**Strong Entropy Concentration,
Coding, Game Theory
and Randomness**

Peter Grünwald

# STRONG ENTROPY CONCENTRATION, CODING, GAME THEORY AND RANDOMNESS

PETER GRÜNWALD

ABSTRACT. We give a characterization of Maximum Entropy/Minimum Relative Entropy inference by providing two 'strong entropy concentration' theorems. These theorems unify and generalize Jaynes' 'concentration phenomenon' and Van Campenhout and Cover's 'conditional limit theorem'. The theorems characterize exactly in what sense a 'prior' distribution $Q$ conditioned on a given constraint and the distribution $\tilde{P}$ minimizing $D(P\|Q)$ over all $P$ satisfying the constraint are 'close' to each other. We show how our theorems are related to 'universal models' for exponential families, thereby establishing a link with Rissanen's MDL/stochastic complexity. We then apply our theorems to establish the relationship (A) between entropy concentration and a game-theoretic characterization of Maximum Entropy Inference due to Topsøe and others; (B) between maximum entropy distributions and sequences that are random (in the sense of Martin-Löf/Kolmogorov) with respect to the given constraint. These two applications have strong implications for the use of Maximum Entropy distributions in *sequential prediction tasks*, both for the logarithmic loss and for general loss functions. We identify circumstances under which Maximum Entropy predictions are almost optimal.

## 1. INTRODUCTION

Jaynes' Maximum Entropy (MaxEnt) Principle is a well-known principle for inductive inference [6, 8, 26, 16, 27, 5, 11, 20]. It has been applied to statistical and machine learning problems ranging from protein modeling so stock market prediction [18]. One of its characterizations (some would say 'justifications') is the so-called *concentration phenomenon* [14, 15]. Here is an informal version of this phenomenon, in Jaynes' words:

> "If the information incorporated into the maximum-entropy analysis includes all the constraints actually operating in the random experiment, then the distribution predicted by maximum entropy is overwhelmingly the most likely to be observed experimentally." [17, Page 1124]

For the case in which a prior distribution over the domain at hand is available, Van Campenhout and Cover [27, 5] have proven the related *conditional limit theorem*. In Part I of this paper, we provide a strong generalization of both the concentration phenomenon and the conditional limit theorem. In Part II we apply this. We first show how our theorems can be used to construct universal models for exponential families, thereby establishing a link with Rissanen's Minimum Description Length Principle. We then extend an existing game-theoretic characterization of Maximum Entropy due to Topsøe [26]. Finally we combine the results of Part I with the theory of algorithmic (Martin-Löf/Kolmogorov) randomness. This allows us to substantiate the often-heard informal claim that 'adopting the Maximum Entropy distribution leads to good predictions if the data are random with respect to the given constraint' and to make precise informal notions like '*all* constraints *actually* operating in an environment' as used in Jaynes' statement above. We end by discussing implications of our results in Part II for (sequential) *prediction*. We identify circumstances in which Maximum Entropy distributions lead to almost optimal predictions.

## 2. INFORMAL OVERVIEW

Before we dig into the mathematical details, let us give an informal overview of the results of this paper.

**Maximum Entropy.** Let $X$ be a random variable taking values in some set $\mathcal{X}$, which (only for the time being!) we assume to be finite: $\mathcal{X} = \{1, \ldots, m\}$. Let $P, Q$ be distributions for $\mathcal{X}$ with probability mass functions $p$ and $q$. We define $\mathbf{H}_Q(P)$, *the Q-entropy of P*, as

$$\mathbf{H}_Q(P) = -E_P[\log \frac{p(x)}{q(x)}] = -D(P||Q) \tag{1}$$

In the usual MaxEnt setting[1], we are given a 'prior' distribution $Q$ and a *moment constraint*:

$$E[T(X)] = \tilde{t} \tag{2}$$

where $T$ is some function $T : \mathcal{X} \to \mathbf{R}^k$ for some $k > 0$. We define, if it exists, $\tilde{P}$ to be the unique distribution over $\mathcal{X}$ that maximizes the $Q$-entropy over all distributions (over $\mathcal{X}$)

---

[1]More general formulations with arbitrary convex constraints exist [6], but here we stick to constraints of form (2).

satisfying (2):

$$\tilde{P} = \underset{\{P:E_P[T(X)]=\tilde{t}\}}{\arg\max} \mathbf{H}_Q(P) = \underset{\{P:E_P[T(X)]=\tilde{t}\}}{\arg\min} D(P||Q) \tag{3}$$

The MaxEnt Principle then tells us that, in absence of any further knowledge about the 'true' or 'posterior' distribution according to which data are distributed, our best guess for it is $\tilde{P}$. In practical problems we are usually not given a constraint of form (2). Rather we are given an *empirical constraint* of the form

$$\frac{1}{n}\sum_{i=1}^{n} T(X_i) = \tilde{t} \quad \text{which we always abbreviate to } `\overline{T^{(n)}} = \tilde{t}` \tag{4}$$

The MaxEnt Principle is then usually applied as follows: suppose we are given an empirical constraint of form (4). We then have to make predictions about new data coming from the same source. In absence of knowledge of any 'true' distribution generating this data, we should make our predictions based on the MaxEnt distribution $\tilde{P}$ for the moment constraint (2) corresponding to empirical constraint (4). $\tilde{P}$ is extended to several outcomes by taking the product distribution.

**The Concentration Phenomenon and The Conditional Limit Theorem.** Why should this procedure make any sense? Here is one justification. If $\mathcal{X}$ is finite, and in the absence of any prior knowledge beside the constraint, one usually picks the uniform distribution for $Q$. In this case, Jaynes' 'concentration phenomenon' applies[2]. It says that for all $\epsilon > 0$,

$$Q^n(\sup_{j\in\mathcal{X}} |\frac{1}{n}\sum_{i=1}^{n} I_j(X_i) - \tilde{P}(X=j)| > \epsilon \mid \overline{T^{(n)}} = \tilde{t}) = O(e^{-cn}) \tag{5}$$

for some constant $c$ depending on $\epsilon$. Here $Q^n$ is the $n$-fold product distribution of $Q$, and $I$ is the indicator function: $I_j(x) = 1$ if $x = j$ and 0 otherwise. In words, for the overwhelming majority among the sequences satisfying the constraint, the empirical frequencies are close to the maximum entropy probabilities. It turns out that (5) still holds if $Q$ is non-uniform. For an illustration we refer to Example 1. A closely related result (Theorem 1, [27]) is the Van Campenhout-Cover conditional limit theorem[3], which says that

$$\lim_{\substack{n\to\infty \\ n\tilde{t}\in\mathbf{N}}} Q^1(\cdot \mid \overline{T^{(n)}} = \tilde{t}) = \tilde{P}^1(\cdot) \tag{6}$$

where $Q^1(\cdot \mid \overline{T^{(n)}} = \tilde{t})$ and $\tilde{P}^1(\cdot)$ refer to the marginal distribution of $X_1$ under $Q(\cdot \mid \overline{T^{(n)}} = \tilde{t})$ and $\tilde{P}$ respectively.

---

[2]We are referring here to the version in [14]. The theorem in [15] extends this in a direction different from the one we consider here.

[3]This theorem too has later been extended in several directions different from the one considered here [7]; see Section 4.3.

**Our Results.** Both theorems above say that for some sets $\mathcal{A}$,

$$Q^n(\mathcal{A} \mid \overline{T^{(n)}} = \tilde{t}) \approx \tilde{P}^n(\mathcal{A}) \tag{7}$$

In the concentration phenomenon, the set $\mathcal{A} \subset \mathcal{X}^n$ is about the frequencies of individual outcomes in the sample. In the conditional limit theorem $\mathcal{A} \subset \mathcal{X}^1$ only concerns the first outcome. One might conjecture that (7) holds asymptotically in a much wider sense, namely for *just about any set whose probability one may be interested in*. For examples of such sets see Example 1. In Theorems 1, 2 and 3 we show that (7) indeed holds for a very large class of sets; moreover, we give an explicit indication of the error one makes if one approximates $Q(\mathcal{A} \mid \overline{T^{(n)}} = \tilde{t})$ by $\tilde{P}(\mathcal{A})$. In this way we unify and strengthen both the concentration phenomenon and the conditional limit theorem. To be more precise, let $\{\mathcal{A}_n\}$, with $\mathcal{A}_i \subset \mathcal{X}^i$ be a sequence of 'typical' sets for $\tilde{P}$ in the sense that $\tilde{P}^n(\mathcal{A}_n)$ goes to 1 sufficiently fast. Then broadly speaking theorems 1 and 2 show that $Q^n(\mathcal{A}_n \mid \overline{T^{(n)}} = \tilde{t})$ goes to 1 too, 'almost' as fast as $\tilde{P}^n(\mathcal{A}_n)$. Theorem 3, our main theorem, says that, if $m$ is an arbitrary increasing sequence with $\lim_{n \to \infty} m/n = 0$, then for *every* (measurable) sequence $\{\mathcal{A}_m\}$ (i.e. not just the typical ones), with $\mathcal{A}_m \subset \mathcal{X}^m$, $\tilde{P}^n(\mathcal{A}_m) \to Q^n(\mathcal{A}_m \mid \overline{T^{(n)}} = \tilde{t})$. Then, in part II of the paper, we first connect this to the notion of 'universal models' as arising in the MDL (Minimum Description Length) approach to inductive inference. We next show in what sense our strong concentration phenomena really provide a 'justification', not just a characterization, of MaxEnt. We show first (Theorem 4) that our concentration phenomenon implies that the MaxEnt distribution $\tilde{P}$ *uniquely* achieves the best minimax logarithmic loss achievable for sequential prediction of samples satisfying the constraint. We also show (Theorem 5) that for sequences that are algorithmically *random* relative to the constraint, $\tilde{P}$ achieves good loss also for loss functions other than the logarithmic loss.

## 3. Mathematical Preliminaries

**The Sample Space.** From now on we assume a sample space $\mathcal{X} \subseteq \mathbf{R}^l$ for some $l > 0$ and let $X$ be the random vector with $X(x) = x$ for all $x \in \mathcal{X}$. We reserve the symbol $Q$ to refer to a distribution for $X$ called the *prior distribution* (formally, $Q$ is a distribution over $(\mathcal{X}, \sigma(X))$ where $\sigma(X)$ is the Borel-$\sigma$-algebra generated by $X$). We will be interested in sequences of i.i.d. random variables $X_1, X_2, \ldots$, all distributed according to $Q$. Whenever no confusion can arise, we use $Q$ also to refer to the joint (product) distribution of $\times_{i \in \mathbf{N}} X_i$. Otherwise, we use $Q^m$ to denote the $m$-fold product distribution of $Q$. The sample $(X_1, \ldots, X_m)$ will also be written as $X^{(m)}$.

