

Recap

Letter Typical Sequences:

- We are looking for a subset $\mathcal{T}^{(n)}(\mathbb{P}) \subseteq \mathcal{X}^n$ such that

$$(i) \text{ "small": } \lim_{n \rightarrow \infty} \frac{|\mathcal{T}^{(n)}(\mathbb{P})|}{|\mathcal{X}^n|} = 0$$

$$(ii) \text{ "high probability": } \lim_{n \rightarrow \infty} \mathbb{P}^{\otimes n}(x^n \in \mathcal{T}^{(n)}(\mathbb{P})) = 1$$

Key Observation:

While for a given $\mathbb{P} \in \mathcal{P}(\mathcal{X})$ it may be the case that $\mathbb{P}^{\otimes n}(\{x^n\}) > 0$, if $x^n \in \mathcal{X}^n$, we expect that for, e.g., $\text{Ber}(\alpha)$ we will most likely draw a sequence w/ roughly $\alpha \cdot n$ ones and $(1-\alpha) \cdot n$ zeros.

Empirical Frequency / PMF

For any $x^n \in \mathcal{X}^n$

$$N_{x^n}(a) := \sum_{i=1}^n \mathbf{1}_{\{x_i = a\}}, \quad a \in \mathcal{X}$$

$$\gamma_{x^n}(a) := \frac{1}{n} N_{x^n}(a)$$

The above $\gamma_{x^n}(a)$ is the empirical frequency of $x^n \in \mathcal{X}^n$.

To quantify "roughly" we will introduce a sleekness parameter $\varepsilon > 0$.

(Note we are only discussing \mathcal{X} for which $|\mathcal{X}| < \infty$)

Definition (Letter Typical Set):

Let X be a finite alphabet, $P \in P(X)$ have pmf p , $n \in \mathbb{N}$ and $\varepsilon > 0$. The ε -letter typical set of n -lengthed sequences with respect to P with slackness ε is

$$T_{\varepsilon}^{(n)}(P) := \left\{ x^n \in X^n : |\nu_{x^n}(a) - p(a)| \leq \varepsilon p(a), \forall a \in X \right\}$$

where $\nu_{x^n}(a)$ is the empirical frequency of a in $x^n = (x_1, \dots, x_n)$.

Example

$$P = \text{Ber}(0.3), n=10, \varepsilon = \frac{1}{2}$$

$$T_{\frac{1}{2}}^{(10)}(\text{Ber}(0.3)) = \left\{ x^{10} \in \{0, 1\}^{10} : \begin{array}{l} |\nu_{x^{10}}(0) - 0.7| \leq 0.35 \\ |\nu_{x^{10}}(1) - 0.3| \leq 0.15 \end{array} \right\}$$

Remark (Initial Observation)

- ① If $p(a) = 0$, then $\nu_{x^n}(a) = 0 \quad \forall x^n \in T_{\varepsilon}^{(n)}(P)$, meaning that zero probability letters cannot appear in letter typical sequences.
- ② Suppose $\varepsilon < 1$. If $p(a) > 0$, then $\nu_{x^n}(a) > 0 \quad \forall x^n \in T_{\varepsilon}^{(n)}(P)$ (indeed, if $\nu_{x^n}(a) = 0$ for some x^n , then we must have $p(a) \leq \varepsilon \cdot p(a)$ which is a contradiction)

③ If $P = \text{Unif}(X)$, then for all $x^n \in T_{\varepsilon}^{(n)}(P)$,

$$\frac{(1-\varepsilon)}{|X|} \leq v_{x^n}(a) \leq \frac{(1+\varepsilon)}{|X|}$$

It follows that for small enough ε , it holds

$T_{\varepsilon}^{(n)}(P) \subsetneq X^n$. For example, if $\varepsilon < |X| - 1$, then

$T_{\varepsilon}^{(n)}(P)$ does not contain any sequence comprised of a single letter.

The strength of letter typical sequences lie in the fact that they lend themselves well for empirical averages, as shown in the following lemma.

Lemma (Typical Averaging Lemma):

Let $g: X \rightarrow \mathbb{R}$ be a nonnegative measurable function such that $\mathbb{E}_P[g(X)] < \infty$. Then, for all $x^n \in T_{\varepsilon}^{(n)}(P)$, we have

$$(1-\varepsilon) \mathbb{E}_P[g(X)] \leq \frac{1}{n} \sum_{i=1}^n g(x_i) \leq (1+\varepsilon) \mathbb{E}_P[g(X)]$$

Proof

If $x^n \in T_{\varepsilon}^{(n)}(P)$, $a \in X$ we have $|v_{x^n}(a) - p(a)| \leq \varepsilon p(a)$

where p is the PMF of X . It follows that

$\frac{1}{n} \sum_{i=1}^n g(x_i) = \sum_{a \in X} v_{x^n}(a) g(a)$. Then

$$\begin{aligned} \left| \frac{1}{n} \sum_{i=1}^n g(x_i) - \mathbb{E}_P[g(X)] \right| &= \left| \sum_{a \in X} (v_{x^n}(a) - p(a)) g(a) \right| \stackrel{\text{Triangle Inequality}}{\leq} \sum_{a \in X} |v_{x^n}(a) - p(a)| g(a) \leq \varepsilon \sum_{a \in X} p(a) g(a) \\ &= \varepsilon \mathbb{E}_P[g] \end{aligned}$$

Based on this lemma we derive further properties of $T_{\varepsilon}^{(n)}(\mathbb{P})$. In particular, we provide upper and lower bounds on the cardinality and the probability of $T_{\varepsilon}^{(n)}(\mathbb{P})$.

