LECTURE 5 METHOD OF TYPES AND APPLICATIONS

Information Theory

Thomas Debris-Alazard

Inria, École Polytechnique

Relation between information theory and probability theory

 \longrightarrow The method of types

Powerful technique:

- Probability of rare events (large deviations)
- Universal source coding
- Testing hypothesis
- etc...

Thought experiment:

Given i.i.d $X_i \in \mathcal{X}$ according to X, we want to estimate the $\mathbb{P}(X = a)$'s

A natural approach, observe a sequence **x** of length *n* and compute the empirical distribution:

$$\mathbb{P}(\mathsf{X}=a) \approx \frac{\sharp\{i \in [1,n]: x_i = a\}}{n}$$

 \rightarrow By the weak law of large number (AEP) we know that our estimation will tend to the right one when *n* large enough

How large should be n?

What is the exact probability to make mistakes?

AEP has been a powerful tool but it does not help for rare events!

- 1. Method of Types
- 2. Alternative Law of Large Numbers
- 3. Universal Coding
- 4. Large Deviation Theory
- 5. Chernoff's Bound
- 6. Sanov's Theorem
- 7. Some Applications of Sanov's Theorem

METHOD OF TYPES

AEP: what are the typical sequences? Their probability to appear is given by the entropy!

→ Crude tool in many situations!

Method of types: split sequences according to their empirical distribution (the type)

----- The event of interest is partitioned into its intersections with the type classes. But,

- 1. The number of types is polynomial
- 2. There are an exponential number of sequences
- 3. All sequences in a type are equiprobable (memoryless source)

The event probability has the same exponential asymptotics as the largest one among the probabilities of the above intersections!

We will consider random variables X_1, \ldots, X_n from an alphabet $\mathcal{X} = \{a_1, \ldots, a_{\sharp \mathcal{X}}\}$

Any vector
$$\mathbf{x}$$
 (bold letter) denotes a sequence $x_1, \ldots, x_n \in \mathcal{X}$

If X_1, \ldots, X_n are i.i.d random variables distributed according to Q, *i.e.* $\mathbb{P}(X_i = a) = Q(a)$, then

$$\forall \mathbf{x} \in \mathcal{X}^n, \quad Q^n(\mathbf{x}) = \prod_{i=1}^n \mathbb{P}(\mathbf{X}_i = x_i) = \prod_{i=1}^n Q(x_i)$$

Furthermore, given some event \mathcal{E} ,

$$\mathbb{P}_{\mathbf{x}\leftarrow Q^n}(\mathbf{x}\in \mathcal{E}) \stackrel{\text{def}}{=} \sum_{\mathbf{x}\in \mathcal{E}} Q^n(\mathbf{x}) = \sum_{\mathbf{x}\in \mathcal{E}} \prod_{i=1}^n Q(x_i)$$

Type:

Given $\mathbf{x} \in \mathcal{X}^n$ for some n > 0, a type $P_{\mathbf{x}}^{emp}$ is a probability distribution over \mathcal{X} defined as:

$$\forall a \in \mathcal{X}, \ P_{\mathbf{x}}^{\mathsf{emp}}(a) \stackrel{\mathsf{def}}{=} \frac{\sharp\{i \in [1, n]: \ x_i = a\}}{n}$$

(A type $P_{\mathbf{x}}^{emp}$ is also called *empirical distribution of* \mathbf{x})

An example:

 $\mathcal{X} = \{0, 1\}$ and $\mathbf{x} = (1, 1, 1, 0, 1, 0, 0, 1) \in \{0, 1\}^8$,

$$P_{\mathbf{x}}^{emp}(0) = \frac{3}{8}$$
 and $P_{\mathbf{x}}^{emp}(1) = \frac{5}{8}$

 \longrightarrow Be careful: a type \textit{P}^{emp}_{x} is defined according to the sequence x

Type of fixed length:

Given n > 0,

$$\mathcal{P}_n \stackrel{\text{def}}{=} \left\{ P_{\mathbf{x}}^{\text{emp}} : \mathbf{x} \in \mathcal{X}^n \right\}$$

The binary case:

Given $\mathcal{X} = \{0, 1\}$, types of length *n*:

$$\mathcal{P}_{n} = \left\{ \left(\underbrace{\frac{0}{n}}_{p(0)}, \underbrace{\frac{1}{n}}_{p(1)}\right), \left(\underbrace{\frac{1}{n}}_{p(0)}, \underbrace{\frac{n-1}{n}}_{p(1)}\right), \ldots, \left(\underbrace{\frac{n}{n}}_{p(0)}, \underbrace{\frac{0}{n}}_{p(1)}\right) \right\}$$

Type class:

Given $P \in \mathcal{P}_n$, its type class is,

$$T(P) \stackrel{\text{def}}{=} \left\{ \mathbf{x} \in \mathcal{X}^n : P_{\mathbf{x}}^{\text{emp}} = P \right\}$$

The type class is the set of vectors having the same empirical distribution

Exercise:

Let $\mathcal{X} = \{0, 1\}$ and $\mathbf{x} = (1, 1, 1, 0, 1, 0, 0, 1) \in \mathcal{X}^8$. Describe $T(P_{\mathbf{x}}^{emp})$.

The number of class type is polynomial!

