phi_n(x^n)$ by concatenating symbol by symbol the representation of $i \in \mathbb{N}$ by a bit string of
 1032 length $\lceil \log \frac{1}{(i+1)(i+2)} \rceil$ coming from a prefix code for \mathbb{N} corresponding to the probability distribution
 1033 assigning probability $\frac{1}{(i+1)(i+2)}$ to $i \in \mathbb{N}$. Then we have

$$\frac{1}{n} E_p l(\phi_n(X^n)) \stackrel{(a)}{\geq} \frac{1}{n} E_p \log \frac{1}{p(X^n)} = \infty, \quad (29)$$

1034 and so

$$\lim_{n \rightarrow \infty} \frac{1}{n} E_p l(\phi_n(X^n)) = H(p),$$

1035 for all $p \in \mathcal{P}$. Here (a) in (29) can be seen by picking a probability measure q_n on \mathbb{N}^n that satisfies
 1036 $l(\phi_n(X^n)) \geq \log \frac{1}{q_n(x^n)}$ and observing that $E_p \log \frac{p(X^n)}{q_n(X^n)} \geq 0$. If there are probability distributions in
 1037 \mathcal{P} with finite entropy, let q be a probability measure on \mathbb{N}^∞ verifying the requirements in Definition 4.
 1038 For $n \geq 1$, let \hat{q}_n denote the probability measure on \mathbb{N}^n resulting from restricting q to \mathbb{N}^n . We can then
 1039 define a compression scheme $(\phi_n, n \geq 1)$ such that $l(\phi_n(\mathbf{x})) = \lceil \log \frac{1}{\hat{q}_n(\mathbf{x})} \rceil$ for all $\mathbf{x} \in \mathbb{N}^n$ for all $n \geq 1$.
 1040 Hence, for every $p \in \mathcal{P}$, we have

$$\frac{1}{n} E_p l(\phi_n(X^n)) = \frac{1}{n} E_p \lceil \log \frac{1}{\hat{q}_n(X^n)} \rceil = \frac{1}{n} E_p \lceil \log \frac{1}{q(X^n)} \rceil.$$

1041 Suppose $H(p) = \infty$. By the same argument as that used in (29) we conclude that $\frac{1}{n}E_p l(\phi_n(X^n)) = \infty$
 1042 for all $n \geq 1$ and so, for all such p , we have

$$\lim_{n \rightarrow \infty} \frac{1}{n} E_p l(\phi_n(X^n)) = H(p).$$

1043 On the other hand, if $H(p) < \infty$ we have

$$\begin{aligned} \frac{1}{n} E_p l(\phi_n(X^n)) &\leq \frac{1}{n} E_p \log \frac{1}{q(X^n)} + \frac{1}{n} \\ &= \frac{1}{n} E_p \log \frac{p(X^n)}{q(X^n)} + H(p) + \frac{1}{n}, \end{aligned}$$

1044 and so, letting $n \rightarrow \infty$, we see that

$$\lim_{n \rightarrow \infty} \frac{1}{n} E_p l(\phi_n(X^n)) = H(p)$$

1045 also holds for such p . We have established that \mathcal{P}^∞ is also weakly compressible in the sense of Kieffer [16],
 1046 irrespective of whether \mathcal{P} is comprised entirely of probability distributions with infinite entropy or also
 1047 contains probability distributions with finite entropy.

1048 For the converse, suppose that \mathcal{P}^∞ is weakly compressible in the sense of Kieffer [16]. For each $n \geq 1$
 1049 we can find a probability measure \hat{q}_n on \mathbb{N}^n such that $\hat{q}_n(\mathbf{x}) \geq 2^{-l(\phi_n(\mathbf{x}))}$ for all $\mathbf{x} \in \mathbb{N}^n$, where $(\phi_n, n \geq 1)$
 1050 is a compression scheme verifying the weak compressibility of \mathcal{P}^∞ in the sense of Kieffer [16]. For each
 1051 $n \geq 1$ we define the probability measure q_n on \mathbb{N}^∞ in terms of \hat{q}_n as in Lemma 29, and we define the
 1052 probability measure q on \mathbb{N}^∞ which, for each $n \geq 1$ and $\mathbf{x} \in \mathbb{N}^n$, assigns to \mathbf{x} the probability

$$q(\mathbf{x}) := \sum_{i=1}^{\infty} \frac{q_i(\mathbf{x})}{i(i+1)}.$$

1053 For each $p \in \mathcal{P}$ with finite entropy, we have

$$\begin{aligned} \frac{1}{n} E_p \log \frac{p(X^n)}{q(X^n)} &\leq \frac{1}{n} E_p \log \frac{p(X^n)}{q_n(X^n)} + \frac{\log n(n+1)}{n} \\ &= \frac{1}{n} E_p \log \frac{p(X^n)}{\hat{q}_n(X^n)} + \frac{\log n(n+1)}{n} \\ &\leq -H(p) + \frac{1}{n} E_p l(\phi_n(X^n)) + \frac{\log n(n+1)}{n}, \end{aligned}$$

1054 and so, from $\lim_{n \rightarrow \infty} \frac{1}{n} E_p l(\phi_n(X^n)) = H(p)$, we conclude that $\limsup_{n \rightarrow \infty} \frac{1}{n} E_p \log \frac{p(X^n)}{q(X^n)} = 0$. This
 1055 proves that \mathcal{P}^∞ is weakly compressible in the sense of Definition 4.

1056 To close this section, we give proofs of two statements that allow us to think about strong compress-
 1057 ibility and weak compressibility respectively in terms of vanishing asymptotic per-symbol redundancy.

