

# Infinite urn models and applications

Collegio Carlo Alberto Seminar-Torino

Stéphane Boucheron

January 27th 2017

# Infinite urn models, missing mass and applications

## Infinite urn models

# Setting

## Probability distribution over $\mathbb{N}$

1. Sequence  $U_1, \dots, U_n, \dots$ , i.i.d. over  $\mathbb{N}_+ := \mathbb{N} \setminus \{0\}$  (the *symbols*)
2. (Unknown) distribution defined by p.m.f. :  $(p_j)_{j=1}^{\infty}$
3. Counts:  $X_{n,j} = \sum_{i=1}^n \mathbb{I}_{U_i=j}$ , number of occurrences of  $j$  among the first  $n$  symbols

## Inference problems (species sampling)

1. Estimation of  $(p_j)_{j=1}^{\infty}$
2. Estimation of missing mass  $M_{n,0} = \sum_{j=1}^{\infty} p_j \mathbb{I}_{X_{n,j}=0}$
3. Estimation of mass of rare symbols  $M_{n,r} = \sum_{j=1}^{\infty} p_j \mathbb{I}_{X_{n,j}=r}$
4. Estimation of discovery probability  $M_{n+m,r}$  from  $U_1, \dots, U_n$

## Occupancy counts, profile

 $K_{n,r}$ 

For  $r \geq 1$ ,  $K_{n,r}$  : number of symbols with  $r$  occurrences in sample

 $K_{n,\bar{r}}$ 

For  $r \geq 1$ ,  $K_{n,\bar{r}}$  : number of symbols with at least  $r$  occurrences in sample

 $K_n$ 

$K_n$  : number of distinct symbols in the sample

### Profile, histogram

 $(K_{n,r})_{r \geq 1}$ 

$(K_{n,r}/K_n)_r$  defines an (empirical) distribution over integers

$$K_n = \sum_{r \geq 1} K_{n,r} \quad n = \sum_{r \geq 1} K_{n,r} \times r \quad \text{obvious!}$$

# Motivations for mass estimators

## Connections between masses and occupancy counts

$$\mathbb{E}M_{n,0} = \mathbb{E} \frac{K_{n+1,1}}{n+1} \quad \Rightarrow \quad \mathbb{E} \frac{K_{n,1}}{n} = \mathbb{E}M_{n-1,0} \geq \mathbb{E}M_{n,0}$$

$$\mathbb{E}M_{n,0} \leq \mathbb{E} \frac{K_{n,1}}{n} = \mathbb{E} \frac{K_{n+1,1}}{n+1} + \frac{2}{n} \frac{\mathbb{E}K_{n+1,2}}{n+1} = \mathbb{E}M_{n,0} + \frac{1}{n} \mathbb{E}M_{n,1}$$

$$\mathbb{E}M_{n,0} \leq \mathbb{E} \frac{K_{n,1}}{n} \leq \mathbb{E}M_{n,0} + \frac{1}{n}$$

## Mass estimation

### Good-Turing estimator of the missing mass

$$\hat{M}_{n,0} := \frac{K_{n,1}}{n}$$

### Good-Turing mass-estimators

$$\hat{M}_{n,r} := (r + 1) \frac{K_{n,r+1}}{n}$$

### Harder problems

- ▶ Discovery probabilities  $\mathbb{E}(M_{n+m,r} \mid \mathcal{F}_n)$
- ▶ Number of new discoveries  $\mathbb{E}(K_{n+m} \mid \mathcal{F}_n) - K_n$

# Assessing Good-Turing estimators

Providing non-asymptotic/distribution-free guarantees for GT estimators

How large is  $\left| \frac{\widehat{M}_{0,n}}{M_{0,n}} - 1 \right|$  ?