Proposition:

$$\sharp \mathcal{P}_n \leq (n+1)^{\sharp \mathcal{X}}$$

Proof:

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

$$P_{\mathbf{x}}^{emp} \in \underbrace{\left\{\frac{i}{n}: 0 \le i \le n\right\}}_{p(a_0)} \times \cdots \times \underbrace{\left\{\frac{i}{n}: 0 \le i \le n\right\}}_{p(a_{\sharp,\mathcal{X}})}$$

≇ X times

Therefore,

$$\sharp \mathcal{P}_n \leq \underbrace{\sharp \left\{ \frac{i}{n} : 0 \leq i \leq n \right\}}_{\sharp \mathcal{X} \text{ times}} \times \cdots \times \sharp \left\{ \frac{i}{n} : 0 \leq i \leq n \right\}}_{\sharp \mathcal{X} \text{ times}} = (n+1)^{\sharp \mathcal{X}}$$

$$\mathcal{P}_n = \left\{ P_{\mathbf{x}}^{\mathsf{emp}} : \mathbf{x} \in \mathcal{X}^n \right\}$$

 $\sharp \mathcal{P}_n \leq (n+1)^{\sharp \mathcal{X}}$

 \longrightarrow There is a polynomial number of types of length n

But there are $\# \mathcal{X}^n = 2^{n \log_2 \# \mathcal{X}}$ sequences (exponential)

Pigeonhole principle:

It exists an exponential number of sequences having the same type!

Exercise:

Let $\mathcal{X} = \{0, 1\}$. How many different types $P_{\mathbf{x}}^{emp}$ exist? How many sequences have a fixed type?

Kullback-Leiber divergence:

$$D_{\mathsf{KL}}(P||Q) = \sum_{a} P(a) \log_2 \frac{P(a)}{Q(a)}$$

Theorem:

Let X_1, \ldots, X_n be i.i.d according to Q, then

$$Q^{n}(\mathbf{x}) = 2^{-n\left(H(P_{\mathbf{x}}^{\text{emp}}) + D_{\text{KL}}(P_{\mathbf{x}}^{\text{emp}}) | Q\right)}$$

Furthermore, if
$$Q \in \mathcal{P}_n$$
 and $\mathbf{x} \in T(Q)$ (where $T(Q) = \{\mathbf{x} \in \mathcal{X}^n : Q_{\mathbf{x}}^{emp} = Q\}$)
 $Q^n(\mathbf{x}) = 2^{-nH(Q)}$

 \longrightarrow It shows that, when considering a sequence **x** with its associated empirical distribution,

```
i.e., P_{\mathbf{x}}^{\text{emp}}, what we loose is D_{\text{KL}}\left(P_{\mathbf{x}}^{\text{emp}}||Q\right)
```

(with the AEP, a typical event happens with probability $2^{-nH(X)}$)

Proof: $Q^{n}(\mathbf{x}) = \prod_{a \in \mathcal{X}} Q(a)^{\sharp \{i \in [1,n]: x_{i}=a\}}$ $= \prod_{a \in \mathcal{X}} Q(a)^{n P_{\mathbf{x}}^{emp}(a)}$ $= \prod_{a \in \mathcal{X}} 2^{n \left(P_{\mathbf{x}}^{emp}(a) \log_{2} Q(a) - P_{\mathbf{x}}^{emp}(a) \log_{2} P_{\mathbf{x}}^{emp}(a) + P_{\mathbf{x}}^{emp}(a) \log_{2} P_{\mathbf{x}}^{emp}(a)\right)}$ Therefore, $Q^{n}(\mathbf{x}) = 2^{-n \left(H(P_{\mathbf{x}}^{emp}) + D_{\mathrm{KL}}(P_{\mathbf{x}}^{emp} | |Q)\right)}$. To conclude use that $D_{\mathrm{KL}}(P||Q) = 0$ if P = Q

Theorem:

For any
$$P_{\mathbf{x}}^{\text{emp}} \in \mathcal{P}_{n}$$
,

$$\frac{1}{(n+1)^{\sharp \mathcal{X}}} 2^{nH(P_{\mathbf{x}}^{\text{emp}})} \leq \sharp T(P_{\mathbf{x}}^{\text{emp}}) \leq 2^{nH(P_{\mathbf{x}}^{\text{emp}})}$$
where $T(P_{\mathbf{x}}^{\text{emp}}) = \{\mathbf{y} \in \mathcal{X}^{n} : P_{\mathbf{y}}^{\text{emp}} = P_{\mathbf{x}}^{\text{emp}}\}$

Proof:

Let $P \stackrel{\text{def}}{=} P_{\mathbf{x}}^{\text{emp}}$,

$$1 \ge \sum_{\mathbf{x} \in T(P)} P^{n}(\mathbf{x}) \underset{\text{prev th.}}{=} \sum_{\mathbf{x} \in T(P)} 2^{-nH(P)} = \#T(P) \ 2^{-nH(P)} \text{ which gives the upper bound}$$

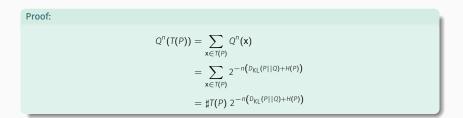
To derive the lower bound, let us admit: $\forall Q \in \mathcal{P}_n, P(T(Q)) \leq P(T(P)).$

$$1 = \sum_{Q \in \mathcal{P}_n} P^n(T(Q)) \le \sum_{Q \in \mathcal{P}_n} P^n(T(P)) \le (n+1)^{\sharp \mathcal{X}} \sum_{\mathbf{x} \in T(P)} P^n(\mathbf{x}) = (n+1)^{\sharp \mathcal{X}} \sum_{\mathbf{x} \in T(P)} 2^{-nH(P)}$$

Theorem:

For any $P \in \mathcal{P}_n$ we have,

$$\frac{1}{(n+1)^{\sharp \mathcal{X}}} 2^{-nD_{\mathsf{KL}}(P||Q)} \leq Q^n \left(T(P)\right) = \underset{\mathbf{x} \leftarrow Q^n}{\mathbb{P}} \left(P_{\mathbf{x}}^{\mathsf{emp}} = P\right) \leq 2^{-nD_{\mathsf{KL}}(P||Q)}$$



Moral: sequences with empirical distribution P appear under the distribution Q with an exponentially small probability, whose exponent is given by $D_{KL}(P, Q)$

- $\sharp \mathcal{P}_n \leq (n+1)^{\sharp \mathcal{X}}$ (polynomial number of types)
- $\sharp T\left(P_x^{emp}\right) \stackrel{(\text{poly})}{=} 2^{nH\left(P_x^{emp}\right)}$ (exponential number of sequences in each type)
- $Q^{n}(\mathbf{x}) = 2^{-n\left(H\left(P_{\mathbf{x}}^{emp}\right) + D_{KL}\left(P_{\mathbf{x}}^{emp}||Q\right)\right)}$

(probability of sequence of some type under Q)

•
$$Q^{n}\left(T\left(P_{\mathbf{x}}^{emp}\right)\right) \stackrel{(\text{poly})}{=} 2^{-nD_{\text{KL}}\left(P_{\mathbf{x}}^{emp}||Q\right)}$$

Exercise:

Explicit and compute the above equations when $\mathcal{X} = \{0,1\}$

 \longrightarrow These results admit many consequences that we will describe now!