**The Constraint Functions $T$.** Let $T = (T_{[1]}, \ldots, T_{[k]})$ be a $k$-dimensional random vector that is $\sigma(X)$-measurable. We refer to the event $\{x \in \mathcal{X} \mid T(x) = t\}$ both as '$T(X) = t$' and as '$T = t$'. Similarly we write $T_i = t$ as an abbreviation of $T(X_i) = t$ and $T^{(n)}$ as short for $(T(X_1), \ldots, T(X_n))$. The *average* of $n$ observations of $T$ will be denoted by $\overline{T^{(n)}} := n^{-1} \sum_{i=1}^n T(X_i)$. We assume that the support of $\mathcal{X}$ is either countable (in which case the prior distribution $Q$ admits a probability mass function) or that it is a connected subset of $\mathbf{R}^l$ for some $l > 1$ (in which case we assume that $Q$ has a bounded continuous

density with respect to Lebesgue measure). In both cases, we denote the probability mass function/density by $q$. If $\mathcal{X}$ is countable, we shall further assume that $T$ is of the *lattice form* (which it will be in most applications):

**Definition 1.** [10, Page 490] *A $k$-dimensional lattice random vector $T = (T_{[1]}, \ldots, T_{[k]})$ is a random vector for which there exists real-valued $b_1, \ldots, b_k$ and $h_1, \ldots, h_k$ such that, for $1 \leq j \leq k$, $\forall x \in \mathcal{X} : T_{[j]}(x) \in \{b_j + sh_j \mid s \in \mathbf{N}\}$. We call the largest $h_i$ for which this holds the* span *of $T_{[i]}$.*

If $X$ is continuous, we shall assume that $T$ is 'regular':

**Definition 2.** *We say a $k$-dimensional random vector is of* regular continuous form *if its distribution under $Q$ admits a bounded continuous density with respect to Lebesgue measure.*

**Maximum Entropy.** Throughout the paper, log is used to denote logarithm to base 2. Let $P, Q$ be distributions for $\mathcal{X}$. We define $\mathbf{H}_Q(P)$, *the $Q$-entropy of $P$*, as

$$\mathbf{H}_Q(P) = -D(P\|Q) \tag{8}$$

This is defined even if $P$ or $Q$ have no densities, see [6]. Assume we are given a constraint of form (2), i.e. $E_P[T(X)] = \tilde{t}$. Here $T = (T_{[1]}, \ldots, T_{[k]})$, $\tilde{t} = (\tilde{t}_{[1]}, \ldots, \tilde{t}_{[k]})$. We define, if it exists, $\tilde{P}$ to be the unique distribution over $\mathcal{X}$ that maximizes the $Q$-entropy over all distributions (over $(\mathcal{X}, \sigma(X))$) satisfying (2). That is, $\tilde{P}$ is given by (3). If Condition 1 below holds, then $\tilde{P}$ exists and is given by the exponential form (9), as expressed in the proposition below. In the condition, the notation $a^{\mathrm{T}}b$ refers to the dot product between $a$ and $b$.

**Condition 1:** There exists a $\tilde{\beta} \in \mathbf{R}^k$ such that $Z(\tilde{\beta}) = \int_{x \in \mathcal{X}} \exp(-\tilde{\beta}^{\mathrm{T}}T(x))dQ(x)$ is finite and the distribution $\tilde{P}$ with density (with respect to $Q$)

$$\tilde{p}(x) \quad := \quad \frac{1}{Z(\tilde{\beta})}e^{-\tilde{\beta}^{\mathrm{T}}T(x)} \tag{9}$$

satisfies $E_{\tilde{P}}[T(X)] = \tilde{t}$.

**Proposition 1** ([6]). *Assume Condition 1 holds for Constraint (2). Then it holds for only one $\tilde{\beta} \in \mathbf{R}^k$ and $\inf\{D(P\|Q) \mid P : E_P[T(X)] = \tilde{t}\}$ is attained by (and only by) the $\tilde{P}$ given by (9).*

If Condition 1 holds, then $\tilde{t}$ determines both $\tilde{\beta}$ and $\tilde{P}$. In our theorems, we shall simply assume that Condition 1 holds. A sufficient (by no means necessary!) requirement for Condition 1 is for example that $Q$ has bounded support; see [6] for a more precise characterization. We will also assume in our theorems the following natural condition:

**Condition 2:** The '$T$-covariance matrix' $\Sigma$ with $\Sigma_{ij} = E_{\tilde{P}}[T_{[i]}T_{[j]}] - E_{\tilde{P}}[T_{[i]}]E_{\tilde{P}}[T_{[j]}]$ is invertible.

$\Sigma$ is guaranteed to exist by Condition 1 (see any book with a treatment of exponential families, for example, [19]) and will be singular only if either $\tilde{t}_j$ lies at the boundary of the range of $T_{[j]}$ for some $j$ or if some of the $T_{[j]}$ are affine combinations of the others. In the first case, the constraint $T_{[j]} = \tilde{t}_j$ can be replaced by restricting the sample space to $\{x \in \mathcal{X} \mid T_{[j]}(x) = \tilde{t}_j\}$ and considering the remaining constraints for the new sample space. In the second case, we can remove some of the $T_{[i]}$ from the constraint without changing the set of distributions satisfying it, making $\Sigma$ once again invertible.

## 4. Part I: The Concentration Theorems

### 4.1. The Concentration Phenomenon for Typical Sets.

**Theorem 1. (the concentration phenomenon for typical sets, lattice case)**
*Assume we are given a constraint of form (2) such that $T$ is of the lattice form and $h = (h_1, \ldots, h_k)$ is the span of $T$ and such that conditions 1 and 2 hold. Then there exists a sequence $\{c_i\}$ satisfying*

$$\lim_{n \to \infty} c_n = \frac{\prod_{j=1}^{k} h_j}{\sqrt{(2\pi)^k \det \Sigma}}$$

*such that*

(1) *Let $\mathcal{A}_1, \mathcal{A}_2, \ldots$ be an arbitrary sequence of sets with $\mathcal{A}_i \subset \mathcal{X}^i$. For all $n$ with $Q(T_n = \tilde{t}) > 0$, we have:*

$$\tilde{P}(\mathcal{A}_n) \geq n^{-k/2} c_n Q(\mathcal{A}_n \mid \overline{T^{(n)}} = \tilde{t}). \tag{10}$$

*Hence if $\mathcal{B}_1, \mathcal{B}_2, \ldots$ is a sequence of sets with $\mathcal{B}_i \subset \mathcal{X}^i$ whose probability tends to 1 under $\tilde{P}$ in the sense that $1 - \tilde{P}(\mathcal{B}_n) = O(f(n)n^{-k/2})$ for some function $f : \mathbf{N} \to \mathbf{R}$; $f(n) = o(1)$, then $Q(\mathcal{B}_n | \overline{T^{(n)}} = \tilde{t})$ tends to 1 in the sense that $1 - Q(\mathcal{B}_n | \overline{T^{(n)}} = \tilde{t}) = O(f(n))$.*

(2) *If for all $n$, $\mathcal{A}_n \subseteq \{x^{(n)} \mid n^{-1} \sum_{i=1}^{n} T(x_i) = \tilde{t}\}$ then (10) holds with equality.*

Theorem 1 has applications for coding/compression, Minimum Description Length inference and prediction. These are discussed in Section 5. The proof of Theorem 1 is in Appendix A. It is based on the 'local' central limit theorem for lattice random variables, which says that the probability mass functions (rather than just the distribution functions) of properly scaled sums of $k$-dimensional random vectors converge to the $k$-dimensional normal distribution. The original derivation of the concentration phenomenon [14] used Stirling's approximation of the factorial rather than the local central limit theorem; the connection to the present Theorem 1 is in Section 4.3 below.

**Example 1.** The 'Brandeis dice example' is a toy example frequently used by Jaynes and others in discussions of the MaxEnt formalism [14]. Let $\mathcal{X} = \{1, \ldots, 6\}$ and $X$ be the outcome in one throw of some given die. We initially believe (e.g. for reasons of symmetry) that the distribution of $X$ is uniform. Then $Q(X = j) = 1/6$ for all $j$ and $E_Q[X] = 3.5$. We are then told that the average number of spots is $E[X] = 4.5$ rather than 3.5. As calculated by Jaynes, the MaxEnt distribution $\tilde{P}$ given this constraint is given by

$$(\tilde{p}(1), \ldots, \tilde{p}(6)) = (0.05435, 0.07877, 0.11416, 0.16545, 0.23977, 0.34749). \tag{11}$$

By the Chernoff bound, for every $j \in \mathcal{X}$, every $\epsilon > 0$, $\tilde{P}(|n^{-1} \sum_{i=1}^{n} I_j(X_i) - \tilde{p}(j)| > \epsilon) < 2 \exp(-nc)$ for some constant $c > 0$ depending on $\epsilon$; here $I_j(X)$ is the indicator function for $X = j$. Theorem 1 then implies that $Q(|n^{-1} \sum_{i=1}^{n} I_j(X_i) - \tilde{p}(j)| > \epsilon | \overline{T^{(n)}} = \tilde{t}) = O(\sqrt{n} e^{-nc}) = O(e^{-nc'})$ for some $c' > 0$. In this way we recover Jaynes' original concentration phenomenon (5): the fraction of sequences satisfying the constraint with frequencies close to MaxEnt probabilities $\tilde{p}$ is overwhelmingly large. Suppose now we receive new information about an additional constraint: $P(X = 4) = P(X = 5) = 1/2$. This can be expressed as a moment constraint by $E[(I_4(X), I_5(X))^{\mathrm{T}}] = (0.5, 0.5)^{\mathrm{T}}$. We can now either use $\tilde{P}$ defined as in (11) in the rôle of prior $Q$ and impose the new constraint $E[(I_4(X), I_5(X))^{\mathrm{T}}] = (0.5, 0.5)^{\mathrm{T}}$, or use uniform $Q$ and impose the combined constraint $E[T] = E[(T_{[1]}, T_{[2]}, T_{[3]})^{\mathrm{T}}] = (4.5, 0.5, 0.5)^{\mathrm{T}}$, with $T_{[1]} = X, T_{[2]} = I_4(X), T_{[3]} =$

$I_5(X)$. In both cases we end up with a new MaxEnt distribution $\tilde{\tilde{p}}(4) = \tilde{\tilde{p}}(5) = 1/2$. This distribution, while still consistent with the original constraint $E[X] = 4.5$, rules out the vast majority of sequences satisfying it. However, we can apply our concentration phenomenon again to the new MaxEnt distribution $\tilde{\tilde{P}}$. Let $\mathcal{I}_{j,j',\epsilon}$ denote the event that

$$|\frac{1}{n}\sum_{i=1}^{n} I_j(X_i) - \frac{\sum_{i=1}^{n-1} I_{j'}(X_i)I_j(X_{i+1})}{\sum_{i=1}^{n-1} I_{j'}(X_i)}| > \epsilon.$$

According to $\tilde{\tilde{P}}$, we still have that $X_1, X_2, \ldots$ are i.i.d. Then by the Chernoff bound, for each $\epsilon > 0$, for $j, j' \in \{4, 5\}$, $\tilde{\tilde{P}}(\mathcal{I}_{j,j',\epsilon})$ is exponentially small. Theorem 1 then implies that $Q^n(\mathcal{I}_{j,j'\epsilon} \mid \overline{T^{(n)}} = (4.5, 0.5, 0.5)^{\mathrm{T}})$ is exponentially small too: for the overwhelming majority of samples satisfying the combined constraint, the sample will look just as if it had been generated by an i.i.d. process, even though $X_1, \ldots, X_n$ are obviously not *completely* independent under $Q^n(\cdot|\overline{T^{(n)}} = (4.5, 0.5, 0.5)^{\mathrm{T}})$.

For completeness, we now give a version of Theorem 1 for continuous-valued random vectors. Unfortunately, we cannot use the proof technique used above to compare $\tilde{P}(\mathcal{A}_n)$ to $Q(\mathcal{A}_n|\overline{T^{(n)}} = \tilde{t})$ in the continuous case. The reason is that $\overline{T^{(n)}} = \tilde{t}$ is a set of $Q$-measure 0 (more on this in Section 4.2). Instead, we will condition on $\overline{T^{(n)}}$ being in a small ball around $\tilde{t}$ which we will let shrink to 0 radius as $n$ increases. For $\tilde{t} \in \mathbf{R}^k$, let $\mathcal{B}_\epsilon(\tilde{t}) := \{t \in \mathbf{R}^k \mid \sup_i |t_{[i]} - \tilde{t}_{[i]}| < \epsilon\}$.

