

# Coding on Countably Infinite Alphabets

## Non-parametric Information Theory



## Lossless Coding on infinite alphabets

- Source Coding

- Universal Coding

- Infinite Alphabets

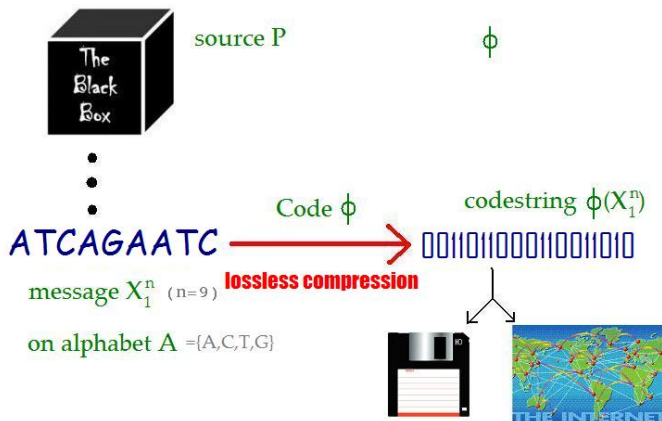
## Enveloppe Classes

- Theoretical Properties

- Algorithms



# Data Compression : Shannon Modelization



# Data Compression : Shannon Modelization



**Source P**

= stationary process on **Alphabet A** = {A,C,T,G}



ATCAGAATC

**Message**  $X_1^n$  ( $n=9$ )

**Code**  $\phi_n : A^n \rightarrow \{0,1\}^*$

**lossless compression**

Winzip, compress, etc.

0011011000110011010

**Codestring**  $\phi_n(X_1^n)$

Goal : minimize average  
codelength

$$E_P [|\phi(X_1^n)|]$$





## Lossless Coding on infinite alphabets

Source Coding

Universal Coding

Infinite Alphabets

## Enveloppe Classes

Theoretical Properties

Algorithms



- Shannon ('48) :

$$\mathbb{E}_P [|\phi_n(x)|] \geq H_n(P) = \mathbb{E}_P [-\log P^n(X_1^n)]$$

$n$ -block entropy and there is a code reaching the bound (within 1 bit).

- Moreover, if  $P$  is stationary and ergodic :

$$\frac{1}{n} H_n(P) \rightarrow H(P)$$

entropy rate of the source  $P$   
= minimal number of bits necessary per symbol.

- Kraft Inequality :

$$\sum_{x_1^n} 2^{-|\phi_n(x)|} \leq 1$$

and reciprocal.

- Code  $\phi \iff$  coding distribution  $Q_n^\phi(\mathbf{x}) = 2^{-|\phi_n(x)|}$ .
- The Shannon '48 theorem expresses that in average the best coding distribution is  $P$ !
- Otherwise, coding distribution  $Q_n$  suffers from the regret  
$$-\log Q_n(X_1^n) - (-\log P^n(X_1^n)) = \log \frac{P^n(X_1^n)}{Q_n(X_1^n)}.$$



## Lossless Coding on infinite alphabets

Source Coding

Universal Coding

Infinite Alphabets

## Enveloppe Classes

Theoretical Properties

Algorithms



- What if the source statistics are unknown ?
- What if we need versatile code ?

⇒ Need a **single coding distribution**  $Q_n$  for a **class of sources**

$$\Lambda = \{P_\theta, \theta \in \Theta\}$$

Ex : memoryless processes, Markov chains, HMM. . .

⇒ **unavoidable redundancy** :

$$\begin{aligned}\mathbb{E}_{P_\theta} [|\phi(X_1^n)|] - H(X_1^n) &= \mathbb{E}_{P_\theta} [\log Q_n(X_1^n) + \log P_\theta(X_1^n)] \\ &= KL(P_\theta, Q_n)\end{aligned}$$

= **Kullback-Leibler divergence** between  $P_\theta$  and  $Q_n$ .