1058 **Lemma 31.** Let \mathcal{P} be a collection of probability distribution on \mathbb{N} and \mathcal{P}^∞ the collection of probability
 1059 measures on \mathbb{N}^∞ induced by *i.i.d.* sampling from the individual probability distributions in \mathcal{P} . Then
 1060 \mathcal{P}^∞ is strongly compressible iff it has zero asymptotic per-symbol redundancy.

1061 **Proof**

1062 If \mathcal{P}^∞ is strongly compressible, then taking the probability measure q on \mathbb{N}^∞ which verifies the strong
 1063 compressibility condition in (1) from Definition 2 as the q in (2) from Definition 3 for each $n \geq 1$
 1064 immediately implies that \mathcal{P}^∞ has zero asymptotic per-symbol redundancy.

1065 Conversely, suppose \mathcal{P}^∞ has zero asymptotic per-symbol redundancy. Given $\epsilon > 0$, for each $n \geq 1$ let
 1066 q_n be a probability measure on \mathbb{N}^∞ for which $\sup_{p \in \mathcal{P}^\infty} E_p \log \frac{p(X^n)}{q_n(X^n)} \leq R_n + \epsilon$, and define the probability
 1067 measure q on \mathbb{N}^∞ by

$$q(\mathbf{x}) := \sum_{i=1}^{\infty} \frac{q_i(\mathbf{x})}{i(i+1)}.$$

1068 Then we have

$$\frac{1}{n} \sup_{p \in \mathcal{P}^\infty} E_p \log \frac{p(X^n)}{q(X^n)} \leq \frac{1}{n} \sup_{p \in \mathcal{P}^\infty} E_p \log \frac{r(X^n)}{q_n(X^n)} + \frac{\log(n(n+1))}{n},$$

1069 and so

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \sup_{p \in \mathcal{P}^\infty} E_p \log \frac{p(X^n)}{q(X^n)} \leq \epsilon.$$

1070 Letting $\epsilon \rightarrow 0$ shows that \mathcal{P}^∞ is strongly compressible. \square

1071 **Lemma 32.** Let \mathcal{P} be a collection of probability distribution on \mathbb{N} and \mathcal{P}^∞ the collection of probability
 1072 measures on \mathbb{N}^∞ induced by *i.i.d.* sampling from the individual probability distributions in \mathcal{P} . Then
 1073 \mathcal{P}^∞ is weakly compressible iff there is a probability measure q on \mathbb{N}^∞ such that for every $p \in \mathcal{P}$ with
 1074 finite entropy the corresponding $p^\infty \in \mathcal{P}^\infty$ has zero asymptotic per-symbol redundancy with respect to
 1075 q .

1076 **Proof**

1077 The claim is vacuously true if all the probability distributions in \mathcal{P} have infinite entropy. If there
 1078 are distributions in \mathcal{P} with finite entropy and \mathcal{P}^∞ is weakly compressible, then consider the probability
 1079 measure q on \mathbb{N}^∞ which verifies the weak compressibility condition in (3) from Definition 4. By definition,
 1080 with respect to this q , every $p \in \mathcal{P}$ with finite entropy is such that the corresponding $p^\infty \in \mathcal{P}^\infty$ has zero
 1081 asymptotic per-symbol redundancy with respect to q . Conversely, if there are distributions in \mathcal{P} with
 1082 finite entropy and there is a probability measure q on \mathbb{N}^∞ such that for every $p \in \mathcal{P}$ the corresponding
 1083 $p^\infty \in \mathcal{P}^\infty$ has zero asymptotic per-symbol redundancy with respect to q then, by definition, this q
 1084 satisfies the condition in (3) from Definition 4 for all $p \in \mathcal{P}$ with finite entropy. This establishes that
 1085 \mathcal{P}^∞ is weakly compressible. \square

1086 APPENDIX II

1087 BASIC PROPERTIES OF RELATIVE ENTROPY AND REDUNDANCY

1088 In this appendix we gather some basic results on the KL divergence and redundancy, which are used
 1089 at various points in the document.

1090 **Proposition 33.** Let p and q be two probability distributions on a countable set \mathcal{X} . Then

$$\sum_{x \in \mathcal{X}} p(x) \left| \log \frac{p(x)}{q(x)} \right| \leq D(p||q) + 2 \frac{\log e}{e}.$$

1091 **Proof** Let $S \subset \mathcal{X}$ be the set of all elements $x \in \mathcal{X}$ such that $p(x) \leq q(x)$. Note that $q(S) > 0$. We

1092 have

$$\begin{aligned}
D(p||q) - \sum_{x \in \mathcal{X}} p(x) \left| \log \frac{p(x)}{q(x)} \right| &= 2 \sum_{x \in \mathcal{S}} p(x) \log \frac{p(x)}{q(x)} \\
&\stackrel{(a)}{\geq} 2p(S) \log \frac{p(S)}{q(S)} \\
&\geq 2p(S) \log p(S) \\
&\geq -2 \frac{\log e}{e},
\end{aligned}$$

1093 where step (a) is from the log sum inequality. The proposition follows. \square

1094 **Proposition 34.** For all probability measures r and q on \mathbb{N}^∞ and all $1 \leq m \leq n$, we have

$$D_m(r||q) \leq D_n(r||q).$$

1095 In particular, for any collection of probability distributions \mathcal{P} on \mathbb{N} , if \mathcal{P}^∞ denotes the associated
1096 collection of *i.i.d.* probability measures on \mathbb{N}^∞ , we will have

$$R_m(\mathcal{P}) := \inf_q \sup_{p \in \mathcal{P}} E_p \log \frac{p(X^m)}{q(X^m)} \leq \inf_q \sup_{p \in \mathcal{P}} E_p \log \frac{p(X^n)}{q(X^n)} = R_n(\mathcal{P}),$$

1097 where the outer infimum on both sides is taken over all probability measures q on \mathbb{N}^∞ and so $R_m(\mathcal{P})$
1098 and $R_n(\mathcal{P})$ are the length- m redundancy and the length- n redundancy of \mathcal{P} , respectively.