**Theorem 2. (the concentration phenomenon for typical sets, continuous case)**
*Assume we are given a constraint of form (2) such that $T$ is of regular continuous form and such that Conditions 1 and 2 hold. Fix some $h > 0$ and let $\epsilon_n := h/n$. Then there exists a sequence $c_1, c_2, \ldots$ satisfying*

$$\lim_{n\to\infty} c_n = e^{-2|\tilde{\beta}|} \frac{h^k}{\sqrt{(2\pi)^k \det \Sigma}}$$

*such that*

(1) *Let $\mathcal{A}_1, \mathcal{A}_2, \ldots$ be an arbitrary sequence of (measurable) sets with $\mathcal{A}_i \subset \mathcal{X}^i$. For all $n$ we have:*
$$\tilde{P}(\mathcal{A}_n) \geq n^{-k/2} c_n Q(\mathcal{A}_n \mid \overline{T^{(n)}} \in \mathcal{B}_{\epsilon_n}(\tilde{t})). \tag{12}$$

(2) *If for all $n$, $\mathcal{A}_n \subseteq \{x^{(n)} \mid n^{-1}\sum_{i=1}^n T(x_i) \in \mathcal{B}_{\epsilon_n}(\tilde{t})\}$ then*
$$\lim_{n\to\infty} \tilde{P}(\mathcal{A}_n) \leq e^{2|\tilde{\beta}|} n^{-k/2} c_n Q(\mathcal{A}_n \mid \overline{T^{(n)}} \in \mathcal{B}_{\epsilon_n}(\tilde{t})).$$

4.2. **The Strong Concentration Phenomenon.** There are a few limitations to Theorems 1 and 2: (1) we must require that $\tilde{P}(\mathcal{A}_n)$ goes to 0 or 1 as $n \to \infty$; (2) the continuous case needed a separate statement, which is caused by the more fundamental (3) the proof technique used cannot be adapted to point-wise conditioning on $\overline{T^{(n)}} = \tilde{t}$ in the continuous case. Theorem 3 overcomes all these problems. The price we pay is that, when conditioning on $\overline{T^{(n)}} = \tilde{t}$, the sets $\mathcal{A}_m$ must only refer to $X_1, \ldots, X_m$ where $m$ is such that $m/n \to 0$; for example, $m = \lceil n/\log n \rceil$ will work. Whenever we write $Q(\cdot \mid \overline{T^{(n)}} = t)$ or $\tilde{P}(\cdot \mid \overline{T^{(n)}} = t)$ we refer to the continuous version of these quantities. These exist by Proposition 2 in Appendix B. Recall that (for $m < n$) $Q^m(\cdot \mid \overline{T^{(n)}} = \tilde{t})$ refers to the marginal distribution of $X_1, \ldots, X_m$ conditioned on $\overline{T^{(n)}} = \tilde{t}$. It is implicitly

understood in the theorem that in the lattice case, $n$ ranges only over those values for which $Q(\overline{T^{(n)} = \tilde{t}}) > 0$.

**Theorem 3. (Main Theorem: the Strong Concentration Phenomenon/ Strong Conditional Limit Theorem)** *Let $\{m_i\}$ be an increasing sequence with $m_i \in \mathbf{N}$, such that $\lim_{n \to \infty} m_n/n = 0$. Assume we are given a constraint of form (2) such that $T$ is of the regular continuous form or of the lattice form and suppose that Conditions 1 and 2 are satisfied. Then as $n \to \infty$, $Q^{m_n}(\cdot \mid \overline{T^{(n)} = \tilde{t}})$ converges weakly[4] to $\tilde{P}^{m_n}(\cdot)$.*

The proof (using the same key idea, but involving much more work than the proof of Theorem 1) is in Appendix B.

**4.3. Related Results.** Theorem 1 is related to Jaynes' original concentration phenomenon, the proof of which is based on Stirling's approximation of the factorial. Another closely related result (also based on Stirling's approximation) is in Example 5.5.8 of [21]. Both results can be easily extended to prove the following weaker version of Theorem 1, item 1: $\tilde{P}(\mathcal{A}_n) \geq n^{-|\mathcal{X}|} c_n Q(\mathcal{A}_n \overline{| T^{(n)} = \tilde{t}})$ where $c_n$ tends to some constant. Note that in this form, the theorem is void for infinite sample spaces. It also cannot be applied to prove (weaker) analogues of Theorem 2. In [15] the original concentration phenomenon is extended in a direction somewhat different from Theorem 1; it would be interesting to study the relations.

Theorem 3 is similar to the original 'conditional limit theorems' (Theorems 1 and 2) of Van Campenhout and Cover [27]. We note that the preconditions for our theorem to hold are weaker and the conclusion is stronger than for the original conditional limit theorems: our theorem is a generalization of theirs which supplies us with an explicit bound on how fast $m$ can grow as $n$ tends to infinity. The conditional limit theorem was later extended by Csiszár [7]. His setting is considerably more general than ours (e.g. allowing for general convex constraints rather than just moment constraints), but his results also lack an explicit estimate of the speed at which $m$ can increase with $n$. Csiszár [7] and Cover and Thomas [5] (where a simplified version of the conditional limit theorem is proved) both make the connection to large deviation results, in particular Sanov's theorem. As shown in the latter reference, weak versions of the conditional limit theorem can be interpreted as immediate consequences of Sanov's theorem.

## 5. Part II: Applications

For simplicity we restrict ourselves in this section to countable sample spaces $\mathcal{X}$ and we identify probability mass functions with probability distributions. Subsections 5.1 and 5.2 make frequent use of coding-theoretic concepts which we now briefly review (Sections 5.3 and 6 can be read without knowledge of coding/information theory).

Recall that by the Kraft Inequality [5], for every prefix code with lengths $L$ over symbols from a countable alphabet $\mathcal{X}^n$, there exists a (possibly sub-additive) probability mass function $p$ over $\mathcal{X}^n$ such that for all $x^{(n)} \in \mathcal{X}^n$, $L(x^{(n)}) = -\log p(x^{(n)})$. We will call this $p$ the 'probability (mass) function corresponding to $L$'. Similarly, for every probability mass function $p$ over $\mathcal{X}^n$ there exists a (prefix) code with lengths $L(x^{(n)}) = \lceil -\log p(x^{(n)}) \rceil$. Neglecting the round-off error, we will simply say that for every $p$, there exists a code with

---

[4]That is, for all sequences $\{\mathcal{A}_m\}$ where each $\mathcal{A}_m$ is a measurable continuity set $\mathcal{A}_m \subseteq \mathcal{X}^m$, $Q^{m_n}(\mathcal{A}_{m_n} \mid \overline{T^{(n)} = \tilde{t}}) \to \tilde{P}^{m_n}(\mathcal{A}_{m_n})$. A 'continuity set' $\mathcal{A}_m$ is a set such that the $\tilde{P}^m$-probability of the *boundary* of the set $\mathcal{A}_m$ is 0; in our case, *all* measurable sets $\mathcal{A}_m$ are continuity sets. See Theorem 2.1 of [4].

lengths $L(x^{(n)}) = -\log p(x^{(n)})$. We call the code with these lengths 'the code corresponding to $p$'. By the information inequality [5], this is also the most efficient code to use if data $X^{(n)}$ were actually distributed according to $p$.

We can now see that Theorem 1, item 2, has important implications for coding. Consider the following special case of Theorem 1, which obtains by taking $\mathcal{A}_n = \{x^{(n)}\}$ and logarithms:

**Corollary 1. (the concentration phenomenon, coding-theoretic formulation)**
*Assume we are given a constraint of form (2) such that $T$ is of the lattice form and $h = (h_1, \ldots, h_k)$ is the span of $T$ and such that conditions 1 and 2 hold. For all $n$, all $x^{(n)}$ with $n^{-1} \sum_{i=1}^n T(x_i) = \tilde{t}$, we have*

$$-\log \tilde{p}(x^{(n)}) =$$

$$-\log q(x^{(n)} \mid \frac{1}{n}\sum_{i=1}^n T(X_i) = \tilde{t}) + \frac{k}{2}\log 2\pi n + \log\sqrt{\det\Sigma} - \sum_{j=1}^k \log h_j + o(1) =$$

$$-\log q(x^{(n)} \mid \frac{1}{n}\sum_{i=1}^n T(X_i) = \tilde{t}) + \frac{k}{2}\log n + O(1). \quad (13)$$

In words, this means the following: let $x^{(n)}$ be a sample distributed according to $Q$, Suppose we are given the information that $n^{-1}\sum_{i=1}^n T(x_i) = \tilde{t}$. Then, by the information inequality, the most efficient code to encode $x^{(n)}$ is the one based on $q(\cdot | \overline{T^{(n)}} = \tilde{t})$ with lengths $-\log q(x^{(n)} | \overline{T^{(n)}} = \tilde{t})$. Yet if we encode $x^{(n)}$ using the code with lengths $-\log \tilde{p}(\cdot)$ (which would be the most efficient had $x^{(n)}$ been generated by $\tilde{p}$) then the number of extra bits we need is only of the order $(k/2)\log n$. That means, for example, that the number of additional bits we need *per outcome* goes to 0 as $n$ increases. These and other consequences of Corollary 1 will be exploited in the next three subsections.

### 5.1. Connection to MDL, Stochastic Complexity, Two-Part Codes. *Universal Models* play a fundamental rôle in modern versions of the MDL (Minimum Description Length) approach to inductive inference and model selection [2, 24]. For details about universal models and codes as well as all coding-theoretic concepts appearing in this section, we refer to [2]. The material in the present section is not needed to understand later sections.

Let $\mathcal{M}_k = \{P_\theta(\cdot) \mid \theta \in \Gamma_k\}$, where $\Gamma_k \subseteq \mathbf{R}^k$ is a $k$-dimensional parametric class of i.i.d. distributions for sample space $\mathcal{X}$. Let $C$ be a code for alphabet $\mathcal{X}^n$, with lengths $L_C$ and define the *regret* $R_C(\cdot)$ such that for all $x^{(n)}$,

$$L_C(x^{(n)}) = -\log p_{\hat{\theta}(x^{(n)})}(x^{(n)}) + R_C(x^{(n)}),$$

where $\hat{\theta}(x^{(n)})$ is the (ML) Maximum Likelihood estimator in $\mathcal{M}_k$ for data $x^{(n)}$, assumed to exist. Roughly speaking, a *universal code* for sequences of length $n$ is a code $C$ such that the regret $r_C(x^{(n)})$ is small uniformly for all or (in some sense) 'most' $x^{(n)}$. A *universal model* is the probability distribution corresponding to a universal code.

