

# Context tree models for source coding

---

## Toward Non-parametric Information Theory



## Lossless Source Coding = density estimation with log-loss

Source Coding and Universal Coding

Bayesian coding

## Context-tree Models and Double Mixture

Rissanen's model

Coding in Context Tree models

The Context-Tree Weighting Method

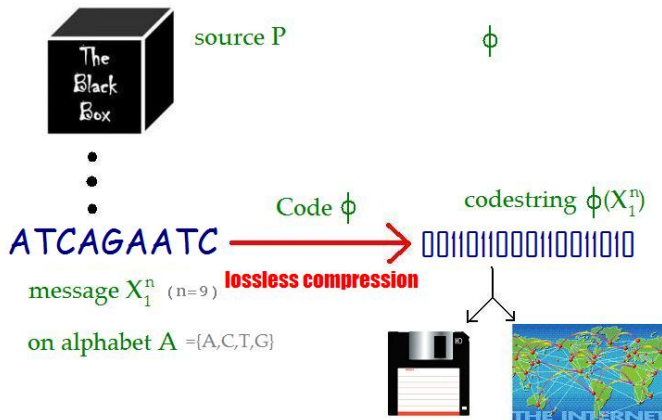
## Context Tree Weighting on Renewal Processes

CTW on Renewal Processes

Perspectives



# 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 Source Coding = density estimation with log-loss

Source Coding and Universal Coding

Bayesian coding

## Context-tree Models and Double Mixture

Rissanen's model

Coding in Context Tree models

The Context-Tree Weighting Method

## Context Tree Weighting on Renewal Processes

CTW on Renewal Processes

Perspectives



# The Problem

- Finite *alphabet*  $A$ ; ex :  $A = \{a, c, t, g\}$
- Set of all finite binary sequences :  $\{0, 1\}^* = \bigcup_{n \geq 0} \{0, 1\}^n$
- If  $x \in A^n$ , its *length* is  $|x| = n$
- $\log =$  base 2 logarithm
- Let  $X$  be a stochastic process on the finite alphabet  $A$ , with stationary ergodic distribution  $\mathbb{P}$
- The marginal distribution of  $X_1^n$  is denoted by  $P^n$
- For  $n \in \mathbb{N}^*$ , the *coding function*  $\phi_n : A^n \rightarrow \{0, 1\}^*$
- Expected code length :

$$\mathbb{E}_{\mathbb{P}} [|\phi_n(X_1^n)|]$$

# Example : DNA Codes $A = \{a, c, t, g\}$

- Uniform concatenation code :

$$\phi_1 = \begin{array}{|c|c|c|c|} \hline a & c & t & g \\ \hline 00 & 01 & 10 & 11 \\ \hline \end{array}, \phi_n(X_1^n) = \phi_1(X_1) \dots \phi_1(X_n)$$

$$\implies \forall X_1^n \in A^n, |\phi_n(X_1^n)| = 2n$$

- Improvement?  $\psi_1 = \begin{array}{|c|c|c|c|} \hline a & c & t & g \\ \hline 0 & 1 & 01 & 10 \\ \hline \end{array}$   
Problem : what is  $\psi_1^{-1}(101)$  :  $ct$  or  $gc$ ?
- *Uniquely decodable* codes :

$$\phi_n(X_1^n) = \phi_m(Y_1^m) \implies m = n \text{ and } \forall i \in \{1, \dots, n\}, X_i = Y_i$$

- *Instant* codes = *prefix* codes :

$$\forall (a, b) \in A^2, \forall j \in \mathbb{N}, \phi_1(a)_{1:j} \neq \phi_1(b)$$

# Kraft's Inequality

**Theorem** (Kraft '49, McMillan '56)

If  $\phi_n$  is a *uniquely decodable* code, then

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

**Proof** for prefix codes : Let  $D = \max_{x \in A^n} |\phi_n(x)|$  and for  $x \in A^n$ , let  $Z(x) = \{w \in \{0, 1\}^D : w_1^{|\phi_n(x)|} = \phi_n(x)\}$ .