1099 **Proof** The first part of the claim follows from convexity, because, for all $y^m \in \mathbb{N}^m$, we have

$$r(y^m) = \sum_{x^n : y^m \preceq x^n} r(x^n) \text{ and } q(y^m) = \sum_{x^n : y^m \preceq x^n} q(x^n).$$

1100 For the second part of the claim, for any $\epsilon > 0$ pick a probability measure q' on \mathbb{N}^∞ such that

$$\sup_{p \in \mathcal{P}} E_p \log \frac{p(X^n)}{q'(X^n)} < R_n(\mathcal{P}) + \epsilon.$$

1101 It then follows from the first part of the claim that

$$R_m(\mathcal{P}) \leq \sup_{p \in \mathcal{P}} E_p \log \frac{p(X^m)}{q'(X^m)} < R_n(\mathcal{P}) + \epsilon.$$

1102 We let $\epsilon \rightarrow 0$ to complete the proof. \square

1103 **Proposition 35.** Let \mathcal{P} be a collection of probability distributions on \mathbb{N} and \mathcal{P}^∞ the corresponding
1104 collection of probability measures on \mathbb{N}^∞ got by *i.i.d.* sampling from the individual probability distri-
1105 butions in \mathcal{P} . For $n \geq 1$, let R_n denote the length- n redundancy of \mathcal{P}^∞ , as defined in (2). Then, for all
1106 $n \geq 1$, the per-symbol length- n redundancy of \mathcal{P}^∞ satisfies $R_n/n \leq R_1$.

1107 **Proof** Let $\epsilon > 0$. Let \tilde{p} be a probability distribution on \mathbb{N} such that the single letter redundancy of
1108 \mathcal{P}^∞ with respect to \tilde{p} is strictly less than $R_1 + \epsilon$. With the usual abuse of notation, let \tilde{p} also denote
1109 the *i.i.d.* probability measure on \mathbb{N}^∞ corresponding to \tilde{p} . Then, for all $p \in \mathcal{P}$, we have

$$\frac{1}{n} E_p \log \frac{p(X^n)}{\tilde{p}(X^n)} = E_p \log \frac{p(X)}{\tilde{p}(X)} < (R_1 + \epsilon).$$

1110 By letting $\epsilon \rightarrow 0$, the proposition follows. \square

1111 **Lemma 36.** Let Λ be a collection of probability measures on \mathbb{N}^∞ . Then we have

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \inf_q \sup_{r \in \Lambda} E_r \log \frac{r(X^n)}{q(X^n)} = \inf_q \limsup_{n \rightarrow \infty} \frac{1}{n} \sup_{r \in \Lambda} E_r \log \frac{r(X^n)}{q(X^n)}, \quad (30)$$

1112 where the infimum is taken over all probability measures q on \mathbb{N}^∞ . Namely, the $\limsup_{n \rightarrow \infty}$ can be
1113 interchanged with the \inf_q in the definition of the asymptotic per-symbol redundancy of Λ .

1114 **Proof** Fix $\epsilon > 0$. For $n \geq 1$, let q_n be a probability measure on \mathbb{N}^∞ such that

$$\frac{1}{n} \sup_{r \in \Lambda} E_r \log \frac{r(X^n)}{q_n(X^n)} < \frac{1}{n} R_n + \epsilon.$$

1115 Define the probability measure \bar{q} on \mathbb{N}^∞ that, for each $n \geq 1$ and $\mathbf{x} \in \mathbb{N}^n$, assigns to \mathbf{x} the probability

$$\bar{q}(\mathbf{x}) := \sum_{i=1}^{\infty} \frac{q_i(\mathbf{x})}{i(i+1)},$$

1116 where, as usual, $q_i(\mathbf{x})$ is the probability under q_i of the event in \mathbb{N}^∞ comprised of the sequences having
1117 the prefix \mathbf{x} . We then have

$$\frac{1}{n} \sup_{r \in \Lambda} E_r \log \frac{r(X^n)}{\bar{q}(X^n)} \leq \frac{1}{n} \sup_{r \in \Lambda} E_r \log \frac{r(X^n)}{q_n(X^n)} + \frac{\log(n(n+1))}{n} < \frac{1}{n} R_n + \epsilon + \frac{\log(n(n+1))}{n}.$$

1118 Thus

$$\inf_q \limsup_{n \rightarrow \infty} \frac{1}{n} \sup_{r \in \Lambda} E_r \log \frac{r(X^n)}{q(X^n)} \leq \limsup_{n \rightarrow \infty} \frac{1}{n} \sup_{r \in \Lambda} E_r \log \frac{r(X^n)}{\bar{q}(X^n)} \leq \limsup_{n \rightarrow \infty} \frac{1}{n} R_n + \epsilon.$$

1119 Letting $\epsilon \rightarrow 0$, we see that the term on the right hand side of (30) is no bigger than the term on its left
1120 hand side. Showing the inequality in the other direction is straightforward, since

$$\frac{1}{n} \inf_q \sup_{r \in \Lambda} E_r \log \frac{r(X^n)}{q(X^n)} \leq \frac{1}{n} \sup_{r \in \Lambda} E_r \log \frac{r(X^n)}{q(X^n)},$$

1121 for each probability measure q on \mathbb{N}^∞ . This completes the proof. \square

1122 For the following lemma, recall the definition of $\mathcal{J}(p, \tilde{p})$ for probability distributions p and \tilde{p} on \mathbb{N} ,
1123 made in (11).