It is well-known [2, 24] that, under mild regularity conditions, there exist universal codes $C$ for $\mathcal{M}_k$ with lengths $L_C(x^{(n)}) = -\log p_{\hat{\theta}(x^{(n)})}(x^{(n)}) + \frac{k}{2}\log n + O(1)$, leading to

regret

$$R_C(x^{(n)}) = \frac{k}{2} \log n + O(1) \tag{14}$$

Usually (14) holds uniformly for all sequences $x_1, x_2, \ldots$. (we sometimes need to restrict ourselves to a compact subset of $\Gamma_k$ in order to make (14) uniformly true). It is also known that (14) is in some sense (up to $O(1)$) the best regret that can be achieved [22, 23]. Therefore, *every* code that achieves (14) is usually called a 'universal code', and its corresponding distribution 'universal model'. Until very recently there were four known ways to construct a universal model for a given class $\mathcal{M}_k$: the two-part code, the Bayesian mixture-code , the Shtarkov-normalized-maximum-likelihood (NML) code and the predictive or 'prequential' code, see [2]. These four methods, while superficially very different, all share the same asymptotic lengths (14). Under further regularity conditions on $\mathcal{M}_k$ and if the code $C$ that is used is allowed to depend on sample size $n$, (14) the Shtarkov-NML and two-part codes can be refined to give [2]:

$$R_C(x^{(n)}) = \frac{k}{2} \log \frac{n}{2\pi} + \log \int_{\Gamma_k} \sqrt{\det I(\theta)} d\theta + o(1), \tag{15}$$

where $I(\theta)$ is the (expected) Fisher information matrix of $\theta$. Quite recently, Rissanen [25] showed that the regret (15) is the best that can be achieved under at least three different definitions of optimality. $L_C(x^{(n)}) = - \log p_{\hat{\theta}(x^{(n)})}(x^{(n)}) + R_C(x^{(n)})$, with $R_C(x^{(n)})$ given by (15), is called the 'stochastic complexity of $x^{(n)}$ relative to $\mathcal{M}_k$'.

In the same recent reference [25], Rissanen implicitly introduced a new type of universal code that achieves regret (15). We illustrate this kind of code for the simple case where $\mathcal{M}_k$ is a $k$-dimensional exponential family with finite sample space $\mathcal{X}$. Let then $\mathcal{M}_k = \{P_\theta(\cdot) \mid \theta \in \Gamma_k\}$ be a $k$-parameter exponential family for $\mathcal{X}$ with $\Gamma_k$ the mean-value parameter space, $q$ the background measure and sufficient statistic $T = (T_{[1]}, \ldots, T_{[k]})$. Then $p_\theta = \tilde{p}$ with $\tilde{p}$ given by (9), and $E_{p_\theta}[T] = \tilde{t} = \theta$.

We will encode $x^{(n)}$ in a way similar to (but, as we shall see, still essentially different from) the two-part coding technique [2]: we first code (describe) a distribution for $\mathcal{X}^n$ and then code the data 'with the help of' this distribution. In our case, for data $x^{(n)}$, we first encode the ML estimator $\hat{\theta}(x^{(n)})$ using some code $C_1$ with lengths $L_1$. We then encode $x^{(n)}$ itself using some code $C_2$, making use of the fact that its ML estimator is $\hat{\theta}(x^{(n)})$. By the Kraft inequality this can be done using $L_2(x^{(n)} \mid \hat{\theta}(x^{(n)})) = -\log q(x^{(n)} \mid \overline{T^{(n)}} = \hat{\theta}(x^{(n)})) = -\log q(x^{(n)} \mid \hat{\theta}(x^{(n)}))$ bits. This leads to a code $C^*$ that allows us to encode all $x^{(n)} \in \mathcal{X}^n$ by concatenating the codewords of $\hat{\theta}(x^{(n)})$ (under $C_1$) and $x^{(n)} | \hat{\theta}(x^{(n)})$ (under $C_2$).

Since $\mathcal{X}$ is finite, $n^{-1} \sum_{i=1}^{n} T_{[j]}(X_i)$ can only take on $n \cdot |\mathcal{X}|$ distinct values. Therefore, we can choose $C_1$ such that $L_1(\hat{\theta}(x^{(n)})) = k \log n + k \log |\mathcal{X}|$. By Corollary 1 the code $L_2(\cdot | \cdot)$ has lengths

$$L_2(x^{(n)} \mid \hat{\theta}(x^{(n)})) = -\log p_{\hat{\theta}(x^{(n)})}(x^{(n)}) - \frac{k}{2} \log n - O(1). \tag{16}$$

Summing $L_1$ and $L_2$, we see that the total code length $L^*(x^{(n)})$ for *arbitrary* $x^{(n)}$ is bounded by $\frac{k}{2} \log n - \log p_{\hat{\theta}(x^{(n)})}(x^{(n)}) + O(1)$. Therefore, the regret satisfies (14) which suggests that $C^*$ is a universal code for $\mathcal{M}_k$. Indeed, we can refine $C^*$ by changing $C_1$

as follows. Let $\Gamma_k^{(n)} := \{\theta \in \Gamma_k \mid \exists x^{(n)} \in \mathcal{X}^n : \theta = \hat{\theta}(x^{(n)})\}$. Rissanen [25] defined the
*canonical prior* $w_n(\hat{\theta})$ as the following probability mass function over $\Gamma_k^{(n)}$:

$$w_n(\theta) := \frac{p_\theta(\overline{T^{(n)} = \theta})}{\sum_{\theta \in \Gamma_k^{(n)}} p_\theta(\overline{T^{(n)} = \theta})} \tag{17}$$

For a sample of length $n$, we let $C_1$ be the code with lengths $L_1(\theta) = -\log w_n(\theta)$, and
we leave $C_2$ unchanged. Then (after some algebra) $L^*(x^{(n)}) = -\log p_{\hat{\theta}(x^{(n)})}(x^{(n)}) + c_n$
where $c_n$ is a 'constant' (still depending on $n$, but not depending on $x^{(n)}$). Moreover, $C^*$
is complete (i.e. it satisfies the Kraft inequality with equality). Therefore it must be a
universal code achieving regret (15). By inspecting (13) and realizing that for exponential
families, the Fisher information $I(\theta)$ is the inverse of $\Sigma$, the $T$-covariance matrix of $p_\theta$
[19], one finds (after quite some algebra) that asymptotically, for $\theta \in \Gamma_k^{(n)}$, $w_n$ becomes a
'discretized Jeffreys' prior'

$$w_n(\theta) \sim \frac{\sqrt{I(\theta)}}{\sum_{\theta \in \Gamma_k^{(n)}} \sqrt{I(\theta)}}.$$

We will omit the details of the argument. Note that the 2-part code described above is
quite different from the usual 2-part code. In that approach, instead of the ML estimator
itself one encodes the ML estimator truncated to a coarser precision of $(1/2)\log n + O(1)$
bits per parameter. Then data $x^{(n)}$ are encoded using the code based on $p_{\bar{\theta}(x^{(n)})}$, where $\bar{\theta}$
is the truncated version of $\hat{\theta}$. Instead of using the unconditional truncated $p_{\bar{\theta}}$, we use the
*conditional, but un-truncated* $p_{\hat{\theta}(x^{(n)})}(\cdot|\overline{T^{(n)} = \hat{\theta}(x^{(n)})}) = q(\cdot|\overline{T^{(n)} = \hat{\theta}(x^{(n)})})$.

Other relations between MDL and Maximum Entropy have been investigated by Feder
[9] and Li and Vitányi [21]. In the next section we will see how Theorem 1 leads to yet
another relation between minimum code length and Maximum Entropy.

## 5.2. Empirical Constraints and Game Theory.
From now on we will only work with
countable $\mathcal{X}$. The $\sigma$-algebra of such $\mathcal{X}$ is always tacitly taken to be the power set of $\mathcal{X}$.
The $\sigma$-algebra thus being implicitly understood, we can define $\mathcal{P}(\mathcal{X})$ to be the set of all
probability distributions over $\mathcal{X}$. For a product $\mathcal{X}^\infty = \times_{i \in \mathbf{N}} \mathcal{X}$ of a countable sample
space $\mathcal{X}$, we define $\mathcal{P}(\mathcal{X}^\infty)$ be the set of all distributions over the product space with the
associated product $\sigma$-algebra.

In [26, 11], a characterization of Maximum Entropy distributions quite different from
the present one was given. It was shown that, under regularity conditions,

$$\mathbf{H}_q(\tilde{p}) = \sup_{p^*: E_{p^*}[T] = \tilde{t}} \inf_p E_{p^*}\left[-\log \frac{p(X)}{q(X)}\right] = \inf_p \sup_{p^*: E_{p^*}[T] = \tilde{t}} E_{p^*}\left[-\log \frac{p(X)}{q(X)}\right] \tag{18}$$

where both $p$ and $p^*$ are understood to be members of $\mathcal{P}(\mathcal{X})$ and $\mathbf{H}_q(\tilde{p})$ is defined as in
(1). By this result, the MaxEnt setting can be thought of as a game between Nature, who
can choose any $p^*$ satisfying the constraint, and Statistician, who only knows that Nature
will choose a $p^*$ satisfying the constraint. Statistician wants to minimize his worst-case
expected logarithmic loss (relative to $q$), where the worst-case is over all choices for Nature.
It turns out that the minimax strategy for Statistician in (18) is given by $\tilde{p}$. That is,

$$\tilde{p} = \arg\inf_p \sup_{p^*: E_{p^*}[T] = \tilde{t}} E_{p^*}\left[-\log \frac{p(x)}{q(x)}\right]. \tag{19}$$

This gives a decision-theoretic justification of using MaxEnt probabilities which seems quite different from our concentration phenomenon. Or is it? Realizing that in practical situations we deal with empirical constraints of form (4) rather than (2) we may wonder what distribution $\hat{p}$ is minimax in the empirical version of problem (19). In this version Nature gets to choose an individual sequence rather than a distribution[5]. To make this precise, let

$$\mathcal{C}_n = \{x^{(n)} \in \mathcal{X}^n \mid n^{-1} \sum_{i=1}^{n} T(x_i) = \tilde{t}\}. \tag{20}$$

Then, for $n$ with $\mathcal{C}_n \neq \emptyset$, $\hat{p}_n$ (if it exists) is defined by

$$\hat{p}_n := \arg \inf_{p \in \mathcal{P}(\mathcal{X}^n)} \sup_{x^{(n)} \in \mathcal{C}_n} -\log \frac{p(x_1, \ldots, x_n)}{q(x_1, \ldots, x_n)} = \arg \sup_{p \in \mathcal{P}(\mathcal{X}^n)} \inf_{x^{(n)} \in \mathcal{C}_n} \frac{p(x_1, \ldots, x_n)}{q(x_1, \ldots, x_n)} \tag{21}$$

$\hat{p}_n$ can be interpreted in two ways: (1) it is the distribution that assigns 'maximum probability' (relative to $q$) to all sequences satisfying the constraint; (2) since $-\log(\hat{p}(x^{(n)})/q(x^{(n)})) = \sum_{i=1}^{n}(-\log \hat{p}(x_i|x_1, \ldots, x_{i-1}) + \log q(x_i|x_1, \ldots, x_{i-1}))$, it is also the $p$ that minimizes cumulative worst-case logarithmic loss relative to $q$ when used for sequentially predicting $x_1, \ldots, x_n$.