- ▶ Needs to get a handle on the tail behavior of  $\frac{K_{n,1}}{n}$  and of  $M_{n,0}$



## Concentration inequalities

## Classical exponential inequalities

$X_1, \dots, X_n$  independent centered random variables

$$Z := \sum_{i=1}^n X_i$$

Hoeffding  $a_i \leq X_i \leq b_i$

$$\mathbb{P}\{Z \geq t\} \leq e^{-\frac{4t^2}{2 \sum_{i=1}^n (b_i - a_i)^2}} \quad \text{for } t > 0$$

Bennett  $X_i \leq b$   $v := \text{var}(Z)$

$$\mathbb{P}\{Z \geq t\} \leq e^{-\frac{v}{b^2} h\left(\frac{bt}{v}\right)} \quad \text{for } t > 0 \text{ with } h(x) = (1+x) \log(1+x) - x$$

Bernstein  $\sum_{i=1}^n \mathbb{E}(X_i)_+^q \leq \frac{q!}{2} v c^{q-2}$   $v \geq \text{var}(Z)$

$$\mathbb{P}\{Z \geq t\} \leq e^{-\frac{t^2}{2(v+ct)}} \quad \text{for } t > 0$$

## Concentration of measure phenomenon

Any function of many independent random variables that does not depend too much on any of them is strongly concentrated around its expectation (or its median)

M. Talagrand

### Gaussian concentration (1975)

$X \sim \mathcal{N}(0, \text{Id}_n)$

$f : \mathbb{R}^n \rightarrow \mathbb{R}$ ,  $L$ -Lipschitz,  $Z := f(X_1, \dots, X_n)$

$$\text{var}(Z) \leq \mathbb{E} [\|\nabla f\|^2] \leq L^2$$

$$\mathbb{P} \{Z \geq \mathbb{E}Z + t\} \leq e^{-\frac{t^2}{2L^2}}$$

## Poisson/binomial scenarii

### Poissonizing sample/message length

- ▶  $(X_j)_{j \geq 1}$  : a collection of independent Poisson Point Processes on  $(0, \infty)$  with intensity  $p_j$  ( $X_j(t) \sim \text{Poi}_{tp_j}$ )
- ▶  $K_r(t)$  : sum of independent Bernoulli random variables  $\mathbb{I}_{X_j(t)=r} \sim \text{Bernoulli}$  with success probability  $\text{Poi}_{tp_j}(r)$
- ▶  $M_r(t)$  : sum of weighted independent Bernoulli random variables  

$$M_r(t) = \sum_{j \geq 1} p_j \mathbb{I}_{X_j(t)=r}$$

### Concentration inequalities come (almost) for free

$K_r(n)$  are sums of independent Bernoulli random variables

$M_r(n)$  are weighted sums of independent Bernoulli random variables

### Poissonized identities for the missing mass

$$\mathbb{E}M_0(t) = \frac{\mathbb{E}K_1(t)}{t} \quad \text{and} \quad \text{var}(M_0(t)) = 2 \frac{\mathbb{E}K_2(t)}{t^2} - \frac{\mathbb{E}K_2(2t)}{2t^2}$$

## Concentration inequalities

## Variance inequalities

For the number of distinct symbols

$$\frac{\mathbb{E}K_1(2t)}{2t} \leq \text{var}(K(t)) \leq \mathbb{E}K_1(t) \quad (\text{Poisson setting})$$

$$\text{var}(K_n) \leq \mathbb{E}[(1 - M_{n,0})K_{n,1}] \leq \mathbb{E}K_{n,1} \leq \mathbb{E}K_n \quad (\text{Binomial setting})$$

Tight in the birthday paradox setting

$$p_j = 1/n^2 \text{ for } 1 \leq j \leq n^2, \quad 1 - M_{n,0} = K_n/n^2$$

$$\text{var}(K_n) \leq \mathbb{E} \left[ K_{n,1} \frac{K_n}{n^2} \right] \leq 1$$

## Efron-Stein-Steele inequalities (1981-86)

Thm. : the expected Jackknife estimate of variance is an upper bound on variance

- ▶  $Z = f(X_1, \dots, X_n)$  where  $X_1, \dots, X_n$  are independent.
- ▶  $Z_i = f_i(X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n)$  ( $f_i$  any function)

$$\text{var}(Z) \leq \mathbb{E} \left[ \sum_{i=1}^n (Z - Z_i)^2 \right]$$

## Efron-Stein-Steele inequalities : sketch of proof

$X_1, \dots, X_n$  independent random variables

$X'_1, \dots, X'_n$  independent "copies" ( $X'_i \sim X_i$ )