Theorem (Properties of Letter Typical Sequences):

Suppose $T_{\varepsilon}^{(n)}(\mathbb{P})$ is an ε -letter typical set. Let $\mathbb{P}^{\otimes n}$ be the n -fold product measure induced by \mathbb{P} , i.e., $\mathbb{P}^{\otimes n}(\{x^n\}) = \prod_{i=1}^n p(x_i)$ for all $x^n \in \mathcal{X}^n$. The following hold:

① For all $x^n \in T_{\varepsilon}^{(n)}(\mathbb{P})$,

$$2^{-n(1+\varepsilon)H(\mathbb{P})} \leq \mathbb{P}^{\otimes n}(\{x^n\}) \leq 2^{-n(1-\varepsilon)H(\mathbb{P})}$$

② We have

$$\lim_{n \rightarrow \infty} \mathbb{P}^{\otimes n}(T_{\varepsilon}^{(n)}(\mathbb{P})) = 1$$

③ For n sufficiently large,

$$(1-\varepsilon)2^{n(1-\varepsilon)H(\mathbb{P})} \leq |T_{\varepsilon}^{(n)}(\mathbb{P})| \leq 2^{n(1+\varepsilon)H(\mathbb{P})}$$

Proof

① Consequence of the lemma. Let $g(x) = -\log(p(x))$.

Then $\mathbb{E}_P[g(x)] = H(P)$ and $\frac{1}{n} \sum_i g(x_i) = \frac{1}{n} \log \frac{1}{P^{\otimes n}(\{x^n\})}$.

Thus

$$(1-\varepsilon)H(P) \leq \frac{1}{n} \log \frac{1}{P^{\otimes n}(\{x^n\})} \leq (1+\varepsilon)H(P)$$

Exponentiating all terms gives

$$2^{-n(1-\varepsilon)H(P)} \geq P^{\otimes n}(\{x^n\}) \geq 2^{-n(1+\varepsilon)H(P)}$$

② By definition we have

$$P^{\otimes n}(T_\varepsilon^{(n)}(P)) = P^{\otimes n}\left(\bigcap_{a \in \mathcal{X}} \{x^n : |\nu_{x^n}(a) - p(a)| \leq \varepsilon p(a)\}\right)$$

By the weak law of large numbers (WLLN), we have that for a set of arbitrary functions $\{f_k(x)\}_{k=1}^K$ (each with finite expectation) and any $\delta > 0$,

$$\lim_{n \rightarrow \infty} P^{\otimes n}\left(\bigcap_{k=1}^K \left\{x : \left| \frac{1}{n} \sum_{i=1}^n f_k(x_i) - \mathbb{E}_P[f_k(x)] \right| \leq \delta \right\}\right) = 1$$

Now set $K = |\mathcal{X}|$ and $f_a = \mathbf{1}_{\{x=a\}}$ if $a \in \mathcal{X}$. Then

$\mathbb{E}_P[f_a(x)] = p(a)$ and $\frac{1}{n} \sum_{i=1}^n f_a(x_i) = \nu_{x^n}(a)$. The WLLN then implies

$$\lim_{n \rightarrow \infty} P^{\otimes n}(T_\varepsilon^{(n)}(P)) = \lim_{n \rightarrow \infty} P^{\otimes n}\left(\bigcap_{a \in \mathcal{X}} \{x^n : |\nu_{x^n}(a) - p(a)| \leq \varepsilon p(a)\}\right) = 1$$

③ Using the fact that $T_\epsilon^{(n)}(\Omega) \subseteq \mathcal{X}^n$, we get

$$1 = P^\otimes(\mathcal{X}^n) \geq P^{\otimes n}(T_\epsilon^{(n)}(\Omega)) = \sum_{x^n \in T_\epsilon^{(n)}(\Omega)} P^{\otimes n}(\{x^n\})$$

$$\geq \sum_{x^n \in T_\epsilon^{(n)}(\Omega)} 2^{-n(1+\epsilon)H(\Omega)}$$

$$(\text{follows from ①}) \quad = |T_\epsilon^{(n)}(\Omega)| 2^{-n(1+\epsilon)H(\Omega)}$$

$$\Rightarrow |T_\epsilon^{(n)}(\Omega)| \leq 2^{n(1+\epsilon)H(\Omega)}$$

For n sufficiently large, by ② we have $1-\epsilon \leq P^{\otimes n}(T_\epsilon^{(n)}(\Omega))$.