1124 **Lemma 37.** Let p and \tilde{p} be probability distributions on \mathbb{N} . Then

$$\frac{\log e}{4} |p - \tilde{p}|_1^2 \leq \mathcal{J}(p, \tilde{p}) \leq |p - \tilde{p}|_1 \log e.$$

1125 If, in addition, p' is a probability distribution on \mathbb{N} , then

$$\mathcal{J}(p, \tilde{p}) + \mathcal{J}(\tilde{p}, p') \geq \mathcal{J}^2(p, p') \frac{1}{8 \log e}.$$

1126 **Proof** The lower bound in the first statement follows from Pinsker's inequality for the KL divergence,
1127 see [21] for example, from which we get

$$D\left(p \parallel \frac{p + \tilde{p}}{2}\right) \geq \frac{\log e}{8} |p - \tilde{p}|_1^2,$$

and similarly for $D\left(\tilde{p} \parallel \frac{p+\tilde{p}}{2}\right)$. For the upper bound in the first statement, since $\log(1+z) \leq z \log e$ for all $z \geq 0$, we may write

$$\begin{aligned} \frac{1}{\log e} \mathcal{J}(p, \tilde{p}) &\leq \sum_{x:p(x) \geq \tilde{p}(x)} p(x) \left(\frac{p(x) - \tilde{p}(x)}{p(x) + \tilde{p}(x)} \right) + \sum_{x:\tilde{p}(x) \geq p(x)} \tilde{p}(x) \left(\frac{\tilde{p}(x) - p(x)}{p(x) + \tilde{p}(x)} \right) \\ &\leq |p - \tilde{p}|_1. \end{aligned}$$

To prove the triangle-like inequality, note that

$$\begin{aligned} \mathcal{J}(p, \tilde{p}) + \mathcal{J}(\tilde{p}, p') &\geq \frac{\log e}{4} (|p - \tilde{p}|_1^2 + |\tilde{p} - p'|_1^2) \\ &\geq \frac{\log e}{8} (|p - \tilde{p}|_1 + |\tilde{p} - p'|_1)^2 \\ &\geq \frac{\log e}{8} (|p - p'|_1)^2 \\ &\geq \frac{1}{8 \log e} \mathcal{J}(p, p')^2, \end{aligned}$$

1128 where the last inequality follows from the upper bound on $\mathcal{J}(p, p')$ already proved in the first part of
1129 the statement. \square

1130 Using Lemma 37, we can prove the following result, which is identical to [2, Lemma 6]. We reproduce
1131 the proof from [2] for completeness.

1132 **Lemma 38.** Let $\epsilon_0 > 0$. If

$$|p_0 - q|_1 \leq \frac{\epsilon_0^2 (\ln 2)^2}{16},$$

1133 then for all $p \in \mathcal{P}$ with $\mathcal{J}(p, p_0) \geq \epsilon_0$, we have

$$\mathcal{J}(p, q) \geq \frac{\epsilon_0^2 \ln 2}{16}. \quad \square$$

1134 **Proof** Since

$$|p_0 - q|_1 \leq \frac{\epsilon_0^2 (\ln 2)^2}{16},$$

1135 Lemma 37 implies that

$$\mathcal{J}(p_0, q) \leq \frac{\epsilon_0^2 \ln 2}{16}.$$

1136 Further, Lemma 37 then implies that

$$\mathcal{J}(p, q) + \frac{\epsilon_0^2 \ln 2}{16} \geq \mathcal{J}(p, q) + \mathcal{J}(p_0, q) \geq \frac{\mathcal{J}^2(p, p_0) \ln 2}{8} \geq \frac{\epsilon_0^2 \ln 2}{8},$$

1137 where the last inequality follows since $\mathcal{J}(p, p_0) \geq \epsilon_0$. This completes the proof. \square

1138 The following result from [2] will be needed to prove the necessity part of Theorem 17.

1139 **Lemma 39.** Fix $\epsilon > 0$. Let p and q be probability distributions on \mathbb{N} with $\mathcal{J}(p, q) \leq \epsilon$. Fix $n \in \mathbb{N}$.
1140 Consider the probability measures on \mathbb{N}^n obtained by *i.i.d.* sampling from p and q respectively, which
1141 we continue to denote by p and q respectively, following our convention.

1142 Suppose $A_n \subset \mathbb{N}^n$ is subset for which $p(A_n) \geq 1 - \alpha$, for some $\alpha > 0$. Then we have

$$q(A_n) > 1 - \alpha - 2n^3 \sqrt{\frac{4\epsilon}{\log e}} - \frac{1}{n}. \quad \square$$

1143 **Proof** Let

$$\mathcal{B}_1 := \left\{ i \in \mathbb{N} : q(i) \leq p(i) \left(1 - \frac{1}{n^2} \right) \right\}, \text{ and } \mathcal{B}_2 := \left\{ i \in \mathbb{N} : p(i) \leq q(i) \left(1 - \frac{1}{n^2} \right) \right\}.$$

1144 Since we have assumed that $\mathcal{J}(p, q) \leq \epsilon$ we have, from Lemma 37, that

$$|p - q|_1 \sqrt{\frac{\log e}{4}} \leq \sqrt{\mathcal{J}(p, q)} \leq \sqrt{\epsilon}.$$

1145 Further, we have

$$|p - q|_1 \geq \sum_{x \in \mathcal{B}_1} (p(x) - q(x)) \geq \frac{p(\mathcal{B}_1)}{n^2} \geq \frac{q(\mathcal{B}_1)}{n^2},$$