$f, g : \prod_{i=1}^n \mathcal{X}_i \rightarrow \mathbb{R}$

$X^{(i)} := X_1, \dots, X_{i-1}, X'_i, X_{i+1}, \dots, X_n$

$X^{(0)} := X_1, \dots, X_n$        $X := X^{(0)}, X' := X^{(n)}$

$$\begin{aligned} \text{cov}(f(X), g(X)) &= \mathbb{E} [f(X)(g(X) - g(X'))] \\ &= \sum_{i=1}^n \mathbb{E} [f(X) (g(X^{(i-1)}) - g(X^{(i)}))] \\ &= \frac{1}{2} \sum_{i=1}^n \mathbb{E} [(f(X) - f(X^{(i)})) (g(X^{(i-1)}) - g(X^{(i)}))] \end{aligned}$$

$$\text{cov}(f(X), f(X)) \leq \frac{1}{2} \sum_{i=1}^n \mathbb{E} [(f(X) - f(X^{(i)}))^2] \quad (\text{Cauchy-Schwarz})$$



## Variance inequalities for occupancy counts and missing mass

$$\text{var}(K_{n,\bar{r}}) \leq r\mathbb{E}K_{n,r}$$

$$\text{var}(K_{n,r}) \leq r\mathbb{E}K_{n,r} + (r+1)\mathbb{E}K_{n,r+1}$$

$$\text{var}(M_{n,0}) \leq \sum_{j=1}^{\infty} p_j^2 \text{var}(\mathbb{I}_{X_{n,j}=0}) \leq 2 \frac{\mathbb{E}K_2(n)}{n^2}$$

## Negative association

### Definition

Real-valued random variables  $Z_1, \dots, Z_K$  are said to be *negatively associated* if, for any two disjoint subsets  $A$  and  $B$  of  $\{1, \dots, K\}$ , and any two real-valued functions  $f : \mathbb{R}^{|A|} \mapsto \mathbb{R}$  and  $g : \mathbb{R}^{|B|} \mapsto \mathbb{R}$  that are both either coordinate-wise non-increasing or coordinate-wise non-decreasing :

$$\mathbb{E} [f(Z_A).g(Z_B)] \leq \mathbb{E} [f(Z_A)] . \mathbb{E} [g(Z_B)] .$$

### Examples

In the Bernoulli scenario  $(X_{n,j})_{j \geq 1}$  are negatively associated

Monotone functions of negatively associated random variables are negatively associated

# Concentration inequalities

$K_n$  is sub-Poisson

$$\log \mathbb{E} e^{\lambda(K_n - \mathbb{E}K_n)} \leq \mathbb{E}K_{n,1} \phi(\lambda) \quad \text{with} \quad \phi(\lambda) = e^\lambda - \lambda - 1$$

$K_n$  satisfies a Bennett inequality with variance factor  $\mathbb{E}K_{n,1}$

$K_{n,\bar{r}}$

Is a sum of negatively associated Bernoulli random variables

$K_{n,\bar{r}}$  satisfies a Bennett inequality with variance factor  $r\mathbb{E}K_{n,r}$

$K_{n,r}$

$$K_{n,r} = K_{n,\bar{r}} - K_{n,\overline{r+1}}$$

# Missing mass

# Generic concentration inequalities for missing mass

Lower tail (sub-Gaussian with sharp bound on variance, McAllester-Ortiz)

$$v_n^- := \frac{\mathbb{E}K_2(n)}{n^2}$$

$$\log \mathbb{E}e^{\lambda(M_{n,0} - \mathbb{E}M_{n,0})} \leq v_n^- \frac{\lambda^2}{2} \quad \lambda \leq 0$$