One immediately verifies that $\hat{p}_n = q^n(\cdot \mid \overline{T^{(n)}} = \tilde{t})$: the solution to the empirical minimax problem is just the conditioned prior, which we know by Theorems 1 and 3 is in some sense very close to $\tilde{p}$. However, for no single $n$, $\tilde{p}$ is exactly equal to $q^n(\cdot \mid \overline{T^{(n)}} = \tilde{t})$. Indeed, $q^n(\cdot \mid \overline{T^{(n)}} = \tilde{t})$ assigns zero probability to any sequence of length $n$ not satisfying the constraint. This means that using $q$ in prediction tasks against the logarithmic loss will be problematic if the constraint only holds approximately (as we will discuss in more detail in the journal version of this paper) and/or if $n$ is unknown in advance. In the latter case, it is impossible to use $q(\cdot \mid \overline{T^{(n)}} = \tilde{t})$ for prediction without modification. The reason is that there exist sequences $x^{(n_2)}$ of length $n_2 > n_1$ satisfying the constraint such that $q(x^{(n_2)}|x^{(n_1)} \in \mathcal{C}_{n_1}) = 0$. We may guess that in this case ($n$ not known in advance), the MaxEnt distribution $\tilde{p}$, rather than $q(\cdot|\overline{T^{(n)}} = \tilde{t})$ is actually the optimal distribution to use for prediction. The following theorem shows that this is indeed so:

**Theorem 4.** *Let $\mathcal{X}$ be a countable sample space. Assume we are given a constraint of form (2) such that $T$ is of the lattice form, and such that Conditions 1 and 2 are satisfied. Let $\mathcal{C}_n$ be as in (20). Then the infimum in*

$$\inf_{p \in \mathcal{P}(\mathcal{X}^\infty)} \sup_{\{n \,:\, \mathcal{C}_n \neq \emptyset\}} \sup_{x^{(n)} \in \mathcal{C}_n} -\frac{1}{n} \log \frac{p(x_1, \ldots, x_n)}{q(x_1, \ldots, x_n)} \tag{22}$$

*is achieved by the Maximum Entropy distribution $\tilde{p}$, and is equal to $\mathbf{H}_q(\tilde{p})$.*

*Proof.* Let $\mathcal{C} = \cup_{i=1}^{\infty} \mathcal{C}_i$. We need to show that for all $n$, for all $x^{(n)} \in \mathcal{C}$,

$$\mathbf{H}_q(\tilde{p}) = -\frac{1}{n} \log \frac{\tilde{p}(x^{(n)})}{q(x^{(n)})} = \inf_{p \in \mathcal{P}(\mathcal{X}^\infty)} \sup_{\{n \,:\, \mathcal{C}_n \neq \emptyset\}} \sup_{x^{(n)} \in \mathcal{C}_n} -\frac{1}{n} \log \frac{p(x^{(n)})}{q(x^{(n)})} \tag{23}$$

Equation (23) implies that $\tilde{p}$ reaches the inf in (22) and that the inf is equal to $\mathbf{H}_q(\tilde{p})$. The leftmost equality in (23) is a standard result about exponential families of form (9); see for example, [12, Proposition 4.1] or [24]. To prove the rightmost equality in (23), let

---

[5]To our knowledge, we are the first to analyze this 'empirical' game.

$x^{(n)} \in \mathcal{C}_n$. Consider the conditional distribution $q(\cdot \mid x^{(n)} \in \mathcal{C}_n)$. Note that, for every distribution $p_0$ over $\mathcal{X}^n$, $p_0(x^{(n)}) \leq q(x^{(n)}|x^{(n)} \in \mathcal{C}_n)$ for at least one $x^{(n)} \in \mathcal{C}_n$. By Theorem 1 (or rather Corollary 1), for this $x^{(n)}$ we have

$$-\frac{1}{n}\log\frac{p_0(x^{(n)})}{q(x^{(n)})} \geq -\frac{1}{n}\log\frac{\tilde{p}(x^{(n)})}{q(x^{(n)})} - \frac{k}{2n}\log n - O(\frac{1}{n}),$$

and we see that for every distribution $p_0$ over $\mathcal{X}^\infty$,

$$\sup_{\{n\,:\,\mathcal{C}_n\neq\emptyset\}} \sup_{x^{(n)}\in\mathcal{C}_n} -\frac{1}{n}\log\frac{p_0(x^{(n)})}{q(x^{(n)})} \geq \sup_{\{n\,:\,\mathcal{C}_n\neq\emptyset\}} \sup_{x^{(n)}\in\mathcal{C}_n} -\frac{1}{n}\log\frac{\tilde{p}(x^{(n)})}{q(x^{(n)})},$$

which shows the rightmost equality in (23). $\qquad\qquad\qquad\qquad\qquad\qquad\square$

### 5.3. Maximum Entropy and Algorithmic Randomness.

In the algorithmic theory of randomness, [21], one (broadly speaking) identifies *randomness of individual sequences* with incompressibility of such sequences. In this section we show that a sequence that is 'random relative to a given constraint' is 'almost' random with respect to the MaxEnt distribution $\tilde{P}$ for the constraint. The reader who is not familiar with Martin-Löf randomness is urged to move on to Theorem 5 which demonstrates the consequences of this fact for prediction based on MaxEnt distributions.

Throughout this section we assume $\mathcal{X}$ to be finite and $Q$ to be uniform, so maximizing the entropy reduces to the 'original' Maximum (Shannon) Entropy formalism. Let $\mathcal{U} := \cup_{i=1}^\infty \mathcal{X}^i$. For $\mathbf{x}, \mathbf{y} \in \mathcal{U}$, $K(\mathbf{x}|\mathbf{y})$ will stand for the prefix Kolmogorov complexity of sequence $\mathbf{x}$ conditional on $\mathbf{y}$; $K(\mathbf{x})$ stands for $K(\mathbf{x}|\lambda)$ where $\lambda$ is the empty sequence. For a finite set $\mathcal{C} \subset \mathcal{U}$, $K(\mathbf{x}|\mathcal{C})$ is the prefix complexity of $\mathbf{x}$ conditional on $\mathbf{x} \in \mathcal{C}$. Kolmogorov complexity is defined here with respect to some fixed universal reference prefix Turing Machine. For precise definitions of all these concepts, see Section 3.1 and Exercise 2.2.12. of [21].

**Theorem** (**Theorem 3.6.1 and Corollary 4.5.2 of** [21]) *An infinite sequence* $(x_1, x_2, \ldots) \in \mathcal{X}^\infty$ *is Martin-Löf random with respect to the uniform distribution iff there exists a constant $c$ such that for all $n$, $K(x_1, \ldots, x_n) \geq n - c$.*

Here, we take this characterization of Martin-Löf randomness as basic. We will extend the notion of randomness to *sequences conditional on constraints* in an obvious manner. Let $\{\mathcal{C}_n\}$ be a sequence of constraints, where $\mathcal{C}_n \subseteq \mathcal{X}^n$ (we identify constraints with the set of sequences satisfying them). The theorem above suggests the following definition:

**Definition 3.** *An infinite sequence* $(x_1, x_2, \ldots) \in \mathcal{X}^\infty$ *is called random with respect to the sequence of constraints $\{\mathcal{C}_n\}$ (relative to the uniform distribution) iff there exists a constant $c$ such that for all $n$ with $\mathcal{C}_n \neq \emptyset$, we have $K(x^{(n)}|\mathcal{C}_n) \geq \log |\mathcal{C}_n| - c$.*

In our situation, the constraint is of form (20). Because of this simple form and since $\mathcal{X}$ is finite, there exists a fixed-length program that, for each $n$, when input $\langle n, \mathbf{x}\rangle$ with $\mathbf{x} \in \mathcal{X}^n$, outputs 1 iff $\mathbf{x} \in \mathcal{C}_n$ and 0 otherwise. Therefore the definition reduces to $(x_1, x_2, \ldots)$ *is random iff* $\exists c\forall n : \mathcal{C}_n \neq \emptyset \Rightarrow K(x^{(n)}|n) \geq \log |\mathcal{C}_n| - c$.

By Theorem 1, if $(x_1, x_2, \ldots)$ is random with respect to the constraints $\{\mathcal{C}_n\}$, then for all $x^{(n)} \in \mathcal{C}_n$,

$$K(x^{(n)}|n) \geq \log |\mathcal{C}_n| - O(1) = -\log\tilde{p}(x^{(n)}) - \frac{k}{2}\log n - O(1). \qquad (24)$$

In words: (see Corollary 4.5.2 of [21]) *If $(x_1, x_2, \ldots)$ is random with respect to the constraints $\{\mathcal{C}_n\}$ (relative to the uniform distribution) then $(x_1, x_2, \ldots)$ is 'almost' Martin-Löf random with respect to the maximum entropy distribution $\tilde{p}$.*

Equation 24 suggests that for the overwhelming majority of sequences satisfying the constraint (namely, those that are random with respect to the constraint), sequentially predicting outcomes in the sequence on the basis of the MaxEnt distribution leads to almost optimal results, no matter what loss function we use. The following theorem shows that this is indeed so. It holds for general prior distributions $Q$ and is proved in Appendix C. Consider a loss function $\textsc{loss} : \mathcal{X} \times \Delta \to [0, \infty]$ where $\Delta$ is some space of *predictions* or *decisions*. A *prediction (decision) strategy* $\delta^*$ is a function $\delta^* : \cup_{i=0}^{\infty} \mathcal{X}^i \to \Delta$. $\delta^*(x_1, \ldots, x_n)$ is to be read as 'the prediction/decision for $X_{n+1}$ based on initial data $(x_1, \ldots, x_n)$'. We assume

**Condition 3.** $\mathcal{X}$ is finite. $\textsc{loss}(x; \cdot)$ is continuous in its second argument for all $x \in \mathcal{X}$. $\Delta$ is a compact convex subspace of $\mathbf{R}^l$ for some $l > 0$.

Under this condition, there exists at least one $\delta$ attaining $\inf E_{\tilde{P}}[\textsc{loss}(X; \delta)]$. Fix any such optimal (under $\tilde{P}$) decision and denote it $\tilde{\delta}$.

**Theorem 5.** *Suppose that $T$ is of lattice form and suppose Conditions 1, 2 and 3 hold. Then (letting $n$ reach over all numbers such that $Q(\overline{T^{(n)}} = \tilde{t}) > 0$), for all decision strategies $\delta^*$, for all $\epsilon > 0$, there exists a $c > 0$ such that*

$$Q(\frac{1}{n}(\sum_{i=1}^{n} \textsc{loss}(x_i; \tilde{\delta}) - \sum_{i=1}^{n} \textsc{loss}(x_i; \delta^*(x_1, \ldots, x_{i-1}))) > \epsilon \mid \overline{T^{(n)}} = \tilde{t}) = O(e^{-cn}). \quad (25)$$

## 6. Consequences for Prediction

We summarize the implications of our results for prediction of individual sequences based on Maximum Entropy distributions. In this section $\mathcal{X}$ is finite and $Q$ stands for the uniform distribution. Suppose then you have to make predictions about a sequence $(x_1, \ldots, x_n)$. You know the sequence satisfies the given constraint (i.e. for some $n$, $x^{(n)} \in \mathcal{C}_n$, with $\mathcal{C}_n$ as in (20)), but you do not know the length $n$ of the sequence in advance. We distinguish between the special case of the log loss function $\textsc{loss}(x; p) = -\log p(x)$ and the general case of arbitrary (computable) loss functions.

(1) **(log loss)** The MaxEnt distribution $\tilde{p}$ is worst-case optimal with respect to log loss, where the worst-case is over all sequences of all lengths satisfying the constraint. This is a consequence of Theorem 4.

(2) **(log loss)** Whatever sample $x^{(n)} \in \mathcal{C}_n$ arrives, the average log loss you make per outcome when you predict outcomes using $\tilde{p}$ is determined in advance and will be exactly equal to $\mathbf{H}_q(\tilde{p}) = E_{\tilde{p}}[-\log \tilde{p}(X)]$. This is also a consequence of Theorem 4.

(3) **(log loss)** For the overwhelming majority of sequences satisfying the constraint, $\tilde{p}$ will be asymptotically almost optimal with respect to log loss in the following sense: the excess loss of $\tilde{p}$ over *every* other prediction strategy (including strategies depending on past data) is at most a sub-linear function of $n$. This is a consequence of Theorem 5. In Example 1, an example of an exceptional sequence for which $\tilde{p}$ is not optimal would be any sequence consisting of 50% fours and 50% fives.