- $\#Z(x) = 2^{D-|\phi_n(x)|}$
- $\phi_n$  prefix  $\implies \forall x, y \in A^n, x \neq y \implies Z(x) \cap Z(y) = \emptyset$
- Hence

$$\sum_{x \in A^n} \#Z(x) = \sum_{x \in A^n} 2^{D-|\phi_n(x)|} \leq \#\{0, 1\}^D = 2^D$$



# Shannon's First Theorem

$\implies$  if  $\ell(x) = |\phi_n(x)|$ , we look for

$$\min_{\ell} \sum_{x \in A^n} \ell(x) P^n(x) \text{ under constraint } \sum_{x \in A^n} 2^{-\ell(x)} \leq 1$$

**Theorem** (Shannon '48)

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

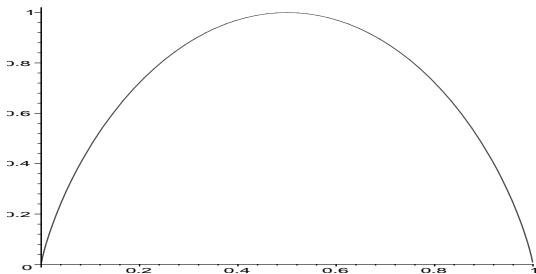
- Misleading notation :  $H_n(X)$  is actually a function of  $P^n$  only, the marginal distribution of  $X_1^n$
- Tight bound : there exist a code  $\phi_n$  such that

$$\mathbb{E}_{\mathbb{P}} [|\phi_n(x)|] \leq H_n(X) + 1$$

# Example : binary entropy

$A = \{0, 1\}$ ,  $X_i \stackrel{iid}{\sim} B(p)$ ,  $0 < p < 1$ .

$$\frac{1}{n} H_n(X) = H_1(X) = h(p) = p \log \frac{1}{p} + (1 - p) \log \frac{1}{1 - p}$$

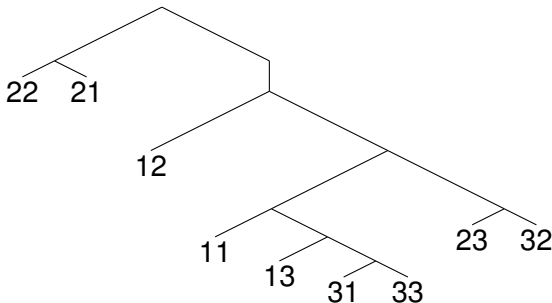


## Example : a 2-block code

$A = \{1, 2, 3\}$ ,  $n = 2$ ,  $P^1 = 0.3\delta_1 + 0.6\delta_2 + 0.1\delta_3$  i.i.d.

$H(X) = -0.3 \times \log(0.3) - 0.6 \times \log(0.6) - 0.1 \times \log(0.1) \approx 1.295$  bits

$X_1^2$	$\phi_2(X_1^2)$
11	1100
12	10
13	11010
21	00
22	01
23	1110
31	110110
32	1111
33	110111



$$\frac{1}{2} \mathbb{E} \left[ \left| \phi_2(X_1^2) \right| \right] \approx 1.335$$

- **Theorem**[Shannon-Breiman-McMillan] if  $P$  is stationary and ergodic, then

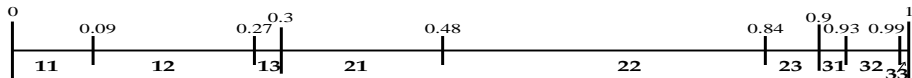
$$\frac{1}{n}H_n(X) \rightarrow H(X)$$

- $\implies$  Interpretation :  $H(X) = \text{minimal number of bits necessary to code each symbol of } X$
- Kraft's inequality : for each code  $\phi_n$  there exists a (sub-)probability distribution  $Q^n$  such that

$$Q^n(x) = 2^{-|\phi_n(x)|} \iff -\log Q^n(x) = |\phi_n(x)|$$

# Arithmetic Coding : a constructive reciprocal

$$\phi_n(x_1^n) = \text{the } \lceil -\log Q^n(x_1^n) \rceil + 2 \text{ first bits of } \frac{Q^n(X_1^n \leq x_1^n) + Q^n(X_1^n < x_1^n)}{2}$$