## 1. Maximal regret :

$$R^*(Q_n, \Lambda) = \sup_{x_1^n \in A^n} \sup_{\theta \in \Theta} \log \frac{P_\theta(x_1^n)}{Q_n(x_1^n)}$$

## 2. Worst case redundancy :

$$R^+(Q_n, \Lambda) = \sup_{\theta \in \Theta} \mathbb{E}_{P_\theta} \left[ \log \frac{P_\theta^n(X_1^n)}{Q_n(X_1^n)} \right] = \sup_{\theta \in \Theta} KL(P_\theta, Q_n)$$

## 3. Expected redundancy with respect to prior $\pi$ :

$$R_\pi^-(Q_n, \Lambda) = \mathbb{E}_\pi \left[ \mathbb{E}_{P_\theta} \left[ \log \frac{P_\theta^n(X_1^n)}{Q_n(X_1^n)} \right] \right] = \mathbb{E}_\pi [KL(P_\theta, Q_n)]$$

$$\implies R^-(Q_n, \Lambda) \leq R^+(Q_n, \Lambda) \leq R^*(Q_n, \Lambda)$$

## 1. Minimax regret :

$$R_n^*(\Lambda) = \inf_{Q_n} R^*(Q_n, \Lambda) = \min_{Q_n} \max_{x_1^n, \theta} \log \frac{P_\theta^n(x_1^n)}{Q_n(x_1^n)}$$

## 2. Minimax redundancy :

$$R_n^+(\Lambda) = \inf_{Q_n} R^+(Q_n, \Lambda) = \min_{Q_n} \max_{\theta} KL(P_\theta^n, Q_n)$$

## 3. Maximin redundancy :

$$R_n^-(\Lambda) = \sup_{\pi} \inf_{Q_n} R_{\pi}^-(Q_n, \Lambda) = \max_{\pi} \min_{Q_n} \mathbb{E}_{\pi} [KL(P_{\theta}^n, Q_n)]$$

$$\implies R_n^-(\Lambda) \leq R_n^+(\Lambda) \leq R_n^*(\Lambda)$$

# General Standard Results

- Minimax result (Haussler '97, Sion) :  $R_n^-(\Lambda) = R_n^+(\Lambda)$ .
- Shtarkov's NML :  $R_n^*(\Lambda) = R^*(Q_{\text{NML}}^n, \Lambda)$  with  $Q_{\text{NML}}^n(\mathbf{x}) = \frac{\hat{p}(\mathbf{x})}{\sum_{\mathbf{y} \in \mathcal{X}^n} \hat{p}(\mathbf{y})}$  and  $\hat{p}(\mathbf{x}) = \sup_{P \in \Lambda} P^n(\mathbf{x})$ , so that and hence

$$R_n^*(\Lambda) = \log \left( \sum_{\mathbf{y} \in \mathcal{X}^n} \hat{p}(\mathbf{y}) \right).$$

- Channel capacity method : if  $d\mu(\theta, x) = d\pi(\theta)dP_\theta(x)$  then

$$\inf_{Q_n} R_{\pi}^-(Q_n, \Lambda) = H_{\mu}(X_1^n) - H_{\mu}(X_1^n | \theta) = H_{\mu}(\theta) - H_{\mu}(\theta | X_1^n).$$

# Finite Alphabets : Standard Results

- **Theorem (Shtarkov, Barron, Spankowski,...)** If  $\mathcal{I}_m$  denotes the class of memoryless processes over alphabet  $\{1, \dots, m\}$ , then

$$R_n^+(\mathcal{I}_m) = \frac{m-1}{2} \log \frac{n}{2e} + \log \frac{\Gamma(1/2)^m}{\Gamma(\frac{m}{2})} + o_m(1)$$

$$R_n^*(\mathcal{I}_m) = \frac{m-1}{2} \log \frac{n}{2} + \log \frac{\Gamma(1/2)^m}{\Gamma(\frac{m}{2})} + o_m(1) \leq \frac{m-1}{2} \log n + 2$$

- **Theorem (Rissanen '84)** If  $\dim \Theta = k$ , and if there exists a  $\sqrt{n}$ -consistent estimator of  $\theta$  given  $X_1^n$ , then

$$\liminf_{n \rightarrow \infty} R_n^-(\Lambda) \geq \frac{k}{2} \log n.$$



## Lossless Coding on infinite alphabets

Source Coding

Universal Coding

Infinite Alphabets

## Enveloppe Classes

Theoretical Properties

Algorithms



**Motivations** : integer coding, lossless image coding, text coding on words, mathematical curiosity.

1. Understand **general structural properties** of minimax redundancy and minimax regret.
2. Characterize those source classes that have **finite** minimax regret.
3. Quantitative **relations** between minimax redundancy or regret and integrability of the envelope function.
4. Develop **effective coding techniques** for source classes with known non-trivial minimax redundancy rate.
5. Develop **adaptive coding schemes** for collections of source classes that are too large to have redundancy rates.