ALTERNATIVE LAW OF LARGE NUMBERS

Types and type classes offer an alternative "statement" of the law of large numbers!

The crucial property:

Polynomial number of types and an exponential number of sequences of each type But the probability of each type class T(P) depends exponentially on $D_{KL}(P||Q_{true \ distrib.})$

New concept of typical sequences: close for the KL-divergence

Typical Set:

Let $\varepsilon > 0$ and a distribution Q,

$$T_{Q}^{(\varepsilon)} = \left\{ \mathbf{x} \in \mathcal{X}^{n} : D_{\mathsf{KL}}(P_{\mathbf{x}}^{\mathsf{emp}} | | Q) \le \varepsilon \right\}$$

 \rightarrow The probability of not being "empirical typical" is exponentially small! (similar to AEP)

Probability of not being typical:

$$\mathbb{P}_{\varsigma \leftarrow Q^{n}}\left(D_{\mathsf{KL}}\left(P_{\mathbf{x}}^{\mathsf{emp}}||Q\right) > \varepsilon\right) = 1 - Q^{n}(T_{Q}^{(\varepsilon)}) \leq (n+1)^{\sharp \mathcal{X}} 2^{-n\varepsilon}$$

Proof:

$$1 - Q^{n}(T_{Q}^{(\varepsilon)}) = \sum_{P \in \mathcal{P}_{n}: \ D_{\mathsf{KL}}(P \mid |Q) \ge \varepsilon} Q^{n}(T(P)) \le \sum_{P \in \mathcal{P}_{n}: \ D_{\mathsf{KL}}(P \mid |Q) \ge \varepsilon} 2^{-n}$$

To conclude the proof, use that there are a polynomial number of types!

Law of large number with KL-divergence (admitted):

 X_1, \ldots, X_n be i.i.d according to Q,

$$\mathbb{P}_{\mathbf{x} \leftarrow Q^n} \left(\mathbf{x} : D_{\mathrm{KL}}(P_{\mathbf{x}}^{\mathrm{emp}} || Q) > \varepsilon \right) \le (n+1)^{\sharp \mathcal{X}} 2^{-n\varepsilon}$$

and, if $\mathbf{x}^{(n)} \leftarrow Q^n$, then $D_{KL}(P_{\mathbf{x}}^{emp}||Q) \xrightarrow[n \to +\infty]{} 0$ almost surely, *i.e.*, $\forall \varepsilon > 0, \quad \mathbb{P}\left(\lim_{n \to +\infty} D_{KL}\left(P_{\mathbf{x}^{(n)}}||Q\right) \le \varepsilon\right) = 1$

UNIVERSAL CODING

Huffman compresses source with known distribution X with an amount of bits given by the entropy H(X)

 \rightarrow if instead a distribution Y is assumed: a penalty of $D_{KL}(X||Y)$ is incurred!

What compression can be achieved if the true distribution X is unknown?

Universal coding:

A symbol code $\varphi : \mathcal{X}^n \to \{0,1\}^{nR}$ is said to be 2^{nR} -universal if we can decode it with probability tending to one, *i.e.* it exists,

$$\mathsf{Dec}: {\{0,1\}}^{nR} o \mathcal{X}^n$$

such that independently of the memoryless source distribution Q,

$$\mathsf{P}_{e}^{(n)} \stackrel{\text{def}}{=} \mathop{\mathbb{P}}_{\mathsf{x} \leftarrow \mathcal{Q}^{n}} \left(\mathsf{Dec}(\varphi(\mathsf{x})) \neq \mathsf{x} \right) \xrightarrow[n \to +\infty]{} 0$$

Consequence of AEP: Shannon source coding theorem but for known source distribution!

Is type method enables to prove the stronger statement that universal source coding exits?

Universal source coding: it is possible!

There exists a sequence (in n) of 2^{nR} -universal code such that $P_e^{(n)} \xrightarrow[n \to +\infty]{} 0$ independently of the memoryless source distribution Q as soon as R > H(Q)

PROOF (I): ENCODING AND DECODING

Proof:

Let $R_n \stackrel{\text{def}}{=} R - \# \mathcal{X} \stackrel{\log_2(n+1)}{n}$. We consider sequences: $A = \{ \mathbf{x} \in \mathcal{X}^n : H(P_{\mathbf{x}}^{emp}) \le R_n \}$ Then, $\# A = \sum_{P \in \mathcal{P}_n : H(P) \le R_n} \# T(P)$ $\leq \sum_{P \in \mathcal{P}_n : H(P) \le R_n} 2^{nH(P)}$ $\leq \sum_{P \in \mathcal{P}_n : H(P) \le R_n} 2^{nR_n}$ $\leq (n+1)^{\# \mathcal{X}} 2^{nR_n} \quad (\text{the method of types in action!})$

$$\varphi(\mathbf{x}) = \begin{cases} \text{index of } \mathbf{x} \text{ in } A & \text{if } \mathbf{x} \in A \\ \bot & \text{otherwise} \end{cases}$$

The compression size asks $\log_2 \# A \le nR$ bits!