Ex :  $A = \{1, 2, 3\}$ ,  $n = 2$ ,  
 $Q^2 = [0.3, 0.6, 0.1] \otimes^2$

•  $I(11) = [0, 0.09[ :$

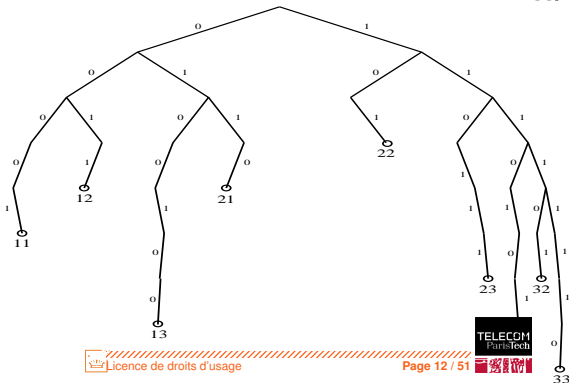
$m(00) = 0.045 = \overline{0.0000101110} \dots$

As  $\lceil \log \frac{1}{0.09} \rceil + 1 = 5$ ,  $\phi_2(00) = 00001$

•  $I(22) = [0.48, 0.84[ :$

$m(22) = 0.66 = \overline{0.1010100011} \dots$

As  $\lceil \log \frac{1}{0.36} \rceil + 1 = 3$ ,  $\phi_2(22) = 101$





# Coding distributions

$\implies -\log Q^n(x) =$  *code length* for  $x$  with *coding distribution*  $Q^n$ .

- Th Shannon : the best coding distribution is  $Q^n = P^n$ .
- Using another distribution causes an expected overhead called *redundancy*

$$\begin{aligned}\mathbb{E}_{\mathbb{P}} [|\phi_n(X_1^n)|] - H(X_1^n) &= \mathbb{E}_{\mathbb{P}} [-\log Q^n(X_1^n) + \log P^n(X_1^n)] \\ &= KL(P^n, Q^n)\end{aligned}$$

= *Kullback-Leibler divergence* of  $P^n$  to  $Q^n$ .

# Universal coding

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

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

$$\Lambda = \{\mathbb{P}_\theta, \theta \in \Theta\}$$

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

⇒ **unavoidable redundancy**, in the worst case :

$$\sup_{\theta \in \Theta} \mathbb{E}_{P_\theta} [|\phi_n(X_1^n)|] - H(X_1^n) = \sup_{\theta \in \Theta} KL(P_\theta^n, Q^n)$$

⇒  $Q^n$  must be close to all the  $\{P_\theta, \theta \in \Theta\}$  for KL divergence

- What if we want to code messages with different or large lengths ?
- Consistent coding family :

$$\forall x_1^n \in A^n, Q^n(x_1^n) = \sum_{a \in A} Q^{n+1}(xa)$$

- Define  $Q(a|x_1^n) = \frac{Q(x_1^n a)}{\sum_{b \in A} Q(x_1^n b)}$ , then

$$Q^n(x_1^n) = \prod_{j=1}^n Q(x_j|x_1^{j-1})$$

- Predictive coding scheme : update  $Q(\cdot|x_1^n)$
- Example : Laplace code  $Q(a|x_1^n) = \frac{n(a, x_1^n) + 1}{n + |A|}$ , where

$$n(a, x_1^n) = \sum_{j=1}^n \mathbb{1}_{x_j=a}$$



# Density estimation with Logarithmic Loss

- Feature space : context  $x_1^{k-1} \in A^*$ ,
- Instant loss : number of bits used  $\ell(Q, x_k) = -\log Q(x_k | x_1^{k-1})$
- Best regressor in average :

$$Q^*(\cdot | x_1^{k-1}) = \mathbb{P}_\theta (X_k = x_k | X_1^{k-1} = x_1^{k-1})$$