(4) **(general loss)** For *every* regular loss function LOSS (satisfying Condition 3), predicting using $\tilde{\delta}$, (that is, *acting as if the sample had been generated by* $\tilde{p}$) leads to almost optimal predictions for the overwhelming majority of sequences satisfying the constraint, in the following sense: the excess loss of $\tilde{\delta}$ over *every* other prediction strategy is at most a sub-linear function of $n$. This is a consequence of Theorem 5.

We stress that the fact that items (3) and (4) hold for the overwhelming majority of sequences certainly does *not* imply that they will hold on actual, real-world sequences! Often these will exhibit more regularity than the observed constraint, and then $\tilde{\delta}$ is not necessarily optimal any more.

There are two important points we have neglected so far: (1) in practice, the given constraints will often only hold approximately. (2) the results of this paper have important implications for Maximum Likelihood and Bayesian prediction of sequences based on model classes that are exponential families [1]. The reason is that the Maximum Entropy model for constraint $E[T] = \tilde{t}$ is the Maximum *Likelihood* model for the (exponential family) model class given by (9) for *every* sequence of data $x^{(n)}$ with $n^{-1} \sum_{i=1}^{n} T(x_i) = \tilde{t}$ (see e.g. [5] or [12]) . The connection will be further discussed in the journal version of this paper.

## References

[1] K. Azoury and M. Warmuth. Relative loss bounds for on-line density estimation with the exponential family of distributions. In *Proceedings of the Fifteenth Conference on Uncertainty in Artificial Intelligence (UAI '99)*, pages 31–40. Morgan Kaufmann, 1999.

[2] A. Barron, J. Rissanen, and B. Yu. The minimum description length principle in coding and modeling. *IEEE Transactions on Information Theory*, 44(6):2743–2760, 1998.

[3] R.N. Bhattacharya and R. Ranga Rao. *Normal Approximation and Asymptotic Expansions*. John Wiley, 1976.

[4] P. Billingsley. *Convergence of Probability Measures*. Wiley, 1968.

[5] T.M. Cover and J.A. Thomas. *Elements of Information Theory*. Wiley Interscience, New York, 1991.

[6] I. Csiszár. *I*-divergence geometry of probability distributions and minimization problems. *The Annals of Probability*, 3(1):146–158, 1975.

[7] I. Csiszár. Sanov property, generalized *i*-projection and a conditional limit theorem. *The Annals of Probability*, 12(3):768–793, 1984.

[8] I. Csiszár. Why least squares and maximum entropy? An axiomatic approach to inference for linear inverse problems. *The Annals of Statistics*, 19(4):2032–2066, 1991.

[9] M. Feder. Maximum entropy as a special case of the minimum description length criterion. *IEEE Transactions on Information Theory*, 32(6):847–849, 1986.

[10] W. Feller. *An Introduction to Probability Theory and Its Applications*, volume 2. Wiley, 1968. Third edition.

[11] P. D. Grünwald. Maximum entropy and the glasses you are looking through. In *Proceedings of the Sixteenth Conference on Uncertainty in Artificial Intelligence (UAI 2000)*. Morgan Kaufmann Publishers, 2000.

[12] P.D. Grünwald. *The Minimum Description Length Principle and Reasoning under Uncertainty*. PhD thesis, University of Amsterdam, The Netherlands, October 1998. Available as ILLC Dissertation Series 1998-03; see www.cwi.nl/~pdg.

[13] W. Hoeffding. Probability inequalities for sums of bounded random variables. *Journal of the American Statistical Association*, 58:13–30, 1963.

[14] E.T. Jaynes. Where do we stand on maximum entropy? In R.D. Levine and M. Tribus, editors, *The Maximum Entropy Formalism*, pages 15–118. MIT Press, Cambridge, MA, 1978.

[15] E.T. Jaynes. On the rationale of maximum-entropy methods. *Proceedings of the IEEE*, 70(939-951), 1982.

[16] E.T. Jaynes. *Papers on Probability, Statistics and Statistical Physics*. Kluwer Academic Publishers, second edition, 1989.

[17] E.T. Jaynes. Probability theory: the logic of science. Available at ftp://bayes.wustl.edu/Jaynes.book/, 1996.

[18] J. N. Kapur and H. K Kesavan. *Entropy Optimization Principles with Applications*. Academic Press, Inc., 1992.

[19] R.E. Kass and P.W. Voss. *Geometrical Foundations of Asymptotic Inference*. Wiley Interscience, 1997.

[20] J. Lafferty. Additive models, boosting and inference for generalized divergences. In *Proceedings of the Twelfth Annual Workshop on Computational Learning Theory (COLT '99)*, 1999.

[21] M. Li and P.M.B. Vitányi. *An Introduction to Kolmogorov Complexity and Its Applications*. Springer-Verlag, New York, revised and expanded second edition, 1997.

[22] N. Merhav and M. Feder. A strong version of the redundancy-capacity theorem of universal coding. *IEEE Transactions on Information Theory*, 41(3):714–722, 1995.

[23] J. Rissanen. Stochastic complexity and modeling. *The Annals of Statistics*, 14:1080–1100, 1986.

[24] J. Rissanen. *Stochastic Complexity in Statistical Inquiry*. World Scientific Publishing Company, 1989.

[25] J. Rissanen. Strong optimality of the normalized ML models as universal codes, 2001. To appear in *IEEE Transactions on Information Theory*.

[26] F. Topsøe. Information theoretical optimization techniques. *Kybernetika*, 15(1), 1979.

[27] J. van Campenhout and T. Cover. Maximum entropy and conditional probability. *IEEE Transactions on Information Theory*, IT-27(4):483–489, 1981.

[28] D. Williams. *Probability with Martingales*. Cambridge Mathematical Textbooks, 1991.

## APPENDIX A. PROOFS OF THE CONCENTRATION PHENOMENA FOR TYPICAL SETS

**Theorem 1. (the concentration phenomenon for typical sets, lattice case)**
*Assume we are given a constraint of form (2) such that $T$ is of the lattice form and $h = (h_1, \ldots, h_k)$ is the span of $T$ and such that conditions 1 and 2 hold. Then there exists a sequence $\{c_i\}$ satisfying*

$$\lim_{n \to \infty} c_n = \frac{\prod_{j=1}^{k} h_j}{\sqrt{(2\pi)^k \det \Sigma}}$$

*such that*

(1) *Let $\mathcal{A}_1, \mathcal{A}_2, \ldots$ be an arbitrary sequence of sets with $\mathcal{A}_i \subset \mathcal{X}^i$. For all $n$ with $Q(T_n = \tilde{t}) > 0$, we have:*

$$\tilde{P}(\mathcal{A}_n) \geq n^{-k/2} c_n Q(\mathcal{A}_n \mid \overline{T^{(n)}} = \tilde{t}). \tag{26}$$

*Hence if $\mathcal{B}_1, \mathcal{B}_2, \ldots$ is a sequence of sets with $\mathcal{B}_i \subset \mathcal{X}^i$ whose probability tends to 1 under $\tilde{P}$ in the sense that $1 - \tilde{P}(\mathcal{B}_n) = O(f(n)n^{-k/2})$ for some function $f : \mathbf{N} \to \mathbf{R}$; $f(n) = o(1)$, then $Q(\mathcal{B}_n | \overline{T^{(n)}} = \tilde{t})$ tends to 1 in the sense that $1 - Q(\mathcal{B}_n | \overline{T^{(n)}} = \tilde{t}) = O(f(n))$.*

(2) *If for all $n$, $\mathcal{A}_n \subseteq \{x^{(n)} \mid n^{-1} \sum_{i=1}^{n} T(x_i) = \tilde{t}\}$ then (26) holds with equality.*

*Proof.* We need the following theorem[6]:

---

[6]Feller gives the theorem only for 1-dimensional lattice random variables with $E[T] = 0$ and $\text{var}[T] = 1$; extending the proof to $k$-dimensional random vectors with arbitrary means and covariances is, however, completely straightforward: see XV.7 (page 494) of [10].

**Theorem. ('local central limit theorem for lattice random variables', [10], page 490)** Let $T = (T_{[1]}, \ldots, T_{[k]})$ be a lattice random vector and $h_1, \ldots, h_k$ be the coresponding spans as in Definition 1; let $E_P[T(X)] = t$ and suppose that $P$ satisfies Condition 2 with $T$-covariance matrix $\Sigma$. Let $X_1, X_2, \ldots$ be i.i.d. with common distribution $P$. Let $V$ be a closed and bounded set in $\mathbf{R}^k$. Let $v_1, v_2, \ldots$ be a sequence in $V$ such that for all $n$, $P(\sum_{i=1}^n (T_i - t)/\sqrt{n} = v_n) > 0$. Then as $n \to \infty$,

$$\frac{n^{k/2}}{\prod_{j=1}^k h_j} P\left(\frac{\sum_{i=1}^n (T_i - t)}{\sqrt{n}} = v_n\right) - \aleph(v_n) \to 0.$$

Here $\aleph$ is the density of a $k$-dimensional normal distribution with mean vector $\mu = t$ and covariance matrix $\Sigma$.

The theorem shows that there exists a sequence $d_1, d_2, \ldots$ with $\lim_{n \to \infty} d_n = 1$ such that, for all $n$ with $P(\sum_{i=1}^n (T_i - t) = \mathbf{0}) > 0$,

$$\frac{\frac{n^{k/2}}{\prod_{j=1}^k h_j} P\left(\frac{\sum_{i=1}^n (T_i - t)}{\sqrt{n}} = \mathbf{0}\right)}{\aleph(0)} = \frac{\sqrt{(2\pi n)^k \det \Sigma}}{\prod_{j=1}^k h_j} P\left(\frac{1}{n} \sum_{i=1}^n T_i = t\right) = d_n \qquad (27)$$

The proof now becomes very simple. First note that $\tilde{P}(\mathcal{A}_n \mid \overline{T^{(n)}} = \tilde{t}) = Q(\mathcal{A}_n \mid \overline{T^{(n)}} = \tilde{t})$ (write out the definition of conditional probability and realize that $\exp(-\tilde{\beta}^{\mathrm{T}} T(x)) = \exp(-\tilde{\beta}^{\mathrm{T}} \tilde{t}) = \text{constant}$ for all $x$ with $T(x) = \tilde{t}$. Use this to show that

$$\begin{aligned} \tilde{P}(\mathcal{A}_n) &\geq \tilde{P}(\mathcal{A}_n, \overline{T^{(n)}} = \tilde{t}) = \tilde{P}(\mathcal{A}_n \mid \overline{T^{(n)}} = \tilde{t}) \tilde{P}(\overline{T^{(n)}} = \tilde{t}) &(28) \\ &= Q(\mathcal{A}_n \mid \overline{T^{(n)}} = \tilde{t}) \tilde{P}(\overline{T^{(n)}} = \tilde{t}). \end{aligned}$$

Clearly, with $\tilde{P}$ in the rôle of $P$, the local central limit theorem is applicable to random vector $T$. Then, by (27), $\tilde{P}(\overline{T^{(n)}} = \tilde{t}) = (\prod_{j=1}^k h_j)/\sqrt{(2\pi n)^k \det \Sigma} d_n$. Defining $c_n := \tilde{P}(\overline{T^{(n)}} = \tilde{t}) n^{k/2}$ finishes the proof of item 1. For item 2, notice that in this case (28) holds with equality; the rest of the proof remains unchanged. $\qquad\square$