 $= 2^{nR}$

Decoding: map an index to its corresponding element

Proof:

But Rn

 X_1, \ldots, X_n be i.i.d according to Q where H(Q) < R. The probability to make a mistake during decoding verifies,

$$P_{e}^{(n)} = 1 - Q^{n}(A)$$

$$= \sum_{P \in \mathcal{P}_{n}: \ H(P) > R_{n}} Q^{n}(T(P))$$

$$\leq (n+1)^{\#\mathcal{X}} \sum_{P: \ H(P) > R_{n}} Q^{n}(T(P))$$

$$\leq (n+1)^{\#\mathcal{X}} 2^{-n} \sum_{P: \ H(P) > R_{n}} D_{\mathsf{KL}}(P||Q) \qquad (the method of types in action!)$$

$$\xrightarrow{n \to +\infty} R \text{ and } H(Q) < R. \text{ Therefore, for } n \text{ sufficiently large } H(Q) < R_{n} \text{ and}$$

 $H(P) > R_n \Longrightarrow H(P) > H(Q)$

from which we conclude that $P \neq Q$ and by Gibb's inequality $D_{KL}(P||Q) > 0$

Some remark:

The error in the decoding tends to 0 exponentially fast in the code-length

Universal source coding is a huge topic

 Elements of Information Theory, Universal Source Coding Chapter 13, by M. Cover & Joy A. Thomas

LARGE DEVIATION THEORY

I repeated too often: a random variable is equal to its expectation...

But, if yes, why? If you don't believe me, how to convince you that I say the truth?

 \longrightarrow Let us study $\mathbb{P}\left(X \gg \mathbb{E}(X)\right)$

First two approaches:

- Markov inequality
- Bienaymé-Tchebychev inequality

Markov's Inequality:

Given $X : \Omega \longrightarrow \mathbb{R}_+$ and t > 0,

$$\mathbb{P}\left(\mathsf{X} \geq t\right) \leq \frac{\mathbb{E}(\mathsf{X})}{t}$$

Bienaymé-Tchebychev

Given t > 0,

$$\mathbb{P}\left(|\mathsf{X} - \mathbb{E}(\mathsf{X})| \ge t\right) \le \frac{\mathbb{V}(\mathsf{X})}{t^2}$$

 \longrightarrow Are these inequalities tight?

We know random variables s.t the above probabilities reach the inequalities (exercise)

Markov and Bienaymé-Tchebychev are worst-case bounds (true for any random variable)

 \longrightarrow Goal of large deviation theory: provide better bounds for a given family of random variables!

Given $X_1, \ldots, X_n \in \{0, 1\}$ be i.i.d with $\mathbb{P}(X_i = 1) = p$.

 $\mathbf{X}^{(n)} \stackrel{\text{def}}{=} \sum_{i=1}^{n} \mathbf{X}_{i}$

$$\mathbb{E}\left(\mathbf{X}^{(n)}\right) = np$$
 and $\mathbb{V}\left(\mathbf{X}^{(n)}\right) = np(1-p)$

Bienaymé-Tchebychev:

Let $\varepsilon > 0$,

$$\mathbb{P}\left(\left|\mathbf{X}^{(n)} - np\right| \geq \varepsilon n\right) \leq \frac{1}{\varepsilon n} p(1-p) \xrightarrow[n \to +\infty]{} 0$$

It tends to 0 as $1/(\varepsilon n)$, is it the best that we can expect?

Given $X_1, \ldots, X_n \in \{0, 1\}$ be i.i.d with $\mathbb{P}(X_i = 1) = p$.

 $\mathbf{X}^{(n)} \stackrel{\text{def}}{=} \sum_{i=1}^{n} \mathbf{X}_{i}$

$$\mathbb{E}\left(\mathbf{X}^{(n)}\right) = np \text{ and } \mathbb{V}\left(\mathbf{X}^{(n)}\right) = np(1-p)$$

Bienaymé-Tchebychev:

Let $\varepsilon > 0$,

$$\mathbb{P}\left(\left|\mathbf{X}^{(n)} - np\right| \geq \varepsilon n\right) \leq \frac{1}{\varepsilon n} p(1-p) \xrightarrow[n \to +\infty]{} 0$$

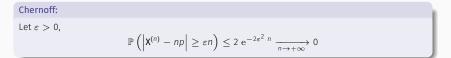
It tends to 0 as $1/(\varepsilon n)$, is it the best that we can expect?

 $\rightarrow No!$

CHERNOFF'S BOUND

G

$$\begin{aligned} \text{ (ven } \mathbf{X}_1, \dots, \mathbf{X}_n \in \{0, 1\} \text{ be i.i.d with } \mathbb{P}(\mathbf{X}_i = 1) &= p \\ \mathbf{X}^{(n)} \stackrel{\text{def}}{=} \sum_{i=1}^n \mathbf{X}_i \\ \mathbb{E}\left(\mathbf{X}^{(n)}\right) &= np \quad \text{and} \quad \mathbb{V}\left(\mathbf{X}^{(n)}\right) = np(1-p) \end{aligned}$$

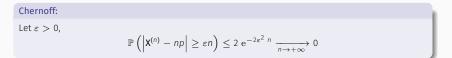


It tends to 0 as $e^{-2\varepsilon^2 n}$: exponentially better than $1/(\varepsilon n)$ But is $-2\varepsilon^2 n$ the best exponent that we can expect?

CHERNOFF'S BOUND

G

Where
$$\mathbf{X}_1, \dots, \mathbf{X}_n \in \{0, 1\}$$
 be i.i.d with $\mathbb{P}(\mathbf{X}_i = 1) = p$
$$\mathbf{X}^{(n)} \stackrel{\text{def}}{=} \sum_{i=1}^n \mathbf{X}_i$$
$$\mathbb{E}\left(\mathbf{X}^{(n)}\right) = np \text{ and } \mathbb{V}\left(\mathbf{X}^{(n)}\right) = np(1-p)$$



It tends to 0 as $e^{-2\varepsilon^2 n}$: exponentially better than $1/(\varepsilon n)$ But is $-2\varepsilon^2 n$ the best exponent that we can expect?