- Excess loss :  $\ell(Q, x_k) - \ell(Q^*, x_k) = \log \frac{\mathbb{P}_\theta(X_k = x_k | X_1^{k-1} = x_1^{k-1})}{Q^k(x_k | x_1^{k-1})}$

$$\mathbb{E}_{\mathbb{P}_\theta} [\ell(Q, x_k) - \ell(Q^*, x_k)] = KL \left( \mathbb{P}_\theta(X_k = \cdot | X_1^{k-1} = x_1^{k-1}), Q^k(\cdot | x_1^{k-1}) \right)$$

- Total excess loss :  $L_n(Q) = \sum_{k=1}^n \ell(Q, x_k) - \ell(Q^*, x_k) = \log \frac{P_\theta^n(x_1^n)}{Q^n(x_1^n)}$
- Total expected excess loss :  $\mathbb{E}_{\mathbb{P}_\theta} L_n(Q) = KL(P_\theta^n, Q^n)$



## Lossless Source Coding = density estimation with log-loss

Source Coding and Universal Coding

Bayesian coding

## Context-tree Models and Double Mixture

Rissanen's model

Coding in Context Tree models

The Context-Tree Weighting Method

## Context Tree Weighting on Renewal Processes

CTW on Renewal Processes

Perspectives

- Let  $\pi$  be any prior on  $\Theta$ , take

$$Q^n(x_1^n) = \int P_\theta^n(x_1^n) \pi(d\theta)$$

- Predictive coding distribution :

$$Q(a|x_1^n) = \int \mathbb{P}_\theta(x_{n+1} = a | X_1^n) \pi(d\theta | x_1^n)$$

- Conjugate priors  $\implies$  easy update rules for  $Q(\cdot|x_1^n)$
- Idea :  $\forall \theta \in \Theta, Q^n(x_1^n) \geq P_{\theta \pm \delta}^n(x_1^n) \pi(\theta \pm \delta)$

## Example : i.i.d. classes

- A stationary memoryless source is parameterized by

$$\theta \in \mathcal{S}_{|A|} = \left\{ (\theta_1, \dots, \theta_{|A|}) : \sum_{i=1}^{|A|} \theta_i = 1 \right\}$$

- Conjugate Dirichlet prior  $\pi = \mathcal{D}(\alpha_1, \dots, \alpha_{|A|})$ .
- Posterior :  $\Pi(d\theta | x_1^n) = \mathcal{D}(\alpha_1 + n_1, \dots, \alpha_{|A|} + n_{|A|})$   
 $\implies Q(a | x_1^n) = \frac{n_a + \alpha_a}{n + \sum_{j=1}^{|A|} \alpha_j}$
- Laplace code = special case  $\alpha_1 = \dots = \alpha_{|A|} = 1$

# The Krichevski-Trofimov Mixture

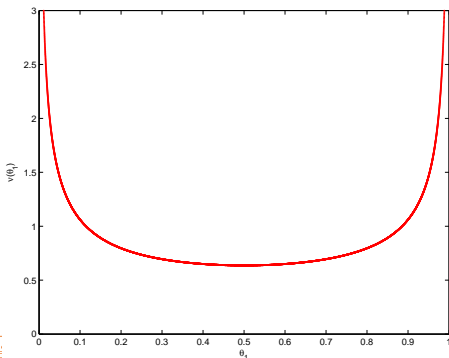
Dirichlet prior with  $\alpha_1 = \dots = \alpha_{|A|} = 1/2$  :

$$Q_{\text{KT}}(x_1^n) = \int_{\theta \in \mathcal{S}_{|A|}} P_{\theta}(x_1^n) \frac{\Gamma\left(\frac{|A|}{2}\right)}{\sqrt{|A|} \Gamma\left(\frac{1}{2}\right)^{|A|}} \prod_{a \in A} \theta_a^{-1/2} d\theta_a.$$