Upper tail (sub-Gamma with looser bound on variance)

$$v_n^+ := \frac{\mathbb{E}K_2(n)}{n^2}$$

$$\log \mathbb{E}e^{\lambda(M_{n,0} - \mathbb{E}M_{n,0})} \leq v_n^+ \frac{\lambda^2}{2(1 - \lambda/n)} \quad \lambda \geq 0$$

## Proof

$$\phi(x) := \exp(x) - x - 1$$

$$\log \mathbb{E} \left[ e^{\lambda(M_{n,0} - \mathbb{E}M_{n,0})} \right] \leq \sum_{j=1}^{\infty} e^{-np_j} \phi(\lambda p_j) \quad \text{neg. assoc and Bennett ineq.}$$

$$\begin{aligned} \sum_{j=1}^{\infty} e^{-np_j} \phi(\lambda p_j) &= \sum_{r=2}^{\infty} \left( \frac{\lambda}{n} \right)^r \sum_{j=1}^{\infty} e^{-np_j} \frac{(np_j)^r}{r!} \\ &= \sum_{r=2}^{\infty} \left( \frac{\lambda}{n} \right)^r \mathbb{E}K_r(n) \\ &\leq \mathbb{E}K_{\frac{1}{2}}(n) \sum_{r=2}^{\infty} \left( \frac{\lambda}{n} \right)^r \\ &= \lambda^2 \frac{\mathbb{E}K_{\frac{1}{2}}(n)/n^2}{1 - \lambda/n}, \end{aligned}$$

## Sub-Gaussian versus Sub-Gamma ?

### Uniformity

A collection  $\mathcal{C}$  of random variables is *uniformly* sub-Gaussian if there exists a constant  $c$ , such that satisfies

$$\forall X \in \mathcal{C}, \quad \log \mathbb{E} e^{\lambda(X - \mathbb{E}X)} \leq c \operatorname{var}(X) \frac{\lambda^2}{2}$$

### Bernoulli random variables are *not* uniformly sub-Gaussian

If  $X_p \sim \operatorname{Ber}(p)$ , the best possible upper-bound for all  $p \in (0, 1)$  is

$$\log e^{\lambda(X_p - \mathbb{E}X_p)} \leq \operatorname{var}(X_p) \frac{\lambda^2}{2(2p(1-p)c_{\text{LS}}(p))} \quad \text{where} \quad c_{\text{LS}}(p) = \frac{\log((1-p)/p)}{1-2p}$$

## Sub-Gaussian bounds for the right-tail of $M_{n,0}$

Best variance proxy

$$w_n := \sum_{j=1}^{\infty} \frac{p_j^2}{2c_{\text{LS}}((1-p_j)^n)} \leq \frac{1}{2n}$$

If  $\max_j p_j \leq 1/4$

$$w_n \geq \frac{3}{32n} \left( 1 - \sum_{j:p_j < 1/n} p_j \right)$$

$w_n$  may be order of magnitudes larger than  $\mathbb{E}K_2(n)/n^2$



## Karlin's setting

# Motivations

Predicting the number of discoveries during next campaign

Ideally

$K_{\lceil nx \rceil} / K_n$  should converge in probability toward a function of  $x$

At least :  $\lim_n \mathbb{E}K_{\lceil nx \rceil} / \mathbb{E}K_n \in (0, \infty)$

Estimating the missing mass

# Counting functions and regular variation

Counting function  $\vec{\nu}$  and counting measure  $\nu$

$$\vec{\nu}(x) := \# \{j : p_j \geq x\} \quad \nu(A) = \sum_{j \geq 1} \mathbb{I}_{p_j \in A}$$

$$\nu_1(0, x) := \sum_{j=1}^{\infty} p_j \mathbb{I}_{p_j \leq x}$$

Karlin's assumption

Regular variation assumption on  $\vec{\nu}$  in the neighborhood of 0

$$\lim_{x \searrow 0} \frac{\vec{\nu}(sx)}{\vec{\nu}(x)} = s^{-\alpha} \quad \text{for some } \alpha \in (0, 1)$$