 \longrightarrow Yes as we show now thanks to the method of types!

CHERNOFF'S BOUND

First approach: central limit theorem

 \rightarrow Poor estimations in many cases (try with the binomial distribution)

Our approach: method of type!

 $\mathbf{x}_1, \ldots, \mathbf{x}_n \in \{0, 1\}$ be i.i.d with $\mathbb{P}(\mathbf{x}_i = 1) = \frac{1}{3}$. Crucial remark: if,

 $\frac{1}{n}\sum_{i=1}^{n}\mathbf{x}_{i}=\frac{3}{4},$ then $P_{\mathbf{x}}^{\text{emp}}=\left(\frac{1}{4},\frac{3}{4}\right)$

 \rightarrow We expect (by method of types) to obtain such sequences with probability $\approx 2^{-nD_{KL}}((\frac{1}{4},\frac{3}{4})||(\frac{2}{3},\frac{1}{3}))$

(exponentially small and we know that the exponent cannot be smaller)

We know the optimal exponent. . .

Theorem: For any $P \in \mathcal{P}_n$ we have, $\frac{1}{(n+1)^{\sharp \mathcal{X}}} 2^{-nD_{\mathsf{KL}}(P||Q)} \le Q^n (T(P)) = \mathbb{P}_{\mathbf{x} \leftarrow Q^n} \left(P_{\mathbf{x}}^{emp} = P \right) \le 2^{-nD_{\mathsf{KL}}(P||Q)}$

$$Q = \begin{cases} 0 \text{ with probability } 1 - p \\ 1 \text{ with probability } p \end{cases} \text{ and } P = \begin{cases} 0 \text{ with probability } 1 - p - \varepsilon \\ 1 \text{ with probability } p + \varepsilon \end{cases}$$

Some computation:

For $\varepsilon > 0$ small enough,

$$D_{\mathrm{KL}}(P||Q) = \frac{2\varepsilon^2}{\ln(2)} + o(1)$$

The crucial remark:

$$\mathbb{P}_{\mathbf{x}\leftarrow Q^n}\left(X_1+\cdots+X_n=np+\varepsilon n\right)=\mathbb{P}_{\mathbf{x}\leftarrow Q^n}\left(P_{\mathbf{x}}^{emp}=P\right)=Q^n(T(P))$$

From the previous theorem,

$$\tfrac{1}{(n+1)^2}\mathrm{e}^{-2n\varepsilon^2(1+o(1))} \leq \underset{\mathbf{x} \leftarrow Q^n}{\mathbb{P}} \left(x_1 + \cdots + x_n = np + \varepsilon n \right) \leq \mathrm{e}^{-2n\varepsilon^2(1+o(1))}$$

We almost recover the optimality of Chernoff's bound! We want to know sufficiently large

deviation, not just the deviation exactly equal to $+\varepsilon n$

• By the previous bound, for all η (integer)

$$\frac{1}{(n+1)^2} \mathrm{e}^{-2n(\varepsilon+\eta/n)^2(1+o(1))} \leq \mathbb{P}_{\mathbf{x}\leftarrow \mathcal{Q}^n} \left(x_1 + \cdots + x_n = np + \varepsilon n + \eta \right) \leq \mathrm{e}^{-2n(\varepsilon+\eta/n)^2(1+o(1))}$$

• By summing all possible η (polynomial number of types)

$$\frac{1}{(n+1)}\mathrm{e}^{-2n\varepsilon^2(1+o(1))} \leq \mathop{\mathbb{P}}_{\mathbf{x}\leftarrow Q^n} \left(x_1 + \cdots + x_n \geq np + \varepsilon n \right) \leq (n+1)\mathrm{e}^{-2n\varepsilon^2(1+o(1))}$$

We can obtain the same bound for

 $\mathbb{P}_{\mathbf{x}\leftarrow Q^n} (x_1 + \dots + x_n \leq np - \varepsilon n)$ by replacing $\varepsilon \longleftrightarrow -\varepsilon$

$$Q = \begin{cases} 0 \text{ with probability } 1 - p \\ 1 \text{ with probability } p \end{cases}$$

$$\frac{1}{n}\log_2 \mathbb{P}_{\mathbf{x}\leftarrow Q^n}\left(|x_1+\cdots+x_n-np|\geq \varepsilon n\right) = -2\varepsilon^2/\ln(2) + o(1)$$

We obtained the best exponent by using method of types!

To obtain the "optimality" of Chernoff's bound we made the following reasoning

- 1. Start from the distribution Q
- 2. Our aim: to get an upper bound on $\mathbb{P}_{\mathbf{x} \leftarrow Q^n} \left(\sum_i x_i = \mathbb{E}(\mathbf{X}) + \alpha \right)$
- 3. To this aim we introduced the distribution P s.t $P_x^{emp} = P$ if and only if $\sum_i x_i = \mathbb{E}(X) + \alpha$

We are going to systematize this approach: Sanov's theorem!