- Jeffrey's prior = Non-informative
- Peaked on extreme parameter values
- Example  $|A|=2$  :

$$\begin{aligned} & Q_{\text{KT}}(001) \\ &= Q_{\text{KT}}(0) Q_{\text{KT}}(0|0) Q_{\text{KT}}(1|00) \end{aligned}$$

$$\frac{1/2}{1} \times \frac{1/2 + 1}{2} \times \frac{1/2}{3}$$



- **Theorem** : (Krichevsky-Trofimov '81) If  $Q_{\text{KT}}^n$  is the KT-mixture code and if  $m = |A|$ , then

$$\max_{\theta \in \Theta} KL(P_{\theta}^n, Q_{\text{KT}}^n) = \frac{m-1}{2} \log \frac{n}{2} + \log \frac{\Gamma(1/2)^m}{\Gamma(\frac{m}{2})} + o_m(1)$$

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

$$\liminf_{n \rightarrow \infty} \min_{Q^n} \max_{\theta \in \Theta} KL(P_{\theta}^n, Q^n) \geq \frac{k}{2} \log n$$

- **Theorem** : Haussler'96  
there exists a Bayesian code  $Q^n$  reaching the minimax redundancy

- **Theorem** : (Krichevsky-Trofimov '81) If  $Q_{\text{KT}}^n$  is the KT-mixture code and if  $m = |A|$ , then

$$\max_{\theta \in \Theta} KL(P_{\theta}^n, Q_{\text{KT}}^n) = \frac{m-1}{2} \log \frac{n}{2} + \log \frac{\Gamma(1/2)^m}{\Gamma(\frac{m}{2})} + o_m(1)$$

- **Theorem** : (Xie-Barron '97)

$$\min_{Q^n} \max_{\theta \in \Theta} KL(P_{\theta}^n, Q^n) = \frac{m-1}{2} \log \frac{n}{2e} + \log \frac{\Gamma(1/2)^m}{\Gamma(\frac{m}{2})} + o_m(1)$$

- Actually, **pointwise** approximation :

$$-\log Q_{\text{KT}}^n(x_1^n) \leq \inf_{\theta} -\log P_{\theta}^n(x_1^n) + \frac{|A|-1}{2} \log n + C$$

- iid sources on finite alphabets are very-well treated. . .
- . . .but they provide poor models for most applications !
- example : Shannon considered **Markov models** for natural languages
- Pb : a Markov Chain of order 3 on an alphabet of size 26 has  $k = 439400$  degrees of freedom !

$$\min_{Q^n} \max_{\theta \in \Theta} KL(P_{\theta}^n, Q^n) \geq \frac{k}{2} \log n$$

$\implies$  need of more *parcimonious* classes





## Lossless Source Coding = density estimation with log-loss

Source Coding and Universal Coding

Bayesian coding

## Context-tree Models and Double Mixture

Rissanen's model

Coding in Context Tree models

The Context-Tree Weighting Method

## Context Tree Weighting on Renewal Processes

CTW on Renewal Processes

Perspectives

## Formal definition :

- **Context Tree** :  $T \subset A^*$  such that

$$\forall x_{-\infty}^0, \exists \text{ unique } s \in T : x_{-|s|}^{-1} = s$$

The unique suffix of  $x_{-\infty}^{-1}$  in  $T$  is denoted by  $T(x)$

- A string  $s$  is a **context** of a stationary ergodic process  $\mathbb{P}$  if  $\mathbb{P}(X_{-|s|}^{-1} = s) > 0$  and if for all  $x_{-\infty}^{-|s|-1}$  :

$$\mathbb{P}(X_0 = a \mid X_{-\infty}^{-1} = x_{-\infty}^{-|s|-1} s) = \mathbb{P}(X_0 = a \mid X_{-|s|}^{-1} = s)$$