# Sanity-check Properties

1.  $R^*(\Lambda^n) < +\infty \iff Q_{\text{NML}}^n$  is well-defined and given by

$$Q_{\text{NML}}^n(\mathbf{x}) = \frac{\hat{p}(\mathbf{x})}{\sum_{\mathbf{y} \in \mathcal{X}^n} \hat{p}(\mathbf{y})} \text{ for } \mathbf{x} \in \mathcal{X}^n$$

where  $\hat{p}(\mathbf{x}) = \sup_{P \in \Lambda} P^n(\mathbf{x})$ .

2.  $R^+(\Lambda^n)$  and  $R^*(\Lambda^n)$  are **non-decreasing**, **sub-additive** (or infinite) functions of  $n$ . Thus,

$$\lim_{n \rightarrow \infty} \frac{R^+(\Lambda^n)}{n} = \inf_{n \in \mathbb{N}_+} \frac{R^+(\Lambda^n)}{n} \leq R^+(\Lambda^1),$$

and

$$\lim_{n \rightarrow \infty} \frac{R^*(\Lambda^n)}{n} = \inf_{n \in \mathbb{N}_+} \frac{R^*(\Lambda^n)}{n} \leq R^*(\Lambda^1).$$



1. Let  $\Lambda$  be a class of stationary memoryless sources over a countably infinite alphabet. Let  $\hat{p}$  be defined by  $\hat{p}(x) = \sup_{P \in \Lambda} P\{x\}$ . The minimax regret with respect to  $\Lambda^n$  is finite if and only if the normalized maximum likelihood (Shtarkov) coding probability is well-defined and :

$$R^*(\Lambda^n) < \infty \Leftrightarrow \sum_{x \in \mathbb{N}_+} \hat{p}(x) < \infty.$$

2. We give an **example of a memoryless class  $\Lambda$**  such that  $R^+(\Lambda^n) < \infty$ , but  $R^*(\Lambda^n) = \infty$ .



## Lossless Coding on infinite alphabets

Source Coding

Universal Coding

Infinite Alphabets

## Enveloppe Classes

Theoretical Properties

Algorithms

## Definition

$$\Lambda_f = \left\{ P : \forall x \in \mathbb{N}, P^1\{x\} \leq f(x) \text{ and } P \text{ is stationary memoryless.} \right\}.$$

## Theorem

Let  $f$  be a non-negative function from  $\mathbb{N}_+$  to  $[0, 1]$ , let  $\Lambda_f$  be the class of stationary memoryless sources defined by envelope  $f$ . Then

$$R^+(\Lambda_f^n) < \infty \iff R^*(\Lambda_f^n) < \infty \iff \sum_{k \in \mathbb{N}_+} f(k) < \infty.$$

# Idea of the Proof

The only thing to prove is that

$$\sum_{k \in \mathbb{N}_+} f(k) = \infty \Rightarrow R^+(\Lambda_f^n) = \infty.$$

Let the sequence of integers  $(h_i)_{i \in \mathbb{N}}$  be defined recursively by  $h_0 = 0$  and

$$h_{i+1} = \min \left\{ h : \sum_{k=h_i+1}^h f(k) > 1 \right\}.$$

We consider the class  $\Lambda = \{P_\theta, \theta \in \Theta\}$  where  $\Theta = \mathbb{N}$ . The memoryless source  $P_i$  is defined by its first marginal  $P_i^1$  which is given by

$$P_i^1(m) = \frac{f(m)}{\sum_{k=h_i+1}^{h_{i+1}} f(k)} \text{ for } m \in \{p_i + 1, \dots, p_{i+1}\}.$$