1146 and similarly

$$|p - q|_1 \geq \sum_{x \in \mathcal{B}_2} (q(x) - p(x)) \geq \frac{q(\mathcal{B}_2)}{n^2} \geq \frac{p(\mathcal{B}_2)}{n^2}.$$

1147 From the preceding inequalities, it follows that

$$p(\mathcal{B}_1 \cup \mathcal{B}_2) \leq 2n^2 \sqrt{\frac{4\epsilon}{\log e}} \text{ and } q(\mathcal{B}_1 \cup \mathcal{B}_2) \leq 2n^2 \sqrt{\frac{4\epsilon}{\log e}}. \quad (31)$$

1148 Let $S := \mathbb{N} - (\mathcal{B}_1 \cup \mathcal{B}_2)$. For all $x \in S$ we have

$$q(x) \geq p(x) \left(1 - \frac{1}{n^2} \right). \quad (32)$$

1149 In addition, from (31) we have

$$p(S) \geq 1 - 2n^2 \sqrt{\frac{4\epsilon}{\log e}}.$$

1150 Let $S_n \subset \mathbb{N}^n$ denote the set of all length- n strings of symbols from S . Clearly

$$p(S_n) \geq (1 - 2n^2 \sqrt{\frac{4\epsilon}{\log e}})^n > 1 - 2n^3 \sqrt{\frac{4\epsilon}{\log e}}.$$

1151 Thus we have

$$p(A_n \cap S_n) > 1 - 2n^3 \sqrt{\frac{4\epsilon}{\log e}} - \alpha.$$

1152 From (32), for all $x^n \in S_n$, we have

$$q(x^n) \geq p(x^n) \left(1 - \frac{1}{n^2} \right)^n > p(x^n) \left(1 - \frac{1}{n} \right).$$

1153 Therefore,

$$q(A_n) \geq q(A_n \cap S_n) > (1 - 2n^3 \sqrt{\frac{4\epsilon}{\log e}} - \alpha) \left(1 - \frac{1}{n} \right) > 1 - \alpha - 2n^3 \sqrt{\frac{4\epsilon}{\log e}} - \frac{1}{n}. \quad \square$$

APPENDIX III

LENGTH- n PER-SYMBOL REDUNDANCY OF \mathcal{M}_h

1154

1155

1156 We construct a probability measure q^* on \mathbb{N}^∞ such that for \mathcal{M}_h we have

$$\sup_{p \in \mathcal{M}_h^n} \frac{1}{n} D_n(p||q) \leq \frac{2h^{\frac{1}{4}}(\sqrt{h} + 1)}{\sqrt{\ln n}} + \pi \sqrt{\frac{2}{3n}} \log e.$$

1157 This implies that the per-symbol length- n redundancy of \mathcal{M}_h diminishes to 0 as $n \rightarrow \infty$. Hence \mathcal{M}_h is
 1158 strongly compressible.

1159 Consider the probability distribution q on \mathbb{N} defined by $q(i) = 1/i(i+1), i \geq 1$. As observed in
 1160 Example 24, we have

$$\sup_{p \in \mathcal{M}_h} E_p \left(\left\lceil \log \frac{1}{q(X)} \right\rceil \right)^2 < 4(\sqrt{h} + 1)^2. \quad (33)$$

1161 We consider a scheme that encodes patterns [22] of symbols (i.e. natural numbers in our case) first,
 1162 followed by an encoding using $\lceil \log \frac{1}{q(x)} \rceil$ bits to describe every symbol x that appeared in the string,
 1163 in the order in which they arrived. To clarify, recall that the pattern of a sequence of symbols from \mathbb{N}
 1164 replaces each symbol by $k \in \mathbb{N}$ if the symbol was the k -th new symbol to appear in the sequence. For
 1165 example, the pattern of the sequence of natural numbers (2, 3, 17, 4, 3, 3, 1, 2, 4) is (1, 2, 3, 4, 2, 2, 5, 1, 4).
 1166 If in addition to the pattern of a finite sequence of natural numbers, in which there are l distinct symbols,
 1167 one knows which symbol was the k -th symbol to appear for each $1 \leq k \leq l$, one learns the sequence of
 1168 symbols.

1169 The expected (not normalized by n) additional number of bits to encode the pattern of a sequence
 1170 of symbols of length n from any $p \in \mathcal{M}_h$ is at most $\pi \sqrt{\frac{2}{3}} n \log e$, using the results in [22], while the
 1171 expected number of bits to describe the symbols of length- n strings using a prefix code based on the
 1172 probability distribution q on \mathbb{N} is at most

$$\sum_{i \in \mathbb{N}} (1 - (1 - p(i))^n) \lceil \log \frac{1}{q(i)} \rceil.$$

Note that the distinct symbols appearing the the string will need to be specified in the order in which they arrived. Let M_n denote the number of distinct symbols that appear in a sequence of length n . Then the expected number of extra bits the scheme uses for length- n strings is (without normalizing by

n) at most $\pi\sqrt{\frac{2}{3}n} \log e$ plus at most

$$\begin{aligned}
& \sum_{i \in \mathbb{N}} (1 - (1 - p(i))^n) \lceil \log \frac{1}{q(i)} \rceil \\
& \stackrel{(a)}{\leq} \sqrt{\sum_{i \in \mathbb{N}} (1 - (1 - p_i)^n) \sum_{j \in \mathbb{N}} (1 - (1 - p_j)^n) \left(\lceil \log \frac{1}{q(j)} \rceil \right)^2} \\
& \leq \sqrt{\sum_{i \in \mathbb{N}} (1 - (1 - p_i)^n) \sum_{j \in \mathbb{N}} (np_j) \left(\lceil \log \frac{1}{q(j)} \rceil \right)^2} \\
& \stackrel{(b)}{\leq} \sqrt{4(\mathbb{E}M_n)n(\sqrt{h} + 1)^2} \\
& \stackrel{(c)}{\leq} \frac{2nh^{1/4}(\sqrt{h} + 1)}{\sqrt{\ln n}}.
\end{aligned}$$