- $T$  the **context tree** of  $\mathbb{P}$  if  $T$  is a context tree, contains all contexts of  $\mathbb{P}$  and if none of its elements can be replaced by a proper suffix without violating one of these properties

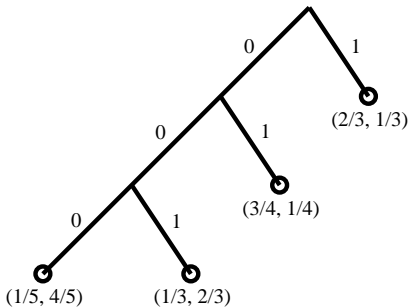
# Context Tree Sources

**Informal Definition** A Context tree Source is a Markov Chain whose order is allowed to depend on the past data.

$$T = \{1, 10, 100, 000\}$$

$$P(X_1^4 = 00110 | X_{-1}^0 = 10)$$

$$\begin{aligned} &= P(X_1 = 0 | X_{-1}^0 = 10) && 3/4 \\ &\times P(X_2 = 0 | X_{-1}^1 = 100) && 1/3 \\ &\times P(X_3 = 1 | X_{-1}^2 = 1000) && 4/5 \\ &\times P(X_4 = 1 | X_{-1}^3 = 10001) && 1/3 \\ &\times P(X_5 = 0 | X_{-1}^4 = 100011) && 2/3 \end{aligned}$$



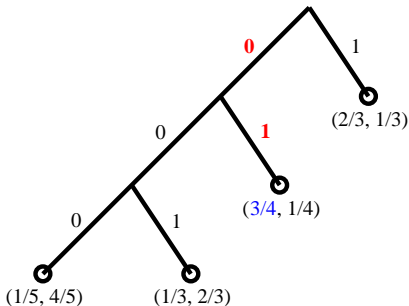
# Context Tree Sources

**Informal Definition** A Context tree Source is a Markov Chain whose order is allowed to depend on the past data.

$$T = \{1, 10, 100, 000\}$$

$$P(X_1^4 = 00110 | X_{-1}^0 = 10)$$

$$\begin{aligned} &= P(X_1 = \mathbf{0} | X_{-1}^0 = \mathbf{10}) && \mathbf{3/4} \\ &\times P(X_2 = 0 | X_{-1}^1 = 100) && 1/3 \\ &\times P(X_3 = 1 | X_{-1}^2 = 1000) && 4/5 \\ &\times P(X_4 = 1 | X_{-1}^3 = 10001) && 1/3 \\ &\times P(X_5 = 0 | X_{-1}^4 = 100011) && 2/3 \end{aligned}$$



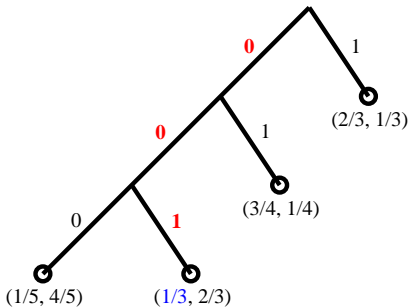
# Context Tree Sources

**Informal Definition** A Context tree Source is a Markov Chain whose order is allowed to depend on the past data.

$$T = \{1, 10, 100, 000\}$$

$$P(X_1^4 = 00110 | X_{-1}^0 = 10)$$

$$\begin{aligned} &= P(X_1 = 0 | X_{-1}^0 = 10) && 3/4 \\ &\times P(X_2 = 0 | X_{-1}^1 = 100) && 1/3 \\ &\times P(X_3 = 1 | X_{-1}^2 = 1000) && 4/5 \\ &\times P(X_4 = 1 | X_{-1}^3 = 10001) && 1/3 \\ &\times P(X_5 = 0 | X_{-1}^4 = 100011) && 2/3 \end{aligned}$$