Taking any **infinite entropy prior**  $\pi$  over  $\Theta$  the  $\Lambda^1 = \{P_i^1; i \in \mathbb{N}_+\}$  shows that

$$R^+(\Lambda^1) \geq R_\pi^-(\Lambda^1) = H(\theta) - H(\theta|X_1) = \infty.$$

## Theorem

If  $\Lambda$  is a class of memoryless sources, let the tail function  $\bar{F}_{\Lambda^1}$  be defined by  $\bar{F}_{\Lambda^1}(u) = \sum_{k>u} \hat{p}(k)$ , then :

$$R^*(\Lambda^n) \leq \inf_{u: u \leq n} \left[ n \bar{F}_{\Lambda^1}(u) \log e + \frac{u-1}{2} \log n \right] + 2.$$

## Corollary

Let  $\Lambda$  denote a class of memoryless sources, then the following holds :

$$R^*(\Lambda^n) < \infty \Leftrightarrow R^*(\Lambda^n) = o(n) \text{ and } R^+(\Lambda^n) = o(n).$$

## Theorem

Let  $f$  denote a non-increasing, summable envelope function. For any integer  $p$ , let  $c(p) = \sum_{k=1}^p f(2k)$ . Let  $c(\infty) = \sum_{k \geq 1} f(2k)$ . Assume furthermore that  $c(\infty) > 1$ . Let  $p \in \mathbb{N}_+$  be such that  $c(p) > 1$ . Let  $n \in \mathbb{N}_+$ ,  $\epsilon > 0$  and  $\lambda \in ]0, 1[$  be such that  $n > \frac{c(p)}{f(2p)} \frac{10}{\epsilon(1-\lambda)}$ . Then

$$R^+(\Lambda_f^n) \geq C(p, n, \lambda, \epsilon) \sum_{i=1}^p \left( \frac{1}{2} \log \frac{nf(2i)(1-\lambda)\pi}{2c(p)e} - \epsilon \right),$$

$$\text{where } C(p, n, \lambda, \epsilon) = \frac{1}{1 + \frac{c(p)}{\lambda^2 nf(2p)}} \left( 1 - \frac{4}{\pi} \sqrt{\frac{5c(p)}{(1-\lambda)\epsilon nf(2p)}} \right).$$

## Theorem

Let  $\alpha$  denote a real number larger than 1, and  $C$  be such that  $C\zeta(\alpha) \geq 2^\alpha$ . The source class  $\Lambda_{C, -\alpha}$  is the envelope class associated with the decreasing function  $f_{\alpha, C} : x \mapsto 1 \wedge \frac{C}{x^\alpha}$  for  $C > 1$  and  $\alpha > 1$ .

Then :

1.

$$n^{1/\alpha} A(\alpha) \log \left[ (C\zeta(\alpha))^{1/\alpha} \right] \leq R^+(\Lambda_{C, -\alpha}^n)$$

where

$$A(\alpha) = \frac{1}{\alpha} \int_1^\infty \frac{1}{u^{1-1/\alpha}} \left( 1 - e^{-1/(\zeta(\alpha)u)} \right) du.$$

2.

$$R^*(\Lambda_{C, -\alpha}^n) \leq \left( \frac{2Cn}{\alpha - 1} \right)^{1/\alpha} (\log n)^{1-1/\alpha} + O(1).$$

## Theorem

Let  $C$  and  $\alpha$  denote positive real numbers satisfying  $C > e^{2\alpha}$ . The class  $\Lambda_{Ce^{-\alpha \cdot}}$  is the envelope class associated with function  $f_\alpha : x \mapsto 1 \wedge Ce^{-\alpha x}$ . Then

$$\begin{aligned} \frac{1}{8\alpha} \log^2 n (1 - o(1)) &\leq R^+(\Lambda_{Ce^{-\alpha \cdot}}^n) \\ &\leq R^*(\Lambda_{Ce^{-\alpha \cdot}}^n) \leq \frac{1}{2\alpha} \log^2 n + O(1) \end{aligned}$$

cf Dominique Bontemps :  $R^+(\Lambda_{Ce^{-\alpha \cdot}}^n) \sim \frac{1}{4\alpha} \log^2 n$





## Lossless Coding on infinite alphabets

Source Coding

Universal Coding

Infinite Alphabets

## Enveloppe Classes

Theoretical Properties

Algorithms

# The CensoringCode Algorithm

## Algorithm CensoringCode

```
K ← 0
counts ← [1/2, 1/2, ...]
for i from 1 to n do
  cutoff ← ⌊ (4 * Ci / (α - 1))1/α ⌋
  if cutoff > K then
    for j ← K + 1 to cutoff do
      counts[0] ← counts[0] -
      counts[j] + 1/2
    end for
    K ← cutoff
  end if
  if x[i] ≤ cutoff then
    ArithCode(x[i], counts[0 : cutoff])
  else
    ArithCode(0, counts[0 : cutoff])
    C1 ← C1 · EliasCode(x[i])
    counts[0] ← counts[0] + 1
  end if
  counts[x[i]] ← counts[x[i]] + 1
end for
C2 ← ArithCode()
C1 · C2
```