SANOV'S THEOREM

Given i.i.d random variables X_i 's distributed as Q, we want to estimate,

$$\mathbb{P}_{\mathbf{x}\leftarrow Q^n}\left(\frac{1}{n}\sum_{i=1}^n g(\mathbf{x}_i) \geq \alpha\right)$$

Fundamental ideal: introduce the $\mathbf{x} \leftarrow Q^n$ such that $\sum_{a \in \mathcal{X}} g(a) P^{\mathsf{emp}}_{\mathbf{x}}(a) \geq lpha$

 \longrightarrow We expect the probability exponent to behave as $D_{\text{KL}}(P_{\mathbf{x}}^{\text{emp}}||Q)$

The fundamental ideal: introduce the $\mathbf{x} \leftarrow Q^n$ such that $\sum_{a \in \mathcal{X}} g(a) \mathcal{P}^{\mathsf{emp}}_{\mathbf{x}}(a) \geq \alpha$

 \longrightarrow We expect the probability exponent to behave as $D_{\text{KL}}\left(P_{\mathbf{x}}^{\text{emp}}||Q\right)$

Issue:
There are many
$$P_x^{emp}$$
's (from different classes) verifying $\sum_{a \in \mathcal{X}} g(a) P_x^{emp}(a) \ge \alpha$

 \longrightarrow We will show that the exponent behaves as $D_{KL}(P^*||Q)$ for P^* minimizing $D_{KL}(P||Q)$ for the

$$P \in \mathcal{P}_n$$
's verifying $\sum_{a \in \mathcal{X}} g(a) P(a) \ge \alpha$

(the minimization is here to "extract" the dominant exponential term)

Goal: find the "closest" P in the constraint set for the KL-divergence to obtain the exponent!

 \longrightarrow We need to define the topology associated to this "KL-distance"

TOPOLOGY

Probability simplex:

Subset of $[0, 1]^{\sharp \mathcal{X}}$ defined as, $\mathcal{P} \stackrel{\text{def}}{=} \left\{ (x_1, \dots, x_{\sharp \mathcal{X}}) \in [0, 1]^{\sharp \mathcal{X}} : x_i \ge 0 \text{ and } \sum_{i=1}^{\sharp \mathcal{X}} x_i = 1 \right\}$

 $\mathcal{P} \subseteq \mathbb{R}^{\sharp \mathcal{X}}$, we will speak of closure, interior, \ldots But for the D_{KL} -divergence

We will identify distributions P over ${\mathcal X}$ with elements of ${\mathcal S}$

Proposition:

 $\bigcup_{n\in\mathbb{N}}\mathcal{P}_n$ is dense in \mathcal{S} the set of all distributions over \mathcal{X}

Proof:

Given
$$P = (p(a_1), \dots, p(a_m)) \in S$$
,
 $\forall i \in [1, \#\mathcal{X} - 1], n_i \stackrel{\text{def}}{=} \frac{\lfloor np(a_i) \rfloor}{n} \text{ and } n_{\#S} = \frac{1 - \sum_{i=1}^{\#\mathcal{X} - 1} n_i}{n}$
Then, $P_n^{\text{emp}} \stackrel{\text{def}}{=} (n_i)_i \in \mathcal{P}_n \text{ and } D_{\text{KL}} \left(P_n^{\text{emp}} ||P \right) \xrightarrow[n \to +\infty]{} 0$

$$P_{\mathbf{x}}^{emp} \longleftrightarrow \mathbf{x}$$

$$\mathbf{F}_{\mathbf{x}}^{emp} \longleftrightarrow \mathbf{x}$$

$$Q^{n}(E) \stackrel{\text{def}}{=} Q^{n}(E \cap \mathcal{P}_{n}) = \sum_{\mathbf{x}: P_{\mathbf{x}}^{emp} \in E \cap \mathcal{P}_{n}} Q^{n}(\mathbf{x})$$

$$Q^{n}(E) = \sum_{\mathbf{x} \leftarrow Q^{n}} (P_{\mathbf{x}}^{emp} \in E \cap \mathcal{P}_{n})$$

THE FORMALISM IN ACTION(I)

But why this formalism?

Given a random variable $\mathbf{X} = (\mathbf{X}_1, \dots, \mathbf{X}_n) \in \mathcal{X}^n$ where the \mathbf{X}_i 's are i.i.d according to Q and some

function g, we are interested in

$$\mathbb{P}_{\mathbf{x} \leftarrow Q^n} \left(\frac{1}{n} \sum_{i=1}^n g(\mathbf{x}_i) \ge \alpha \right)$$

We introduce:

$$E \stackrel{\text{def}}{=} \{ P \in \mathcal{P} : \sum_{a \in \mathcal{X}} g(a) p(a) \ge \alpha \}$$

$$\frac{1}{n}\sum_{i=1}^{n}g(x_i) \ge \alpha \iff \sum_{a\in\mathcal{X}}g(a)P_{\mathbf{x}}^{emp}(a) \ge \alpha$$
$$\iff P_{\mathbf{x}}^{emp} \in E \cap \mathcal{P}_n$$

Conclusion:

$$\mathbb{P}_{\mathbf{X}}\left(\frac{1}{n}\sum_{i=1}^{n}g(\mathbf{X}_{i})\geq\alpha\right)=Q^{n}(E)=\underset{\mathbf{x}\leftarrow Q^{n}}{\mathbb{P}}\left(P_{\mathbf{x}}^{emp}\in E\cap\mathcal{P}_{n}\right)$$

We are interested in

$$\mathbb{P}_{\mathbf{x}\leftarrow Q^n}\left(\frac{1}{n}\sum_{i=1}^n g(\mathbf{x}_i) \geq \alpha\right) = Q^n(E) = \mathbb{P}_{\mathbf{x}\leftarrow Q^n}\left(\mathsf{P}_{\mathbf{x}}^{\mathsf{emp}} \in E \cap \mathcal{P}_n\right)$$

How does Qⁿ(E) behave?