**Theorem 2. (the concentration phenomenon for typical sets, continuous case)** *Assume we are given a constraint of form (2) such that $T$ is of regular continuous form and such that Conditions 1 and 2 hold. Fix some $h > 0$ and let $\epsilon_n := h/n$. Then there exists a sequence $c_1, c_2, \ldots$ satisfying*

$$\lim_{n \to \infty} c_n = e^{-2|\tilde{\beta}|} \frac{h^k}{\sqrt{(2\pi)^k \det \Sigma}}$$

*such that*

(1) *Let $\mathcal{A}_1, \mathcal{A}_2, \ldots$ be an arbitrary sequence of (measurable) sets with $\mathcal{A}_i \subset \mathcal{X}^i$. For all $n$ we have:*

$$\tilde{P}(\mathcal{A}_n) \geq n^{-k/2} c_n Q(\mathcal{A}_n \mid \overline{T^{(n)}} \in \mathcal{B}_{\epsilon_n}(\tilde{t})). \qquad (29)$$

(2) *If for all $n$, $\mathcal{A}_n \subseteq \{x^{(n)} \mid n^{-1} \sum_{i=1}^n T(x_i) \in \mathcal{B}_{\epsilon_n}(\tilde{t})\}$ then $\lim_{n \to \infty} \tilde{P}(\mathcal{A}_n) \leq e^{2|\tilde{\beta}|} n^{-k/2} c_n Q(\mathcal{A}_n \mid \overline{T^{(n)}} \in \mathcal{B}_{\epsilon_n}(\tilde{t}))$.*

*Proof.* The proof is completely analogous to the discrete case, except that now we use the 'local central limit theorem for continuous random variables'. The 1-dimensional case,

along with a simple proof, can be found in [10] (Theorem 2, page 489). The general case can be found in [3] (Theorem 19.1). We cite it explicitly:

**Theorem. (uniform local central limit theorem for random variables in $\mathbf{R}^k$, [3])**
Let $T = (T_{[1]}, \ldots, T_{[k]})$ be a random vector; let $P$ be a distribution so that $T(X)$ has a bounded continuous density with respect to Lebesgue measure. let $E_P[T(X)] = t$ and suppose that $P$ satisfies Condition 2 with $T$-covariance matrix $\Sigma$. Let $X_1, X_2, \ldots$ be i.i.d. with common distribution $P$. Let $P_n^*$ be the distribution of $\frac{\sum_{i=1}^{n}(T_i - t)}{\sqrt{n}}$. Then $P_n^*$ has a density $p_n^*$ and we have
$$\lim_{n \to \infty} \sup_{t \in \mathbf{R}^k} |p_n^*(t) - \aleph(t)| = 0.$$
Here $\aleph$ is the density of a $k$-dimensional normal distribution with mean vector $\mu = t$ and covariance matrix $\Sigma$.  □

## Appendix B. Proof of Theorem 3

Before giving the proof, we first establish some facts about the conditional distributions $Q(\cdot \mid \overline{T^{(n)}} = t)$ and $\tilde{P}(\cdot \mid \overline{T^{(n)}} = t)$. Recall that in the measure-theoretic framework, these can be arbitrarily chosen for $t$ in a null set [28]. If we want to speak about these quantities for *arbitrary* $t$, we need to make sure that continuous (in $t$) versions of the conditional distributions exist; we then define the 'canonical' conditional distribution to be the continuous version.

**Proposition 2.** *Suppose $T$ appearing in (2) is of regular continuous form and such that Condition 1 and 2 both hold. Then there exists an open ball $\mathcal{B}_\epsilon(\tilde{t})$ around $\tilde{t}$ such that, for $t \in \mathcal{B}_\epsilon(\tilde{t})$, there exists a continuous (in the weak topology) version of both $Q(\cdot \mid \overline{T^{(n)}} = t)$ and $\tilde{P}(\cdot \mid \overline{T^{(n)}} = t)$.*

*Proof.* Condition 2 ensures that $\tilde{t}$ lies in the interior of the range of $T$, so we can take $\epsilon_0$ so that $\mathcal{B}_{\epsilon_0}(\tilde{t})$ falls within this range. $T$ has a bounded and continuous density (with respect to Lebesgue measure) under $Q$. It follows that there exists a version of $Q(\cdot|\overline{T^{(n)}} = t)$ such that for all $t_0 \in \mathcal{B}_{\epsilon_0}(\tilde{t})$, for every bounded continuous function $g : \mathcal{X}^n \to \mathbf{R}$, $E_Q[g(X)|\overline{T^{(n)}} = t_0]$ is given by a fraction of two Riemann-integrals which are uniquely defined for each $t_0$. It is then straightforward to show that, for this version of $E_Q[\cdot|\overline{T^{(n)}} = t]$, $\lim_{t \to t_0} E_Q[g(X)|\overline{T^{(n)}} = t] = E_Q[g(X)|\overline{T^{(n)}} = t_0]$ for all $t_0 \in \mathcal{B}_{\epsilon_0}(\tilde{t})$. Since this holds for all continuous bounded $g$, our version of $Q(\cdot|\overline{T^{(n)}} = t_0)$ is continuous in the weak topology; see Theorem 2.1 of [4]. Existence of a continuous version of $\tilde{P}(\cdot|\overline{T^{(n)}} = t)$ is shown in the same way.  □

We now restate and prove Theorem 3:

**Theorem 3.** Let $\{m_i\}$ be an increasing sequence with $m_i \in \mathbf{N}$, such that $\lim_{n \to \infty} m_n/n = 0$. Assume we are given a constraint of form (2) such that $T$ is of the regular continuous form or of the lattice form and suppose that Conditions 1 and 2 are satisfied. Then as $n \to \infty$, $Q^{m_n}(\cdot \mid \overline{T^{(n)}} = \tilde{t})$ converges weakly to $\tilde{P}^{m_n}(\cdot)$.

*Proof.* We give the proof only for the regular continuous case where $\mathcal{X} = \mathbf{R}$ and $T = (T_{[1]})$ is 1-dimensional. All other cases have analogous proofs. For ease of notation, we omit the subscript $n$ from $m_n$ whenever $n$ is clear from the context.

Note that there exists some function $h : \mathbf{N} \to \mathbf{R}$ with $h(n) = o(1)$ such that we can write $m = nh(n)$. By the definition of weak convergence [4], it is sufficient to prove convergence of the distribution functions of $Q^m(\cdot | \overline{T^{(n)}} = \tilde{t})$ to $\tilde{P}^m$. I.e. we must show that for all sequences $\{R_i\}$ with $R_i \in \mathbf{R}$,

$$\lim_{n \to \infty} Q(X_1 \le R_1, \dots, X_m \le R_m \mid \overline{T^{(n)}} = \tilde{t}) = \tilde{P}(X_1 \le R_1, \dots, X_m \le R_m) \qquad (30)$$

Let us abbreviate $\mathcal{R}_m := \{(x_1, x_2, \dots) \in \mathcal{X}^\infty \mid x_1 \le R_1, \dots, x_m \le R_m\}$. The following equalities both follow, with some work, by Proposition 2 and the definition of weak convergence; we omit the details.

$$Q(\mathcal{R}_m \mid \overline{T^{(n)}} = \tilde{t}) = \lim_{\epsilon \to 0} Q(\mathcal{R}_m \mid | \sum_{i=1}^{n} (T_i - \tilde{t})| < \epsilon) = \lim_{\epsilon \to 0} \tilde{P}(\mathcal{R}_m \mid | \sum_{i=1}^{n} (T_i - \tilde{t})| < \epsilon) \quad (31)$$

The last equality is the analogue of (28) in Theorem 1. By definition, $\tilde{P}(\mathcal{R}_m \mid | \sum_{i=1}^{n} (T_i - \tilde{t})| < \epsilon) = \tilde{P}_n^{(\text{num})} / \tilde{P}_n^{(\text{den})}$, where we abbreviated

$$\tilde{P}_n^{(\text{num})} \quad := \quad \tilde{P}(\mathcal{R}_m \, ; | \sum_{i=1}^{n} (T_i - \tilde{t})| < \epsilon) \qquad (32)$$

$$\tilde{P}_n^{(\text{den})} \quad := \quad \tilde{P}(| \sum_{i=1}^{n} (T_i - \tilde{t})| < \epsilon) \qquad (33)$$

The strategy of the proof will be to rewrite $\tilde{P}_n^{(\text{num})}$ and $\tilde{P}_n^{(\text{den})}$ so that the local central limit theorem can be applied to them both. In our previous theorems for 'typical sets', we only had to apply the local central limit theorem to $\tilde{P}_n^{(\text{den})}$; the fact that $m/n \to 0$ as $n \to \infty$ allows us to apply it to $\tilde{P}_n^{(\text{num})}$ too. In Stage 2 we combine the results, take the limits $\lim_{n \to \infty} \lim_{\epsilon \to 0} \frac{\tilde{P}_n^{(\text{num})}}{\tilde{P}_n^{(\text{den})}}$ and by (31) the result will follow.

**Stage 1.** Let $\delta > 0$. We partition the sample space $\mathcal{X}^m = \mathbf{R}^m$ into hyper-rectangles $\mathcal{H}$, which, when mapped to '$T$-space' by the transformation $T(\mathcal{H}) := \{t^{(m)} \in \mathbf{R}^m \mid \exists x^{(m)} \in \mathcal{H} : T(x_1) = t_1, \dots, T(x_m) = t_m\}$, become hyper-cubes with side length $\delta$. We define $a_0 = 0$ and

$$a_{j+1} \quad = \quad \inf \{a \mid a > a_j \, ; \, |T(a) - T(a_j)| = \delta\} \text{ if } j \ge 0 \qquad (34)$$

$$a_{j-1} \quad = \quad \sup \{a \mid a < a_j \, ; \, |T(a) - T(a_j)| = \delta\} \text{ if } j \le 0. \qquad (35)$$

If the inf in (34) does not exist, $a_{j+1} := a_j + 1$. If the sup in (35) does not exist, $a_{j-1} := a_j - 1$. Both $T$ and $X$ have a bounded continuous density with respect to Lebesgue-measure, which implies that for all $j \in \mathbf{Z}$, $a_j > a_{j-1}$ and that $\lim_{j \to (-)\infty} a_j = (-)\infty$. Therefore we can cover $\mathcal{X}$ by the sets $\mathcal{H}(j) := (a_j, a_{j+1}]$. For $j^{(m)} = (j_1, \dots, j_m) \in \mathbf{Z}^m$, define $\mathcal{H}(j^{(m)}) := (a_{j_1}, a_{j_1+1}] \times (a_{j_2}, a_{j_2+1}] \times \dots \times (a_{j_m}, a_{j_m+1}]$. Clearly, the sets $\mathcal{H}(j^{(m)})$ cover $\mathcal{X}^m$.