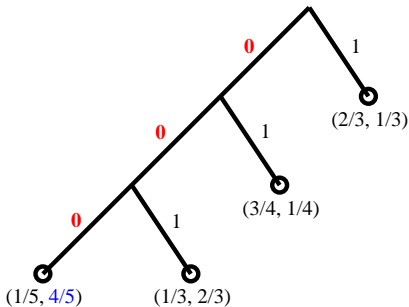
# Context Tree Sources

**Informal Definition** A Context tree Source is a Markov Chain whose order is allowed to depend on the past data.

$$T = \{1, 10, 100, 000\}$$

$$P(X_1^4 = 00110 | X_{-1}^0 = 10)$$

$$\begin{aligned} &= P(X_1 = 0 | X_{-1}^0 = 10) && 3/4 \\ &\times P(X_2 = 0 | X_{-1}^1 = 100) && 1/3 \\ &\times P(X_3 = \mathbf{1} | X_{-1}^2 = \mathbf{1000}) && \mathbf{4/5} \\ &\times P(X_4 = 1 | X_{-1}^3 = 10001) && 1/3 \\ &\times P(X_5 = 0 | X_{-1}^4 = 100011) && 2/3 \end{aligned}$$



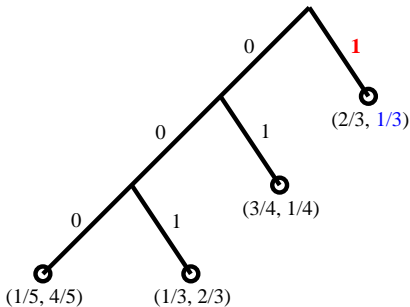
# Context Tree Sources

**Informal Definition** A Context tree Source is a Markov Chain whose order is allowed to depend on the past data.

$$T = \{1, 10, 100, 000\}$$

$$P(X_1^4 = 00110 | X_{-1}^0 = 10)$$

$$\begin{aligned} &= P(X_1 = 0 | X_{-1}^0 = 10) && 3/4 \\ &\times P(X_2 = 0 | X_{-1}^1 = 100) && 1/3 \\ &\times P(X_3 = 1 | X_{-1}^2 = 1000) && 4/5 \\ &\times P(X_4 = \mathbf{1} | X_{-1}^3 = 1000\mathbf{1}) && \mathbf{1/3} \\ &\times P(X_5 = 0 | X_{-1}^4 = 10001\mathbf{1}) && 2/3 \end{aligned}$$



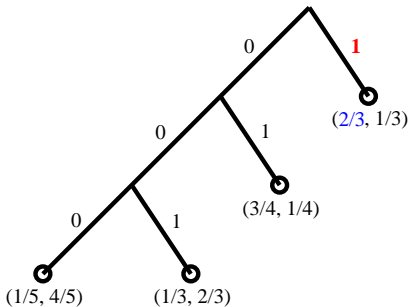
# Context Tree Sources

**Informal Definition** A Context tree Source is a Markov Chain whose order is allowed to depend on the past data.

$$T = \{1, 10, 100, 000\}$$

$$P(X_1^4 = 00110 | X_{-1}^0 = 10)$$

$$\begin{aligned} &= P(X_1 = 0 | X_{-1}^0 = 10) && 3/4 \\ &\times P(X_2 = 0 | X_{-1}^1 = 100) && 1/3 \\ &\times P(X_3 = 1 | X_{-1}^2 = 1000) && 4/5 \\ &\times P(X_4 = 1 | X_{-1}^3 = 10001) && 1/3 \\ &\times P(X_5 = \mathbf{0} | X_{-1}^4 = 10001\mathbf{1}) && \mathbf{2/3} \end{aligned}$$