Here (a) follows from the Cauchy-Schwarz inequality, while (b) follows from (33) and the definition of M_n . As for (c), a result similar to (c) can be found in [23], but we justify (c) below for completeness. We observe that for all $i \in \mathbb{N}$ we have

$$\begin{aligned}
1 - (1 - p_i)^n &= p_i \sum_{j=0}^{n-1} (1 - p_i)^j \\
&\leq p_i \left(\sum_{j=0}^{n-1} (1 - p_i)^j \right) \frac{\sum_{k=1}^n \frac{1}{k}}{\ln n} \\
&\stackrel{(a)}{\leq} \frac{np_i}{\ln n} \sum_{j=0}^{n-1} \frac{(1 - p_i)^j}{j} \\
&\leq \frac{np_i \log \frac{1}{p_i}}{\ln n}.
\end{aligned}$$

1173 Combining the above with the fact that the entropy of any $p \in \mathcal{M}_h$ is at most \sqrt{h} , which was shown
1174 in Example 24, proves (c) in the previous set of equations. In the above set of equations, inequality
1175 (a) follows from Minkowski's inequality which says that if x_i and y_i ($0 \leq i \leq n - 1$) are both
1176 decreasing positive sequences, then $n \sum x_i y_i \geq \sum x_j \sum y_k$. Minkowski's inequality is easily proved by
1177 noting $\sum x_j \sum y_k = \sum_m \sum x_i y_{(i+m) \bmod n}$ and that $\sum x_i y_i \geq \sum x_i y_{(i+m) \bmod n}$ for all $0 \leq m \leq n - 1$.

1178 The claim about the per-symbol length- n redundancy of \mathcal{M}_h follows after normalization by n .

1179 APPENDIX IV

1180 TYPICALITY OF EMPIRICAL DISTRIBUTIONS THAT ARE NOT TOO SPREAD OUT

1181 In this section we prove a useful result quantifying how close the empirical distribution of a sample
1182 drawn *i.i.d.* from a probability distribution p on \mathbb{N} is to p , when the alphabet of symbols showing up in
1183 the sample is not too spread out. There is a lemma that looks somewhat similar in [24]. The difference
1184 of the result in Lemma 40 from that in [24] is that the right side of the inequality in (34) does *not*
1185 depend on p . The result of Lemma 40 will be used in the sufficiency proof in Appendix VI and this
1186 property is crucial for its use.

1187 **Lemma 40.** Let p be any probability distribution on \mathbb{N} . Let $\gamma > 0$ and let $k \geq 2$ be an integer. Let
 1188 X_1^n be a sequence generated *i.i.d.* with marginals p and let $t(X^n)$ be the empirical distribution of X_1^n .
 1189 Then

$$p\left(|t(X^n) - p|_1 > \gamma \text{ and } 2\dot{F}_t^{-1}(1 - \gamma/6) \leq k\right) \leq (2^k - 2) \exp\left(-\frac{n\gamma^2}{18}\right). \quad (34)$$

1190 **Proof** From [25, Proposition 1] we know that for any probability distribution p' on \mathbb{N} with finite
 1191 support of size L we have

$$p'(|t(X^n) - p'|_1 \geq \alpha) \leq (2^L - 2) \exp\left(-\frac{n\alpha^2}{2}\right), \quad (35)$$

1192 where $t(X^n)$ is the type of X^n generated *i.i.d.* with marginal distribution p' .

1193 Consider the probability distributions p' and t' on A obtained from p and t respectively via the
 1194 mapping from \mathbb{N} to $A := \{1, \dots, k-1\} \cup \{-1\}$ that maps i to i for $0 \leq i \leq k-1$ and maps all the other
 1195 natural numbers to -1 . Thus, we have

$$p'(i) = \begin{cases} p(i), & \text{if } 1 \leq i \leq k-1, \\ \sum_{j=k}^{\infty} p(j), & \text{if } i = -1. \end{cases}$$

1196 Further, sequences of natural numbers generated *i.i.d.* with marginal distribution p and with empirical
 1197 distribution t are mapped to sequences from A that are *i.i.d.* with probability distribution p' and have
 1198 empirical distribution t' .

1199 Applying (35) to p' , we have

$$p'(|p' - t'|_1 > \gamma/3) \leq (2^k - 2) \exp\left(-\frac{n\gamma^2}{18}\right). \quad (36)$$

1200 We first argue that all sequences generated by p with empirical distributions t satisfying

$$|p - t|_1 > \gamma \text{ and } 2\dot{F}_t^{-1}(1 - \gamma/6) \leq k$$

1201 are mapped into sequences generated by p' with empirical t' satisfying

$$|p' - t'|_1 > \gamma/3 \text{ and } t'(-1) \leq \gamma/3.$$

This follows from writing

$$\begin{aligned} |p - t|_1 &= \sum_{i=1}^{k-1} |p(i) - t(i)| \\ &\leq \sum_{j=k}^{\infty} (p(j) - t(j)) + 2 \sum_{j=k}^{\infty} t(j) \\ &\leq |p'(-1) - t'(-1)| + \gamma/3, \end{aligned}$$

1202 where the last inequality above follows from the fact that $2\dot{F}_t^{-1}(1 - \gamma/6) \leq k$ implies $F_t(k-1) \geq 1 - \gamma/6$,
 1203 i.e. $\sum_{j=k}^{\infty} t(j) \leq \gamma/6$. Hence we have

$$|p' - t'|_1 = \sum_{i=1}^{k-1} |p(i) - t(i)| + |p'(-1) - t'(-1)| \geq |p - t|_1 - \gamma/3 > \gamma/3,$$

1204 because $|p - t|_1 > \gamma$.