## Theorem

Let  $C$  and  $\alpha$  be positive reals. Let the sequence of cutoffs  $(K_i)_{i \leq n}$  be given by

$$K_i = \left\lfloor \left( \frac{4C_i}{\alpha - 1} \right)^{1/\alpha} \right\rfloor.$$

Then

$$R^+(\text{CensoringCode}, \Lambda_{C, -\alpha}^n) \leq \left( \frac{4Cn}{\alpha - 1} \right)^{\frac{1}{\alpha}} \log n (1 + o(1)).$$

# Approximately Adaptive Algorithms

A sequence  $(Q^n)_n$  of coding probabilities is said to be *approximately asymptotically adaptive* with respect to a collection  $(\Lambda_m)_{m \in \mathcal{M}}$  of source classes if for each  $P \in \cup_{m \in \mathcal{M}} \Lambda_m$ , for each  $\Lambda_m$  such that  $P \in \Lambda_m$  :

$$D(P^n, Q^n) / R^+(\Lambda_m^n) \in O(\log n).$$

- A modification of **CensoringCode** choosing the  $n + 1$ th cutoff  $K_{n+1}$  according to the number of *distinct symbols* in  $\mathbf{x}$  is approximately adaptive on

$$\mathcal{W}_\alpha = \left\{ P : P \in \Lambda_{-\alpha}, \right. \\ \left. 0 < \liminf_k k^\alpha P^1(k) \leq \limsup_k k^\alpha P^1(k) < \infty \right\}.$$

- **Pattern coding** is approximately adaptive if  $1 < \alpha \leq 5/2$ .

The information conveyed in a message  $x$  can be separated into

1. a **dictionary**  $\Delta = \Delta(x)$  : the sequence of distinct symbols occurring in  $x$  in order of appearance ;
2. a **pattern**  $\psi = \psi(x)$  where  $\psi_i$  is the rank of  $x_i$  in dictionary  $\Delta$ .

Example :

Message	$x$	=	<i>a</i>	<i>b</i>	<i>r</i>	<i>a</i>	<i>c</i>	<i>a</i>	<i>d</i>	<i>a</i>	<i>b</i>	<i>r</i>	<i>a</i>
Pattern	$\psi(x)$	=	1	2	3	1	4	1	5	1	2	3	1
Dictionary	$\Delta(x)$	=	<i>a</i>	<i>b</i>	<i>r</i>		<i>c</i>		<i>d</i>				

$\implies$  A random process  $(X_n)_n$  with distribution  $P$  induces a random pattern process  $(\Psi_n)_n$  on  $\mathbb{N}_+$  with distribution :

$$P^\Psi (\Psi_1^n = \psi_1^n) = \sum_{x_1^n: \psi(x_1^n) = \psi_1^n} P(X_1^n = x_1^n).$$

- The pattern entropy rate exists and coincides with the process entropy rate :

$$\frac{1}{n} H(\Psi_1^n) = \frac{1}{n} \mathbb{E}_{P^\Psi} \left[ -\log P^\Psi(\Psi_1^n) \right] \rightarrow H(\Psi) = H(X).$$

- Pattern redundancy satisfies :

$$1.84 \left( \frac{n}{\log n} \right)^{\frac{1}{3}} \leq R_{\Psi,n}^-(\mathcal{I}_\infty) \leq R_{\Psi,n}^*(\mathcal{I}_\infty) \leq \left( \pi \sqrt{\frac{2}{3}} \log e \right) \sqrt{n}.$$

The proof of the lower-bound uses fine combinatorics on integer partitions with small summands (see Garivier '09). Upper-bounds in  $O(n^{2/5})$  have been given (see Shamir '07), but there is still a gap between lower- and upper-bounds.

- For memoryless coding, the pattern redundancy is neglectible wrt. the codelength of the dictionary whenever  $n < 5/2$ .