Let us use what we already know: our alternative law of large numbers (see Slide 20)

Let $E \subseteq \mathcal{P}$,

Suppose that E contains a relative entropy neighbourhood of Q, i.e.,

 $\exists \varepsilon > 0$, such that $\{P \in \mathcal{P} : D_{KL}(P||Q) < \varepsilon\} \subseteq E$

$$Q^{n}(E) = \underset{\mathbf{x} \leftarrow Q^{n}}{\mathbb{P}} \left(P_{\mathbf{x}}^{emp} \in E \cap \mathcal{P}_{n} \right) \geq \underset{\mathbf{x} \leftarrow Q^{n}}{\mathbb{P}} \left(D_{KL} \left(P_{\mathbf{x}}^{emp} | | Q \right) < \varepsilon \right) \geq 1 - (n+1)^{\sharp \mathcal{X}} 2^{-n\varepsilon} \xrightarrow[n \to +\infty]{} 1$$

Suppose that E does not contain Q or any element of some neighbourhood of Q, i.e.,

 $\exists \varepsilon \geq 0$ such that $\{P \in \mathcal{P} : D_{KL}(P||Q) \leq \varepsilon\} \cap E = \emptyset$

$$Q^{n}(E) = \underset{\mathbf{x} \leftarrow Q^{n}}{\mathbb{P}} \left(P_{\mathbf{x}}^{\text{emp}} \in E \cap \mathcal{P}_{n} \right) \leq \underset{\mathbf{x} \leftarrow Q^{n}}{\mathbb{P}} \left(D_{\text{KL}} \left(P_{\mathbf{x}}^{\text{emp}} | | Q \right) \geq \varepsilon \right) \leq (n+1)^{\sharp \mathcal{X}} 2^{-n\varepsilon} \xrightarrow[n \to +\infty]{} 0$$

It tends to zero exponentially fast, but what is the exponent?

Sanov's Theorem:

Let X_1, \ldots, X_n be i.i.d. according to Q. Let $E \subseteq \mathcal{P}$. Then,

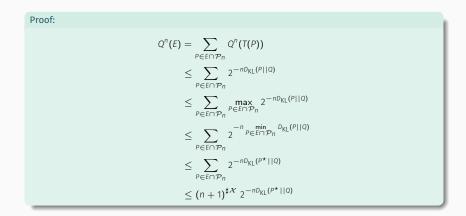
$$\mathbb{P}_{\mathbf{x}\leftarrow Q^n}\left(P_{\mathbf{x}}^{emp}\in E\cap\mathcal{P}_n\right)=Q^n(E)\leq (n+1)^{\sharp\mathcal{X}}2^{-nD_{\mathsf{KL}}\left(P^{\star}\mid\mid Q\right)}\quad\text{where }P^{\star}\stackrel{\text{def}}{=}\arg\min_{P\in E}D_{\mathsf{KL}}(P\mid\mid Q)$$

Furthermore, if E is the closure of its interior,

$$\frac{1}{n}\log_2 Q^n(E) \xrightarrow[n \to +\infty]{} D_{\mathrm{KL}}(P^*||Q)$$

Important remark:

 $D_{KL}(P^*||Q)$ is the exponent of the probability when E is the closure of its interior!



→ The last inequality shows the method of types in action: exponential probability of events versus polynomial number of events!

Proof:

Suppose that *E* is the closure of its interior. In particular, interior E° is not empty. As $\bigcup_{n} \mathcal{P}_{n}$ is dense in \mathcal{P} , then $E^{\circ} \cap \mathcal{P}_{n}$ is non-empty for *n* large enough $\geq n_{0}$. By using the density, it exists $P_{n} \in E^{\circ} \cap \mathcal{P}_{n}$ such that $D_{\mathsf{KL}}(P_{n}||Q) \xrightarrow[n \to +\infty]{} D(P^{\star}||Q)$. Furthermore,

$$\begin{aligned} Q^{n}(E) &= \sum_{P \in E \cap \mathcal{P}_{n}} Q^{n}(T(P)) \\ &\geq Q^{n}(T(P_{n})) \\ &\geq \frac{1}{(n+1)^{\sharp \mathcal{X}}} 2^{-nD_{\mathsf{KL}}(P_{n} \mid \mid Q)} \qquad (\text{the method of types in action!}) \end{aligned}$$

Therefore, combining with the upper bound,

$$-D_{\mathsf{KL}}\left(P_{n}||Q\right) - \frac{\sharp\mathcal{X}\log_{2}(n+1)}{n} \leq \frac{1}{n}\log_{2}Q^{n}(E) \leq \frac{\sharp\mathcal{X}\log_{2}(n+1)}{n} - D_{\mathsf{KL}}\left(P_{n}||Q\right)$$

APPLICATIONS OF SANOV'S THEOREM

We want to tightly derive an upper bound on:

$$\forall j \in [1, k], \mathbb{P}_{\mathbf{x} \leftarrow Q^n} \left(\frac{1}{n} \sum_{i=1}^n g_j(x_i) \ge \alpha_j \right)$$

How to proceed: use Sanov theorem!

$$E = \{P \in \mathcal{P} : \forall j \in [1, k], \sum_{a} P(a)g_j(a) \ge \alpha_j\}$$
 (closure of its interior)

The optimal exponent is given by $D_{KL}(P^*||Q)$ where P^* is given by

$$P^{\star}(a) \stackrel{\text{def}}{=} \frac{Q(a) \mathrm{e}^{\sum_{i} \lambda_{i} g_{i}(a)}}{\sum_{x \in \mathcal{X}} Q(x) \mathrm{e}^{\sum_{i} \lambda_{i} g_{i}(a)}}$$

where the λ_i 's are chosen such that

$$\forall j \in [1, k], \sum_{a} P^{\star}(a)g_j(a) = \alpha_j$$

PROOF IN A SIMPLER CONTEXT: LAGRANGE MULTIPLIERS

We want to tightly derive an upper-bound on

$$\mathbb{P}_{\mathbf{x}\leftarrow Q^n}\left(\frac{1}{n}\sum_{i=1}^n g(\mathbf{x}_i) \geq \alpha\right)$$