Thus, from (36), we will have

$$\begin{aligned} & p(|t(X^n) - p|_1 > \gamma \text{ and } 2\dot{F}_t^{-1}(1 - \gamma/6) \leq k) \\ & \leq p'(|t' - p'|_1 > \gamma/3 \text{ and } t'(-1) \leq \gamma/3) \\ & \leq (2^k - 2) \exp\left(-\frac{n\gamma^2}{18}\right). \end{aligned}$$

1205 This completes the proof of the lemma. □

1206 APPENDIX V

1207 τ ENTERS WITH PROBABILITY 1

1208 We reproduce the argument from [2] here for completeness.

1209 Every probability distribution $p \in \mathcal{P}$ is contained in at least one of the elements of the cover $(Q_{p,m} \cap$
1210 $\mathcal{P}, \tilde{p} \in \tilde{\mathcal{P}}_m)$, where $Q_{p,m}$ denotes the zone of $\tilde{p} \in \tilde{\mathcal{P}}_m$. Recall the enumeration of $\tilde{\mathcal{P}}_m$. Let p' be centroid
1211 with the smallest index among all centroids in $\tilde{\mathcal{P}}_m$ whose zones contain p . With probability 1, sequences
1212 generated by p will eventually have their type (empirical distribution) entirely within $Q_{p',m}$. (see [26]
1213 for a proof).

1214 Next note that for all n sufficiently large the analog of (26), (which makes sense for all $p' \in \tilde{\mathcal{P}}_m$) will
1215 hold. This follows since the right hand side of (26) diminishes to zero polynomially with n while the
1216 left hand side diminishes to zero exponentially fast in n .

1217 Next, (27) will also hold eventually with probability 1, since, if t denotes the empirical probability of
1218 a sequence generated by p , then

$$\dot{F}_t^{-1}(1 - \sqrt{D_{p'}/6}) \rightarrow \dot{F}_{p'}^{-1}(1 - \sqrt{D_{p'}/6}) \quad (37)$$

with probability 1 as $n \rightarrow \infty$, where we note that the quantity on the left hand side of (37) is actually
a random variable and t determines n . Furthermore, we also have

$$\begin{aligned} 2\dot{F}_p^{-1}(1 - \sqrt{D_{p',m}/6}) &< 3 \left(\sup_{r \in B(p', \epsilon_{p',m}; \mathcal{P})} \dot{F}_r^{-1}(1 - \sqrt{D_{p',m}/6}) \right) \\ &= \log C(p', m), \end{aligned}$$

1219 where the first inequality follows since p is in the $\frac{1}{m}$ -reach of p' .

1220 Therefore, both (26) and (27) will eventually hold with probability 1. Furthermore, long enough
1221 sequences generated by p fall into the zone of p' with probability 1. This implies in turn that $\tau_{\eta,m}$ enters
1222 with probability 1. Note that it is entirely possible that some other probability measure traps strings
1223 before they can be trapped by p' , but that does not take away from the fact that $\tau_{\eta,m}$ will enter with
1224 probability 1.

1225 APPENDIX VI

1226 PROBABILITY OF FALLING INTO BAD TRAPS

1227 Let t be any length- n empirical distribution trapped by \hat{p} , which we recall has $\frac{1}{m}$ -reach $\epsilon_{\hat{p},m}$, such
1228 that $p \notin B(\hat{p}, \epsilon_{\hat{p},m}; \mathcal{P})$. Then we have

$$\mathcal{J}(\hat{p}, p) \geq \epsilon_{\hat{p},m},$$

1229 because $p \notin B(\hat{p}, \epsilon_{\hat{p},m}; \mathcal{P})$, and we have

$$|\hat{p} - t|_1 < \frac{\epsilon_{\hat{p},m}^2 (\ln 2)^2}{16},$$

1230 because t has to be in the zone $Q_{\hat{p},m}$ in order to be captured by \hat{p} . From Lemma 38, which is a
1231 consequence of the pseudo-triangle inequality in Lemma 37, we get

$$\mathcal{J}(p, t) \geq \frac{\epsilon_{\hat{p},m}^2 \ln 2}{16}.$$

1232 Hence, for all types t that are trapped by \hat{p} , by the first part of Lemma 37 we get

$$|p - t|_1^2 \geq \mathcal{J}^2(p, t) (\ln 2)^2 \geq \frac{\epsilon_{\hat{p},m}^4 (\ln 2)^4}{256} = D_{\hat{p},m}^2.$$

This means that for every $p \in \mathcal{P}$, the probability that length- n sequences with empirical distribution t are trapped by a bad \hat{p} can be bounded from above as

$$\begin{aligned} &\leq p \left(|t - p|_1^2 \geq D_{\hat{p},m} \text{ and } 2\dot{F}_t^{-1} \left(1 - \frac{\sqrt{D_{\hat{p},m}}}{6} \right) \leq \log C(\hat{p}, m) \right) \\ &= p \left(|t - p|_1 \geq \sqrt{D_{\hat{p},m}} \text{ and } 2\dot{F}_t^{-1} \left(1 - \sqrt{D_{\hat{p},m}}/6 \right) \leq \log C(\hat{p}, m) \right) \\ &\stackrel{(a)}{\leq} (C(\hat{p}, m) - 2) \exp \left(-\frac{n D_{\hat{p},m}}{18} \right) \\ &\stackrel{(b)}{\leq} \frac{\eta (C(\hat{p}, m) - 2)}{2C(\hat{p}, m) \iota(\hat{p})^2 n(n+1)} \\ &\leq \frac{\eta}{2\iota(\hat{p})^2 n(n+1)}, \end{aligned}$$

