

Lecture 13: Method of Types

Lecturer: Tsachy Weissman

Scribe: Fang Cai, Rob Jones, Yi Sun, Can Wang

In last lecture we framed the problem of lossy compression and gave the main theorem that characterizes the tradeoff between the rate and the distortion in the context of lossy compression. With that as our motivation, this week we are going to talk about the method of types, expand our tools related to typical sequences, the notion of strong typicality and the notion of conditional types, which are not only interesting in their own right, but also serve very well when we go back to establish the main result in lossy compression. We are also going to talk about some concrete schemes for lossy compression and how they are related to clustering and machine learning problems. Today we are going to talk about Method of Types.

1 Notation

Denote: $x^n = \{x_1, \dots, x_n\}$ with $x_i \in \mathcal{X} = \{1, \dots, r\}$ and

$$N(a|x^n) = \sum_{i=1}^n \mathbb{I}_{\{x_i=a\}},$$

$$P_{x^n}(a) = \frac{N(a|x^n)}{n}.$$

2 Empirical distribution and type class

Definition 1 (Empirical distribution, type class). The **empirical distribution** of x^n is the probability vector $(P_{x^n}(1), \dots, P_{x^n}(r))$. \mathbb{P}_n denotes the collection of all empirical distribution of sequences of length n , i.e. $\mathbb{P}_n = \{P_{x^n} : x^n \in \mathcal{X}^n\}$. For $P \in \mathbb{P}_n$, the **type class** or **type** of P is $T(P) = \{x^n : P_{x^n} = P\}$. The type class of x^n is $T_{x^n} = T(P_{x^n}) = \{\tilde{x}^n : P_{\tilde{x}^n} = P_{x^n}\}$.

Example 2. If $\mathcal{X} = \{0, 1\}$, then $\mathbb{P}_n = \{(1, 0), (\frac{n-1}{n}, \frac{1}{n}), (\frac{n-2}{n}, \frac{2}{n}), \dots, (0, 1)\}$

Example 3. If $\mathcal{X} = \{a, b, c\}$, $n = 5$ and $x^n = (a, a, c, b, a)$, then $P_{x^n} = (\frac{3}{5}, \frac{1}{5}, \frac{1}{5})$,

$T_{x^n} = \{(a, a, a, b, c), (a, a, a, c, b), \dots, (c, b, a, a, a)\}$ and $|T_{x^n}| = \binom{5}{3 \ 1 \ 1} = \frac{5!}{3! 1! 1!} = 20$.

In the following we show that the number of different type classes induced by x^n , $|\mathbb{P}_n|$, can be upper bounded by something which is polynomial in n , which doesn't increase exponentially with n .

Theorem 4. $|\mathbb{P}_n| \leq (n+1)^{r-1}$.

Proof

Every empirical distribution P_{x^n} is determined by vector $(N(1|x^n), N(2|x^n), \dots, N(r-1|x^n))$, where $N(a|x^n)$ means the number of times that the symbol a appears in the sequence x^n . Since $0 \leq N(a|x^n) \leq n$, each of $N(a|x^n)$ can take on no more than $n+1$ values.

Thus we have a vector of length $r-1$ and each element can take no more than $n+1$ values. Therefore there are at most $(n+1)^{r-1}$ possibilities. \square

Note that for the case $r = 2$, the bound is tight. But the bound is not tight for the cases $r \geq 3$ because we didn't incorporate the constraint that $\sum_{a=1}^{r-1} N(a|x^n)$ must be less than or equal to n when calculating the upper bound.

Further notation

- For probability mass function (PMF) $Q = \{Q(x)\}_{x \in \mathcal{X}}$, we write $H(Q)$ for $H(X)$ when X is distributed according to Q .

- $Q(x^n) = \prod_{i=1}^n Q(x_i)$. For $S \subset \mathcal{X}^n$, $Q(S) = \sum_{x^n \in S} Q(x^n)$

Theorem 5. $\forall x^n, Q(x^n) = 2^{-n[H(P_{x^n}) + D(P_{x^n} \| Q)]}$, where $H(P_{x^n})$ is referred to as empirical entropy of x^n .

Proof

$$\begin{aligned}
 Q(x^n) &= \prod_{i=1}^n Q(x_i) \\
 &= 2^{\sum_{i=1}^n \log Q(x_i)} \\
 &= 2^{\sum_{a \in \mathcal{X}} N(a|x^n) \log Q(a)} \\
 &= 2^{-n \left[\sum_{a \in \mathcal{X}} \frac{N(a|x^n)}{n} \log \frac{1}{Q(a)} \right]} \\
 &= 2^{-n \left[\sum_{a \in \mathcal{X}} P_{x^n}(a) \log \left(\frac{1}{Q(a)} \frac{P_{x^n}(a)}{P_{x^n}(a)} \right) \right]} \\
 &= 2^{-n[H(P_{x^n}) + D(P_{x^n} \| Q)]}
 \end{aligned}$$

□

The next result is about the size of the type class associated with the empirical distribution P .

Theorem 6. $\forall P \in \mathbb{P}_n, \frac{1}{(n+1)^{r-1}} 2^{nH(P)} \leq |T(P)| \leq 2^{nH(P)}$.