## Why?

For  $\alpha \in (0, 1)$  Karlin's assumption implies (among other things)

- ▶  $\mathbb{E}K_n \sim \Gamma(1 - \alpha)\bar{\nu}(1/n)$
- ▶  $\mathbb{E}K_{n,r} \sim \frac{\alpha\Gamma(r-\alpha)}{r!}\bar{\nu}(1/n)$
- ▶  $\text{var}(M_{n,0}) \sim \alpha\Gamma(2 - \alpha)(1 - 2^{\alpha-2}\frac{\bar{\nu}(1/n)}{n^2})$

Chinese restaurant processes (Poisson-Dirichlet processes)

A Poisson-Dirichlet process with parameters  $\alpha \in (0, 1)$  and  $\theta > -\alpha$  generates a distribution over  $\mathbb{N} \setminus \{0\}$  that is a.s.  $\alpha$ -regularly varying

## Challenging concentration inequalities for $M_{n,0}$ : $\alpha \in (0, 1)$

If  $\vec{\nu}$  satisfies the regular variation condition with index  $\alpha \in (0, 1)$ , then the missing mass  $M_{n,0}$  (or  $M_0(n)$ ) is sub-Gaussian on the left tail with variance factor  $v_n^- = 2\mathbb{E}K_2(n)/n^2$  and sub-gamma on the right tail with variance factor  $v_n^+ = 2\mathbb{E}K_{\frac{1}{2}}(n)/n^2$ .

The variance factors satisfy

$$\lim_n \frac{v_n^-}{\text{var}(M_{n,0})} = \frac{1}{1 - 2^{\alpha-2}},$$

$$\lim_n \frac{v_n^+}{\text{var}(M_{n,0})} = \frac{2}{\alpha(1 - 2^{\alpha-2})},$$

and thus

$$\lim_n \frac{v_n^-}{v_n^+} = \frac{\alpha}{2}.$$

## Challenging concentration inequalities for $M_{n,0}$ : $\alpha = 0$

Two settings :  $\mathbb{E}K_{n,1} \rightarrow \infty$  or  $\mathbb{E}K_{n,1} = O(1)$

- ▶  $\mathbb{E}K_{n,1} = O(1)$  for Geometric, negative binomial, Poisson, ... distributions
- ▶  $\nu_1(0, x) \sim_{x \rightarrow 0} x l_1(1/x)$  with  $l_1 \in RV_0$  and  $l_1(t) \nearrow \infty$

de Haan's regular variation (Extended Regular Variation)

$\bar{\nu}(1/\cdot) \in \Pi_\alpha$  where  $\alpha$  is slowly varying if

$$\lim_{t \rightarrow \infty} \frac{\bar{\nu}(1/(tx)) - \bar{\nu}(1/t)}{\alpha(t)} = \log(x) \quad \text{for } x > 0$$

If  $\vec{\nu}(1/\cdot) \in \Pi_\alpha$  with  $\alpha \in RV_0$  and  $\alpha(t) \nearrow \infty$

For each  $r \geq n$

- ▶  $K_n \stackrel{\mathbb{P}}{\sim} \mathbb{E}K_n \underset{+\infty}{\sim} \ell(n)$
- ▶  $K_{n,r} \stackrel{\mathbb{P}}{\sim} \mathbb{E}K_{n,r} \underset{+\infty}{\sim} \frac{\alpha(n)}{r}$
- ▶  $K_{n,\bar{r}} \stackrel{\mathbb{P}}{\sim} \mathbb{E}K_{n,\bar{r}} \underset{+\infty}{\sim} \ell(n)$
- ▶  $M_{n,r} \stackrel{\mathbb{P}}{\sim} \mathbb{E}M_{n,r} \underset{+\infty}{\sim} \frac{\alpha(n)}{n}$ .