1233 where the inequality (a) follows from Lemma 40 and (b) from (26). Therefore, the probability of
1234 sequences falling into bad traps is bounded above by

$$\leq \sum_{n \geq 1} \sum_{\hat{p} \in \tilde{\mathcal{P}}} \frac{\eta}{2\iota(\hat{p})^2 n(n+1)} \leq \frac{\pi^2}{12} \eta < \eta,$$

1235 since $\sum_{\hat{p} \in \tilde{\mathcal{P}}} \frac{1}{\iota(\hat{p})^2} = \frac{\pi^2}{6}$ and $\sum_{n \geq 1} \frac{1}{n(n+1)} = 1$.

1236 APPENDIX VII

1237 A FAKE PROOF

1238 In this section we give a fake proof of the following mistaken claim: *if \mathcal{P}_1 and \mathcal{P}_2 are d.w.c., then*
1239 *$\mathcal{P}_1 \cup \mathcal{P}_2$ is also d.w.c..* We then explain why it is wrong. In the concluding remarks in [2] it was stated,
1240 in passing, that if \mathcal{P}_1 and \mathcal{P}_2 are insurable then $\mathcal{P}_1 \cup \mathcal{P}_2$ is also insurable. This statement is false, for
1241 the reasons explained in this section. This does not affect any of the results in [2].

1242 The argument proceeds as follows. Since \mathcal{P}_i is *d.w.c.* for each $i = 1, 2$, there is a probability measure
1243 q_i on \mathbb{N}^∞ for each $i = 1, 2$ such that for every $m \geq 1$, $0 < 1 - \eta < 1$ and $i = 1, 2$ there is a universal
1244 stopping time $\tau_{\eta,m}^{(i)}$ such that, for all $p \in \mathcal{P}_i$, we have

$$p \left(\exists n \text{ such that } \frac{1}{n} D_n(p || q_i) > \frac{1}{m} \text{ and } \tau_{\eta,m}^{(i)}(X^n) = 1 \right) < \eta.$$

1245 Let $q := (q_1 + q_2)/2$ and, for accuracy $\frac{1}{m} > 0$ and confidence $0 < 1 - \eta < 1$, define

$$\tau_{\eta,m}(\mathbf{x}) := \mathbb{1}(\tau_{\eta,2m}^{(1)}(\mathbf{x}) = 1)\mathbb{1}(\tau_{\eta,2m}^{(2)}(\mathbf{x}) = 1)\mathbb{1}(|\mathbf{x}| > 2m). \quad (38)$$

1246 Now, suppose $p \in \mathcal{P}_1 \cup \mathcal{P}_2$. Without loss of generality, assume that $p \in \mathcal{P}_1$. Now, if $n > 2m$ and we
1247 have

$$\frac{1}{n}D_n\left(p\left\|\frac{q_1 + q_2}{2}\right.\right) > \frac{1}{m},$$

1248 then we have

$$\frac{1}{n}D_n(p\|q_1) > \frac{1}{m} - \frac{1}{n} > \frac{1}{2m}$$

Further, from (38), if $\tau_{\eta,m}(\mathbf{x}) = 1$, then we have $\tau_{\eta,2m}^{(1)}(\mathbf{x}) = 1$ as well. Therefore

$$\begin{aligned} & p\left(\exists n \text{ such that } \frac{1}{n}D_n\left(p\left\|\frac{q_1 + q_2}{2}\right.\right) > \frac{1}{m} \text{ and } \tau_{\eta,m}(X^n) = 1\right) \\ & \leq p\left(\exists n \text{ such that } n > 2m, \frac{1}{n}D_n(p\|q_1) > \frac{1}{2m} \text{ and } \tau_{\eta,2m}^{(1)}(X_1^n) = 1\right) < \eta, \end{aligned}$$

1249 where we have used (38) to see that the event whose probability is being evaluated on the left hand side
1250 of the preceding equation cannot occur unless $n > 2m$. Since the above holds for all $p \in \mathcal{P}_1$ and we can
1251 use a similar argument for all $p \in \mathcal{P}_2$, we are “done”.

1252 The flaw in the above “proof” is that $\tau_{\eta,m}$, as defined in (38), does not necessarily eventually equal
1253 1 almost surely for all sources in $\mathcal{P}_1 \cup \mathcal{P}_2$, which would mean that it is not a universal stopping time
1254 for the model class $\mathcal{P}_1 \cup \mathcal{P}_2$. To see why this issue might arise, note that $\tau_{\eta,2m}^{(i)}$ is known to eventually
1255 equal 1 almost surely only for sources in \mathcal{P}_i . Thus, if it happens to be the case that there is some event
1256 $A \subsetneq \mathbb{N}^\infty$ and $p_1 \in \mathcal{P}_1$ with $p_1(A) > 0$ for which we have $p_2(A) = 0$ for every source $p_2 \in \mathcal{P}_2$, then $\tau_{\eta,2m}^{(2)}$
1257 might never stop waiting on the sequences in A . This doesn’t stop \mathcal{P}_2 from being *d.w.c.*. But when we
1258 introduce sources from \mathcal{P}_1 , in particular p_1 , we find that $\tau_{\eta,m}$, as defined in (38), will never stop waiting
1259 under p_1 . The stopping rule $\tau_{\eta,m}$ would then not be a universal stopping rule for the model class $\mathcal{P}_1 \cup \mathcal{P}_2$.