Then

$$1-\epsilon \leq \sum_{x^n \in T_\epsilon^{(n)}(\Omega)} P^{\otimes n}(\{x^n\}) \leq |T_\epsilon^{(n)}(\Omega)| 2^{-n(1-\epsilon)H(\Omega)}$$

$$\Rightarrow |T_\epsilon^{(n)}(\Omega)| \geq (1-\epsilon) 2^{-n(1-\epsilon)H(\Omega)}$$

□

Remark (Refined Result)

A finite sample bound on the probability of $T_\epsilon^{(n)}(\Omega)$ can be derived using a refined argument (based on Chernoff bounds). This gives

$$1 - \delta_\epsilon(\Omega, n) \leq P^{\otimes n}(T_\epsilon^{(n)}(\Omega)) \leq 1$$

where $\delta_\epsilon(\Omega, n) := 2|\Omega|e^{-2n\epsilon^2\mu^2}$ and $\mu := \min_{\alpha \in \mathcal{X}: p(\alpha) > 0} p(\alpha)$.

Note that $\lim_{n \rightarrow \infty} \delta_\varepsilon(P, n) = 0$, for any fixed $\varepsilon > 0$.

It can thus be shown

$$(1 - \delta_\varepsilon(P, n)) 2^{n(1-\varepsilon)H(P)} \leq |\mathcal{T}_\varepsilon^{(n)}(P)| \leq 2^{n(1+\varepsilon)H(P)}$$

$\forall n \in \mathbb{N}$. This strengthens (1), (2) in Theorem above.

Remark (Interpretation & Illustration):

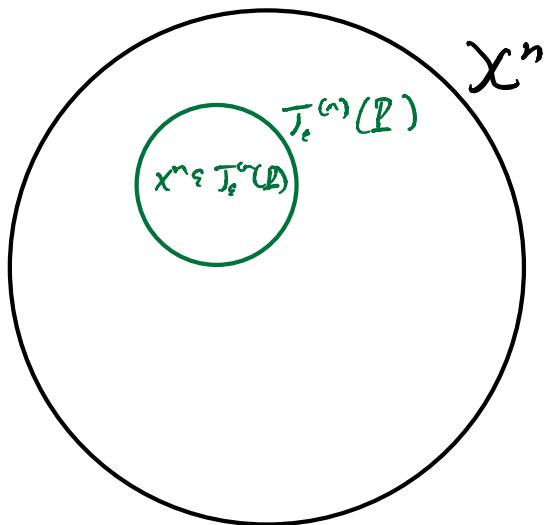
As $n \rightarrow \infty$, the set $\mathcal{T}_\varepsilon^{(n)}(P)$ takes a smaller and smaller portion of the whole space \mathbb{R}^n , while the probability tends to concentrate on this set, uniformly distributed across its elements. The following diagram is a useful illustration

We can further simplify the statement by introducing the following notation. Denote

by $a_n = 2^{nb}$, the inequality $2^{n(b-\delta(\varepsilon))} \leq a_n \leq 2^{n(b+\delta(\varepsilon))}$ if

there exists a $\delta(\varepsilon) = o(1)$ as

$\varepsilon \rightarrow 0$. Then by taking $\delta(\varepsilon) = 2\varepsilon H(P)$, we can write $|\mathcal{T}_\varepsilon^{(n)}(P)| = 2^{nH(P)}$, and $P^{\otimes n}(\{x^n\}) = 2^{-nH(P)}$. The properties in Theorem 1 can now be understood as follows:



① The cardinality of $T_{\epsilon}^{(n)}(\underline{P}) \subsetneq \mathcal{X}^n$ is much smaller than that of \mathcal{X}^n . More precisely, if $\underline{P} \neq \text{Unif}(\mathcal{X})$ and sufficiently small ϵ , it holds that

$$\lim_{n \rightarrow \infty} \frac{|T_{\epsilon}^{(n)}(\underline{P})|}{|\mathcal{X}^n|} = \lim_{n \rightarrow \infty} \frac{2^{nH(\underline{P})}}{2^{n \log |\mathcal{X}|}} = 0$$

② Despite being "small", the probability of $T_{\epsilon}^{(n)}(\underline{P})$ is arbitrarily close to 1 for large n . This means that if we draw an iid sequence of length n with respect to \underline{P} , then with arbitrarily high probability this sequence lands in the set $T_{\epsilon}^{(n)}(\underline{P})$.

③ Inside this "small and high probability" set, all sequences are roughly equiprobable (no typical sequence is favorable over another). Indeed, if $x^n \in T_{\epsilon}^{(n)}(\underline{P})$

$$P^{(\otimes)}(\{x^n\}) \approx 2^{-nH(\underline{P})} = \frac{1}{|T_{\epsilon}^{(n)}(\underline{P})|}$$

which is approximately the uniform distribution on $T_{\epsilon}^{(n)}(\underline{P})$.