Missing masses are uniformly sub-gamma

- ▶  $v(n) := \frac{12\alpha(n)}{n^2}$
- ▶  $\text{var}(M_{n,0}) \sim \frac{3\alpha(n)}{4n^2}$
- ▶  $\exists n_0, \forall n > n_0,$

$$\log \mathbb{E} \left[ e^{\lambda(M_{n,0} - \mathbb{E}M_{n,0})} \right] \leq \frac{v_n \lambda^2}{2(1 - \lambda/n)} \quad \lambda > 0$$

## Excessively slow variation

Geometric distribution ( $p_j = q(1 - q)^{j-1}$ )

No hope to obtain uniform sub-gamma tail behavior

What is the missing mass made of?

- ▶ Essentially symbols with probability of order  $1/n$
- ▶ The *concentration on two points* phenomenon: with probability tending to 1,  $M_{0,n}$  tends to be distributed over a pair of values (that depend on  $n$ )



# Estimation

## Estimating the missing mass

One needs to make assumptions on the sampling distribution to guarantee the consistency of the Good-Turing estimator.

There is no hope to find a universally consistent estimator of the missing mass without any such restrictions, as shown recently by Mossel & Ohannessian (2015)

If  $\mathbb{E}K_{n,2}/\mathbb{E}K_{n,1}$  remains bounded and  $\mathbb{E}K_{n,1} \rightarrow +\infty$

$$\frac{M_{n,0}}{\mathbb{E}M_{n,0}} \xrightarrow{\mathbb{P}} 1,$$

and

$$\frac{\widehat{M}_{n,0}}{M_{n,0}} \xrightarrow{\mathbb{P}} 1.$$

## Poissonized setting

$$\text{var}(\widehat{M}_0(t) - M_0(t)) = \frac{1}{t^2} (\mathbb{E}K_1(t) + 2\mathbb{E}K_2(t)) . \quad (1)$$

$\widehat{M}_0(t) - M_0(t)$  is sub-gamma

1. on the right tail with variance factor  $\text{var}(\widehat{M}_0(t) - M_0(t))$  and scale factor  $1/t$ ,
2. on the left tail with variance factor  $3\mathbb{E}K(t)/t^2$  and scale factor  $1/t$ .

For all  $\lambda \geq 0$ ,

1.  $\log \mathbb{E}e^{\lambda(\widehat{M}_0(t) - M_0(t))} \leq \text{var}(\widehat{M}_0(t) - M_0(t))t^2 \phi\left(\frac{\lambda}{t}\right)$
2.  $\log \mathbb{E}e^{\lambda(M_0(t) - \widehat{M}_0(t))} \leq \frac{3\mathbb{E}K(t)}{2t^2} \frac{\lambda^2}{1 - \lambda/t}$

## Lossless Compression over countable alphabets

# Statistical coding

## Redundancy of coding probability $Q^n$ with respect to source $P^n$

Expected difference between codelengths obtained by feeding an arithmetic coder with  $Q^n(\mathbf{x})$  rather than with the correct source statistics  $P^n(\mathbf{x})$

$$D(P^n, Q^n) = \mathbb{E}_{P^n} \log \frac{P^n(X_{1:n})}{Q^n(X_{1:n})}$$

$\Lambda^n$ : collection of probability distributions over messages of length  $n$ . Each probability distribution is called a **source**.

### Minimax redundancy

$$\bar{R}(\Lambda^n) = \inf_Q \sup_{P \in \Lambda} D(P^n, Q^n)$$

MinMax Theorem

### Maximin redundancy

$\pi$ : prior distribution on sources

$$\underline{R}(\Lambda^n) = \sup_{\pi} \inf_Q \mathbb{E}_{\pi} D(P^n, Q^n)$$

$$\underline{R}(\Lambda^n) = \bar{R}(\Lambda^n)$$

# Envelope classes

## Envelop function

$f: \mathbb{N} \rightarrow \mathbb{R}_+$  with  $1 < \sum_{j>0} f(j) < \infty$ .