Note: We could calculate the size of type class $|T(P)|$ exactly, which is

$$|T(P)| = \binom{n}{n \cdot P(1), n \cdot P(2), \dots, n \cdot P(r)}.$$

But for our purposes, what we care about are (1) the behavior of this quantity for n large on an exponential scale and (2) how it is related to quantities that are familiar and important to us, such as entropy.

Proof of upper bound in Theorem 6:

$$\begin{aligned}
 1 &\geq P(T(P)) \\
 &= \sum_{x^n \in T(P)} P(x^n) \\
 &= \sum_{x^n \in T(P)} 2^{-n[H(P_{x^n}) + D(P_{x^n} \| P)]} && \text{(by Theorem 5, with } Q=P) \\
 &= \sum_{x^n \in T(P)} 2^{-n[H(P) + D(P \| P)]} && \text{(all elements } x^n \in T(P) \text{ have empirical distribution } P) \\
 &= |T(P)| \cdot 2^{-nH(P)}
 \end{aligned}$$

Thus by simple algebraic manipulation we have:

$$|T(P)| \leq 2^{nH(P)}$$

□

Before proving the lower bound, we prove two lemmas.

Lemma 7. For non-negative integers m, n , $\frac{m!}{n!} \geq n^{m-n}$.

Proof:

If $m \geq n$,

$$\frac{m!}{n!} = \underbrace{m(m-1) \cdots (n+1)}_{(m-n) \text{ factors, each } \geq n} \geq n^{m-n}.$$

If $m < n$,

$$\frac{m!}{n!} = \frac{1}{\underbrace{n(n-1) \cdots (m+1)}_{(n-m) \text{ factors, each } \leq n}} \geq \frac{1}{n^{n-m}} = n^{m-n}.$$

□

Lemma 8. $\forall P, Q \in \mathbb{P}_n, P(T(P)) \geq P(T(Q))$.

Proof:

$$\begin{aligned} \frac{P(T(P))}{P(T(Q))} &= \frac{|T(P)| \left(\prod_{a \in \mathcal{X}} P(a)^{nP(a)} \right)}{|T(Q)| \left(\prod_{a \in \mathcal{X}} P(a)^{nQ(a)} \right)} \\ &= \frac{\binom{n}{nP(1), nP(2), \dots, nP(r)}}{\binom{n}{nQ(1), nQ(2), \dots, nQ(r)}} \prod_{a \in \mathcal{X}} P(a)^{nP(a) - nQ(a)} \\ &= \prod_{a \in \mathcal{X}} \frac{(nQ(a))!}{(nP(a))!} P(a)^{n[P(a) - Q(a)]} \\ &\geq \prod_{a \in \mathcal{X}} (nP(a))^{nQ(a) - nP(a)} P(a)^{n[P(a) - Q(a)]} && \text{(by Lemma 7)} \\ &= \prod_{a \in \mathcal{X}} n^{n[Q(a) - P(a)]} \\ &= n^{n \sum_{a \in \mathcal{X}} (Q(a) - P(a))} = 1 \end{aligned}$$

□

Proof of lower bound in Theorem 6:

$$\begin{aligned} 1 &= P(\mathcal{X}^n) \\ &= \sum_{Q \in \mathbb{P}_n} P(T(Q)) \\ &\leq |\mathbb{P}_n| \cdot \max_{Q \in \mathbb{P}_n} P(T(Q)) \\ &= |\mathbb{P}_n| \cdot P(T(P)) && \text{(by Lemma 8)} \\ &= |\mathbb{P}_n| \cdot |T(P)| \cdot 2^{-n[H(P) + D(P||P)]} && \text{(by Theorem 5)} \\ &\leq (n+1)^{r-1} \cdot |T(P)| \cdot 2^{-nH(P)} && \text{(by Theorem 4)} \end{aligned}$$

Thus by simple algebraic manipulation we have:

$$\frac{1}{(n+1)^{r-1}} \cdot 2^{nH(P)} \leq |T(P)|$$

□

Noting that by Theorem 5, we have that, for any probability mass function Q and any empirical distribution $P \in \mathbb{P}_n$,

$$Q(T(P)) = |T(P)|2^{-n[H(P)+D(P||Q)]}.$$

Together with Theorem 6, we obtain the following theorem.

Theorem 9. $\forall PMF Q, \forall P \in \mathbb{P}_n, \frac{1}{(n+1)^{r-1}}2^{-nD(P||Q)} \leq Q(T(P)) \leq 2^{-nD(P||Q)}.$

This shows that up to an insignificant polynomial factor ($\frac{1}{(n+1)^{r-1}}$), on an exponential scale, the probability that the sequence looks like it came from source P , if the data is generated i.i.d. from distribution Q , is exponentially unlikely. The farther away P is from Q , the more unlikely it is. Note in the expression above, the relative entropy $D(P||Q)$ is between P , the “wrong” source, and Q , the true source, unlike in the cost of mismatch in lossless compression $D(p||q)$ (see lecture 6), p is the true source while q is the “wrong” source.