 $E = \left\{ P \in \mathcal{P} : \sum_{a} P(a)g(a) \ge \alpha \right\}$

We wan to minimize $D_{KL}(P||Q)$ over $P \in E!$ To this aim introduce the constraint functions:

$$c_0(\mathbf{p}) = \sum_{i=1}^n p_i - 1$$
 and $c_1(\mathbf{p}) \stackrel{\text{def}}{=} \sum_{i=1}^{\#\mathcal{X}} p_i g(a_i) - \alpha$ $(\mathbf{p} \in \mathbb{R}^{\#\mathcal{X}} \text{ the distribution vector of } P)$

► First step: minimize $f(\mathbf{p}) \stackrel{\text{def}}{=} D_{KL}(P||Q)$ for $P \in \widetilde{E} \stackrel{\text{def}}{=} \{P \in \mathcal{P} : \sum_{a} P(a)g(a) \ge \alpha\} \cap \mathbb{R}_{>0}^{\#\mathcal{X}} \subseteq E$, $c_0(\mathbf{p}) = c_1(\mathbf{p}) = 0$

Use Lagrange Multiplier Theorem:

t exists
$$\lambda, \mu$$
 such that for all $i \in [1, \sharp \mathcal{X}]$,

$$\frac{\partial f}{\partial p_i}(\mathbf{p}) = \log_2 \frac{p_i}{Q(a_i)} + \frac{1}{\ln 2} = \mu \frac{\partial c_0}{\partial p_i}(\mathbf{p}) + \frac{\partial c_1}{\partial p_i}(\mathbf{p}) = \mu + \lambda g(a_i)$$

 \longrightarrow We deduce that $p_i = Q(a_i)2^{-1/\ln(2)+\mu+\lambda g(a_i)}$ is a minimum of $D_{KL}(P||Q)$ for $P \in \widetilde{E}$

Does
$$p_i = \frac{Q(a_i)2^{\lambda g(a_i)}}{C}$$
 where $C = \sum_i Q(a_i)2^{\lambda g(a_i)}$ give the minimum of $D_{KL}(P||Q)$ for $P \in E$ as expected?

First computation:

$$D_{\mathsf{KL}}(p||Q) = \sum_{i} p_i \log_2 \frac{Q(a_i)2^{\lambda g(a_i)}}{CQ(a_i)} = \lambda \underbrace{\sum_{i} p_i g(a_i) - \log_2 C}_{=\alpha \text{ by def of } \tilde{\mathsf{E}}} C = \lambda \alpha - \log_2 C = D_{\mathsf{KL}}(p||Q)$$

$$\sum_{i} R(i) \log_2 \frac{p_i}{Q(i)} \ge \lambda \alpha - \log_2 C$$

Conclusion:

 $D_{\text{KL}}(R||Q) - D_{\text{KL}}(p||Q) \ge D_{\text{KL}}(R||Q) - \sum_{i} R(i) \log_2 \frac{p_i}{Q(i)} = D_{\text{KL}}(R||p) \ge 0$ by Gibb's inequality.

The p_i 's gives the minimum of $D_{KL}(P||Q)$ for $P \in E!$

We toss a fair dice *n* times, what is the probability that the average of the throws is greater than or equal to 4?

By Sanov's theorem,

$$Q^n(E) \stackrel{(\text{poly})}{=} 2^{-nD_{\text{KL}}(P^* ||Q)}$$

where P^* minimizes $D_{KL}(P||Q)$ over all distributions P that satisfy,

 $\sum_{i=1}^{6} iP(i) \geq 4$

$$\forall i \in [1, 6], \ P^{\star}(i) = \frac{2^{\lambda i}}{\sum_{j=1}^{6} 2^{\lambda j}} \text{ where } \lambda \text{ such that } \sum_{i=1}^{6} i P^{\star}(i) = 4$$

Solving numerically we obtain $\lambda = 0.2519$, therefore $D_{KL}(P^*||Q) = 0.0624$. After 1000 coin tosses, the probability that the average is greater than or equal to 4 is $\approx 2^{-624}$

We want to estimate the probability of observing more than 700 heads in a series of 1000 coin tosses of a fair coin

$$\mathbb{P}\left(\overline{\mathsf{X}}_{n} \geq 0.7\right) \stackrel{(\text{poly})}{=} 2^{-nD_{\mathsf{KL}}(P^{\star}||Q)}$$

where $P^* = (0.3, 0.7)$. In that case $D_{KL}(P^*||Q) = 1 - H(0.7) = 0.119$. Our probability is $\approx 2^{-119}$

We are given a joint distribution $(X, Y) \leftarrow Q(x, y)$. Let Q(x), Q(y) be the associated distributions formed by the marginals

Let Q_0 be the distribution product of marginals ("as if X and Y were independent")

 $Q_0(x,y) \stackrel{\text{def}}{=} Q(x)Q(y)$

We want to estimate the probability that $(x^n, y^n) \leftarrow Q_0^n(\cdot, \cdot)$ looks to be picked according to $Q(x^n, y^n)$

Estimating the probability that $P_{x^n,y^n} \in E \cap \mathcal{P}_n(\mathbf{X},\mathbf{Y})$ when $(x^n,y^n) \leftarrow Q^n(\cdot,\cdot)$ and where

$$E \stackrel{\text{def}}{=} \left\{ P(x, y) : \left| -\sum_{x, y} P(x, y) \log_2 Q(x) - H(X) \right| \le \varepsilon, \\ \left| -\sum_{x, y} P(x, y) \log_2 Q(y) - H(Y) \right| \le \varepsilon, \\ \left| -\sum_{x, y} P(x, y) \log_2 Q(x, y) - H(X, Y) \right| \le \varepsilon \right\}$$

57

Using Sanov's theorem,

$$Q_0^n(E) = 2^{-nD_{\text{KL}}(P^*||Q_0)}$$

where P^* is the closest distribution satisfying the constraints. Then $P^* \xrightarrow[\varepsilon \to 0]{} Q$ (exercise session). The probability becomes:

 $2^{-nD_{KL}(Q(x,y)||Q(x)Q(y))} = 2^{-nI(X,Y)}$

EXERCISE SESSION