## Envelope class

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

## Envelope distribution

1.  $F(k) := 1 - \sum_{j>k} f(j)$  for  $k \geq l_f := \max\{k: \sum_{j \geq k} f(j) \geq 1\}$  envelope distribution
2.  $\bar{F} = 1 - F$  tail envelope function
3.  $\vec{v}(x) = |\{j: f(j) \geq x\}|$  for  $x \in ]0, \infty[$  counting function

## PCC Code (I)

{Mixture coding of the censored sequence}

$$\tilde{Q}^{n+1}(\tilde{X}_{1:n}) = \prod_{i=0}^{n-1} \tilde{Q}_{i+1}(\tilde{X}_{i+1} \mid \tilde{X}_{1:i}),$$

Predictive probabilities  $\tilde{Q}_{i+1}(\cdot \mid \tilde{X}_{1:i})$

given by **Krichevsky-Trofimov** mixtures on the alphabet  $\{0, 1, \dots, K_i\}$ ,

$$\tilde{Q}_{i+1}(\tilde{X}_{i+1} = k \mid \tilde{X}_{1:i} = \tilde{X}_{1:i}) = \frac{\tilde{n}_i^k + \frac{1}{2}}{j + \frac{K_i+1}{2}},$$

where, for  $0 \leq k \leq K_i$ ,  $\tilde{n}_i^k$  is the number of occurrences of symbol  $k$  in  $\tilde{X}_{1:i}$ .

For  $0 \leq k \leq K_i$ ,

$$\tilde{n}_i^k = \begin{cases} K_i & \text{if } k = 0, \\ n_i^j - 1 & \text{if } 1 \leq k \leq K_i \text{ and } \langle j, k \rangle \in \mathcal{D}_i, \end{cases}$$

where  $n_i^j$  is the number of occurrences of symbol  $j$  in  $x_{1:i}$ .

## PCC Code (II)

1. **Dictionary** initialization:  $\mathcal{D}_0 \leftarrow \{\langle 0, 0 \rangle\}$   
For  $i > 0$ ,  $\mathcal{D}_i$  contains 0 and symbols from  $x_{1:i}$  with rank of insertion
2. Create a **censored sequence**  $\tilde{x}_{1:n}$

$$\tilde{x}_i = \begin{cases} k & \text{if } \exists k \leq K_{i-1}, \langle x_i, k \rangle \in \mathcal{D}_{i-1} \\ 0 & \text{otherwise, (newly discovered symbol)} \end{cases}$$

where  $K_i :=$  number of distinct symbols in  $x_{1:i}$ .

3. Dictionary update: when  $\tilde{x}_i = 0$ ,  $\mathcal{D}_i \leftarrow \mathcal{D}_{i-1} \cup \langle x_i, K_i \rangle$ , with  $K_i = K_{i-1} + 1$ .
4. Let  $K :=$  number of *censored* input symbols and let  $i_{1:K}$  be their indices.  
**Extract subsequence**  $x_{i_{1:K}}$
5.  $C_M :=$  **Mixture/arithmetic** coding of  $\tilde{x}_{1:n}$
6.  $C_E :=$  **Elias** / integer coding of each redacted symbol in  $x_{i_{1:K}}$  individually
7. **Interleave** the coded censored symbols of  $C_E$  just after each coded 0 symbol of  $C_M$ , to form the overall PCC-code.



## Weak adaptivity of PCC codes

$$\mathbf{R}_f(n) := \log(e) \int_1^n \frac{\vec{\nu}_f(1/t)}{2t} dt.$$

Let  $\Lambda_f$  be an envelope source class, with  $f \in \text{RV}_\alpha$  and  $\alpha \in (0, 1[$ . Then

$$\bar{R}(\Lambda_f^n) \asymp \mathbf{R}_f(n).$$

### Theorem

Let  $(Q_n)$  : coding distribution associated to the Pattern Censoring Code.  $\forall \alpha \in (0, 1[$ ,  $\forall f$  with  $\vec{\nu}_f(1/\cdot) \in \text{RV}_\alpha$ ,  $\exists a_f, b_f > 0$  :

$$(a_f + o_f(1))\mathbf{R}_f(n) \leq \bar{R}(\Lambda_f^n) \leq \bar{R}(Q_n, \Lambda_f^n) \leq (b_f + o_f(1))\mathbf{R}_f(n) \log \log n.$$