Let $S_m := \sum_{i=1}^{m}(T_i - \tilde{t})$ and define

$$\mathcal{S}_m := \{(x_1, x_2, \ldots) \in \mathcal{X}^\infty \mid |S_m| < (h(n))^{1/3}\sqrt{n-m}\}. \tag{36}$$

We have

$$\tilde{P}_n^{(\mathrm{num})} \geq \tilde{P}(\mathcal{R}_m \cap \mathcal{S}_m \,;\, |\sum_{i=1}^{n}(T_i - \tilde{t})| < \epsilon) =$$

$$\sum_{j^{(m)} \in \mathbf{Z}^m} \tilde{P}(\mathcal{H}(j^{(m)}) \cap \mathcal{R}_m \cap \mathcal{S}_m \,;\, |\sum_{i=1}^{n}(T_i - \tilde{t})| < \epsilon) \geq$$

$$\sum_{j^{(m)} \in \mathbf{Z}^m} \tilde{P}(\mathcal{H}(j^{(m)}) \cap \mathcal{R}_m \cap \mathcal{S}_m \,;\, |\sum_{i=m+1}^{n}(T_i - \tilde{t}) + \sum_{i=1}^{m}(t_i - \tilde{t})| < \epsilon + m\delta) \quad (37)$$

where we have defined $t_i := (1/2)(T(a_{j_i}) + T(a_{j_i+1}))$. The last line holds for all $\delta > 0$ (note that we let $\mathcal{H}(\cdot)$ depend on $\delta$). In particular for $\delta = \epsilon/m^2$ we get, letting $s_m := \sum_{i=1}^{m}(t_i - \tilde{t})$,

$$\tilde{P}_n^{(\mathrm{num})} \geq \sum_{j^{(m)} \in \mathbf{Z}^m} \tilde{P}(\mathcal{H}(j^{(m)}) \cap \mathcal{R}_m \cap \mathcal{S}_m \,;\, |\sum_{i=m+1}^{n}(T_i - \tilde{t}) + s_m| < \epsilon(1 + m^{-1})) =$$

$$\sum_{j^{(m)} \in \mathbf{Z}^m} \tilde{P}(\mathcal{H}(j^{(m)}) \cap \mathcal{R}_m \cap \mathcal{S}_m)\tilde{P}(|\sum_{i=1}^{n-m}(T_i - \tilde{t}) + s_m| < \epsilon(1 + m^{-1})) =$$

$$\sum_{j^{(m)} \in \mathbf{Z}^m \,;\, \mathcal{H}(j^{(m)}) \cap \mathcal{S}_m \neq \emptyset} \tilde{P}(\mathcal{H}(j^{(m)}) \cap \mathcal{R}_m \cap \mathcal{S}_m)\tilde{P}(|\sum_{i=1}^{n-m}(T_i - \tilde{t}) + s_m| < \epsilon(1 + m^{-1})). \quad (38)$$

We are now ready to apply the local central limit theorem (see Appendix A) to the right hand side factors in the terms of the sum (38). After some rewriting we see that for all these factors we can write:

$$\tilde{P}(|\sum_{i=1}^{n-m}(T_i - \tilde{t}) + s_m| < \epsilon(1 + m^{-1})) = \frac{1}{\sqrt{2\pi\sigma^2}}\int_{x_l}^{x_r} c_n(x)e^{-x^2/2\sigma^2}\,dx \tag{39}$$

where $\sigma^2 = \Sigma_{11}$ and $x_l = \frac{-s_m - \epsilon(1 + m^{-1})}{\sqrt{n-m}}$, $x_r = \frac{-s_m + \epsilon(1 + m^{-1})}{\sqrt{n-m}}$, $c_n$ some function of $x$ and $s_m$ is the only ingredient that can be different for different terms in (38). Now fix some $\epsilon_0$. Since the sum in (38) is only over terms with $\mathcal{H}(j^{(m)} \cap \mathcal{S}_m) \neq \emptyset$, and by the definition (36) of $\mathcal{S}_m$, for all those terms both $x_l$ and $x_r$ are uniformly (over all $n$, $s_m$, $0 < \epsilon < \epsilon_0$) bounded. Therefore, $c_n(x)$ tends uniformly to 1 for all $x \in [x_l, x_r]$ and for all terms in (38), as long as $\epsilon < \epsilon_0$. It follows that, uniformly for all terms in (38) and all $0 < \epsilon < \epsilon_0$:

$$\tilde{P}(|\sum_{i=1}^{n-m}(T_i - \tilde{t}) + s_m| < \epsilon(1 + m^{-1})) \geq d_n \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \frac{2\epsilon(1 + m^{-1})}{\sqrt{n-m}} \exp\left(-\frac{(|s_m| + \epsilon(1 + m^{-1}))^2}{2\sigma^2(n-m)}\right) \tag{40}$$

where $d_n$ tends to 1 as $n$ tends to $\infty$.

Using the local central limit theorem once more, in the same way as in (39), (40), we can derive the following upper bound for $\tilde{P}_n^{(\text{den})}$:

$$\tilde{P}_n^{(\text{den})} \leq \frac{e_n}{\sqrt{2\pi\sigma^2}} \cdot \frac{2\epsilon}{\sqrt{n}} \tag{41}$$

where $e_n$ tends to 1 as $n$ tends to $\infty$, uniformly for all $\epsilon < \epsilon_0$. Combining (40), (41), (38), we find that

$$\tilde{P}(\mathcal{R}_m \mid \mid \sum_{i=1}^n (T_i - \tilde{t})| < \epsilon) \geq$$

$$\sum_{\substack{j^{(m)} \in \mathbf{Z}^m \\ \mathcal{H}(j^{(m)}) \cap \mathcal{S}_m \neq \emptyset}} \tilde{P}(\mathcal{H}(j^{(m)}) \cap \mathcal{R}_m \cap \mathcal{S}_m) f_n \frac{\frac{2\epsilon(1+m^{-1})}{\sqrt{n-m}} \exp\left(-\frac{(|s_m|+\epsilon(1+m^{-1}))^2}{2\sigma^2(n-m)}\right)}{(2\epsilon/\sqrt{n})} \geq$$

$$\tilde{P}(\mathcal{R}_m \cap \mathcal{S}_m) g_n(\epsilon) \quad (42)$$

where $f_n$ tends to 1 and ($\epsilon_0$ can be chosen such that) $\lim_{n\to\infty} g_n(\epsilon) = 1$ uniformly for all $\epsilon < \epsilon_0$. This follows because, by our definition of $\mathcal{S}_m$, $s_m$ and $\delta$, $s_m/\sqrt{n-m} = o(1)$ for all $\epsilon < \epsilon_0$, uniformly for all terms in the second line of (42). Now (letting $\overline{\mathcal{S}_m} = \mathcal{X}^\infty \setminus \mathcal{S}_m$)

$$\tilde{P}(\mathcal{R}_m \cap \mathcal{S}_m) \geq \tilde{P}(\mathcal{R}_m) - \tilde{P}(\overline{\mathcal{S}_m}) = \tilde{P}(\mathcal{R}_m) - \tilde{P}(|\sum_{i=1}^m (T_i - t)| \geq h(n)^{1/3}\sqrt{n-m}) \geq$$

$$\tilde{P}(\mathcal{R}_m) - \frac{\sigma^2 m}{h(n)^{2/3}(n-m)} = \tilde{P}(\mathcal{R}_m) - \frac{\sigma^2 h(n)^{1/3}}{(1-h(n))} = \tilde{P}(\mathcal{R}_m) - o(1). \quad (43)$$

where we have used Chebyshev's inequality. Combining (42) and (43), we find that

$$\liminf_{n\to\infty} \lim_{\epsilon\to 0} \{ \tilde{P}(\mathcal{R}_m \mid \mid \sum_{i=1}^n (T_i - \tilde{t})| < \epsilon) - \tilde{P}(\mathcal{R}_m) \} \geq 0. \tag{44}$$

**Stage 2.** Stage 2 is now very simple: we repeat exactly the same argument as above with the sets $\overline{\mathcal{R}_m} = \mathcal{X}^\infty \setminus \mathcal{R}_m$. These are all continuity sets which implies that (31) still holds with $\mathcal{R}_m$ replaced by $\overline{\mathcal{R}_m}$. $\overline{\mathcal{R}_m}$. All other steps go through without modification. This repetition of the argument gives

$$\limsup_{n\to\infty} \lim_{\epsilon\to 0} \{ \tilde{P}(\mathcal{R}_m \mid \mid \sum_{i=1}^n (T_i - \tilde{t})| < \epsilon) - \tilde{P}(\mathcal{R}_m) \} \leq 0. \tag{45}$$

Together, (44), (45) and (31) prove the theorem. $\qquad\square$

## Appendix C. Proof of Theorem 5

*Proof.* One easily establishes that the theorem holds trivially if there does not exist an $\alpha > 0$ and a $\delta' \in \Delta$ such that $E_{\tilde{p}}[\text{LOSS}(X; \delta')] - E_{\tilde{p}}[\text{LOSS}(X; \tilde{\delta})] > \alpha$. Suppose then that such a $\delta'$ and $\alpha$ exist. Let $\Delta(\alpha) := \{\delta \in \Delta \mid E_{\tilde{p}}[\text{LOSS}(X; \delta')] - E_{\tilde{p}}[\text{LOSS}(X; \tilde{\delta})] \leq \alpha\}$. By convexity and continuity of $\Delta$, there exist $\alpha_0 > 0$ such that $\Delta(\alpha)$ is non-empty for all $0 < \alpha < \alpha_0$. By compactness of $\Delta$, we can choose an $\alpha$ so small such that for all $\delta \in \Delta(\alpha)$, there exist a $\delta_0 \notin \Delta(\alpha)$ such that for all $x \in \mathcal{X}$, $|\text{LOSS}(x; \delta_0) - \text{LOSS}(x; \delta)| < \epsilon/2$. Fix $\alpha_0$ small enough so that this holds. We now change the prediction strategy $\delta^*$ mentioned in the theorem to a new strategy $\delta^{**}$ as follows: for all $n$, $x^{(n)}$, if $\delta^*(x^{(n)}) \notin \Delta(\alpha_0)$, then

$\delta^{**}(x^{(n)}) = \delta^*(x^{(n)})$. If $\delta^*(x^{(n)}) \in \Delta(\alpha_0)$, then $\delta^{**}(x^{(n)})$ is chosen so that $\delta^{**}(x^{(n)}) \notin \Delta(\alpha_0)$ but for all $x$, $|\text{LOSS}(x; \delta^{**}(x^{(n)})) - \text{LOSS}(x; \delta^*(x^{(n)}))| < \epsilon/2$. We have that

$$Q(\frac{1}{n}(\sum_{i=1}^{n} \text{LOSS}(x_i; \tilde{\delta}) - \sum_{i=1}^{n} \text{LOSS}(x_i; \delta^*(x_1, \ldots, x_{i-1}))) > \epsilon \mid \overline{T^{(n)}} = \tilde{t}) \leq$$

$$Q(\frac{1}{n}(\sum_{i=1}^{n} \text{LOSS}(x_i; \tilde{\delta}) - \sum_{i=1}^{n} \text{LOSS}(x_i; \delta^{**}(x_1, \ldots, x_{i-1}))) > \frac{\epsilon}{2} \mid \overline{T^{(n)}} = \tilde{t}) \quad (46)$$

while at the same time, for all $i$, $x^{(i)}$,

$$E_{\tilde{p}}[\text{LOSS}(X; \delta^{**}(x^{(i)}))] - E_{\tilde{p}}[\text{LOSS}(X; \tilde{\delta})] > \alpha_0 > 0. \quad (47)$$

By the Hoeffding bound for random variables that are bounded from below [13, Theorem 3], (47) implies

$$\tilde{P}(\frac{1}{n}(\sum_{i=1}^{n} \text{LOSS}(x_i; \tilde{\delta}) - \sum_{i=1}^{n} \text{LOSS}(x_i; \delta^{**}(x_1, \ldots, x_{i-1}))) > \frac{\epsilon}{2}) = O(e^{-cn})$$

for some $c > 0$ depending on $\alpha$. Together with Theorem 1 and (46) this implies the theorem. $\qquad \square$