{The Pattern Censoring Code is adaptive, within  $\log \log n$ , w.r.t.  $(\Lambda_f)_{f \in \text{RV}_\alpha} \}_{\alpha \in (0, 1)}$ }

## Redundancy of PCC code

Redundancy of the PCC code with respect to a source  $P$  over  $\mathbb{N} \setminus \{0\}$

$$D(P^n | \mathcal{Q}^n) = \sum_{i=1}^n \mathbb{E} \left[ \mathbb{E}_P \left[ \log \frac{P(X_{i+1})}{Q(X_{i+1}|X_{1:i})} \middle| X_{1:i} \right] \right]$$

Decomposing instantaneous redundancy

$$\begin{aligned} & \mathbb{E}_P \left[ \log \frac{P(X_{i+1})}{Q(X_{i+1}|X_{1:i})} \middle| X_{1:i} \right] \\ &= \underbrace{\sum_{j \geq 1} p_j \mathbb{I}_{N_j^i > 0} \log \left( \frac{p_j \left( i + \frac{K_j+1}{2} \right)}{N_j^i - \frac{1}{2}} \right)}_{(i)} \\ & \quad + \underbrace{\sum_{j \geq 1} p_j \mathbb{I}_{N_j^i = 0} \log \left( \frac{p_j \left( i + \frac{K_j+1}{2} \right)}{K_j + \frac{1}{2}} \right)}_{(ii)} \\ & \quad + \underbrace{\sum_{j \geq 1} p_j \mathbb{I}_{N_j^i = 0} \left( \log(j+1) + 2 \log \log(j+1) \right)}_{(iii)}. \end{aligned}$$

## Redundancy of PCC code}

### Distribution-free bounds

$$\begin{aligned} & \mathbb{E}_P \left[ \log \frac{P(X_{i+1})}{Q(X_{i+1}|X_{1:i})} \right] \\ & \leq \kappa \left( \nu_1(0, 1/i) + \frac{\bar{\nu}(1/i)}{i} \right) \\ & \quad + \sum_{j \geq 1} p_j \mathbb{P}\{N_i^j = 0\} \left( \log \left( \frac{j+1}{\mathbb{E}K_j} \right) + 2 \log \log(j+1) \right), \end{aligned}$$

where  $\kappa$  is (another) universal constant

### Using regular variation of the envelope distribution

For each Regularly varying envelope,  $\exists i_0 \in \mathbb{N}$  such that, for all  $i \geq i_0$

$$\begin{aligned} & \sum_{j \geq 1} p_j \mathbb{I}_{N_i^j=0} (\log \log(j+1)) \\ & \leq \frac{5}{2} \left( \Gamma(1-\alpha) + \frac{1}{1-\alpha} \right) \frac{\log \log(i) \bar{\nu}_f(1/i)}{i} \end{aligned}$$

## Techniques

Everything boils down to controlling inverse binomials, inverse of number of distinct symbols  $K_i$ , missing mass  $M_{i,0}$ , and some random integrals

The distribution and moments of these stochastic quantities ultimately depend on details of the sampling distribution, but ...

Some non-asymptotic inequalities are distribution-free

- ▷  $\mathbb{E}M_{i,0} \leq \mathbb{E}\frac{K_{i,1}}{i}$
- ▷  $\text{var}(K_i) \leq \mathbb{E}K_{i,1}$
- ▷  $\mathbb{E}\frac{1}{K_i} \leq \frac{3}{\mathbb{E}K_i}$
- ▷ ...

The analysis is made modular by delaying invocation of regular variation arguments

## References

- ▶ A. Ben-Hamou, B., M. Ohannessian: *Concentration inequalities in the infinite urn scheme with occupancy counts and the missing mass, with applications*. Bernoulli. **23**, 249-287 (2017)
- ▶ A. Ben-Hamou, B., E. Gassiat: *Censoring meets pattern-coding*. ArXiv.