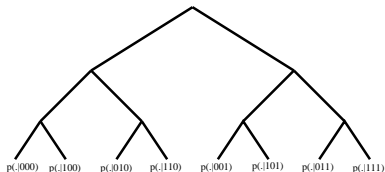


# Context Trees versus Markov Chains

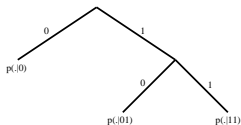
- Markov chain of order  $r$  = context tree source corresponding to a complete tree of depth  $r$

Markov chain of order 3

$$M = \begin{pmatrix} p(\cdot|000) \\ p(\cdot|100) \\ \vdots \\ p(\cdot|111) \end{pmatrix} \Rightarrow$$



- Finite context tree source of depth  $d$  = Markov Chain of order  $d$



$\Rightarrow$

$$M = \begin{pmatrix} p(\cdot|0) \\ p(\cdot|10) \\ p(\cdot|101) \\ p(\cdot|11) \end{pmatrix}$$

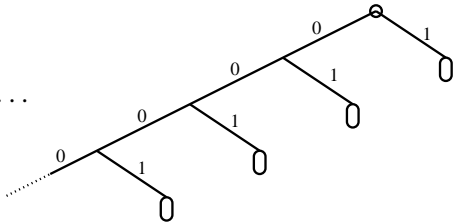
$\Rightarrow$  much more flexibility : large number of models per dimension

There exist (useful and simple) non-markovian context tree sources, see Galves, Ferrari, Comets,...

## Ex : Renewal Processes

$$X_{-\infty}^{\infty} = \dots 1 \underbrace{00 \dots 1}_{N_i} \underbrace{00 \dots 1}_{N_{i+1}} \dots$$

$$N_i \stackrel{iid}{\sim} \mu \in \mathfrak{M}_1(\mathbb{N}_+)$$



$$\mathbb{P} \left( X_0 = 1 \mid X_{-\infty}^{-1} = \dots 10^{j-1} \right) = \frac{\mu(j)}{\mu([j, \infty[)}$$



## Lossless Source Coding = density estimation with log-loss

Source Coding and Universal Coding

Bayesian coding

## Context-tree Models and Double Mixture

Rissanen's model

Coding in Context Tree models

The Context-Tree Weighting Method

## Context Tree Weighting on Renewal Processes

CTW on Renewal Processes

Perspectives

- For a given context tree  $T$ , let

$$\Lambda_T = \{\mathbb{P} \text{ stationary ergodic on } A : \mathcal{T}(\mathbb{P}) = T\}$$

- $\mathbb{P} \in \Lambda_T$  is fully characterized by  $\theta = (\theta_s)_{s \in T} \in \Theta_T = \mathcal{S}_{|A|}^T \sim \mathbb{R}^{|T|(|A|-1)}$  :

$$\mathbb{P} \left( X_0 = \cdot | X_{-\infty}^{-1} = x_{-\infty}^{-1} \right) = \theta_{T(x_{-\infty}^{-1})}(\cdot)$$

- Bayesian coding : prior  $\pi_T$  on  $\Theta_T$ ,

$$Q_{\pi_T}(x) = \int \mathbb{P}_{\theta}(x_1^n) d\pi_T(\theta)$$

- Prior : the  $\theta_s$  are independent,  $\mathcal{D}(\frac{1}{2}, \dots, \frac{1}{2})$ - distributed :

$$\pi_T(d\theta) = \bigotimes_{s \in T} \pi_{KT}(d\theta_s)$$

- For  $s \in T$  let  $T(x_1^n, s)$  = substring of  $x_1^n$  that occurs in context  $s$   
Example : for  $T = \{00, 10, 1\}$  and  $x_0^8 = 1 \text{ 0010001011}$ ,

$$T(x_1^n, 00) = 0101 \quad T(x_1^n, 10) = 001 \quad T(x_1^n, 1) = 0001$$

- The KT-mixture on model  $T$  is :

$$Q_T^n(x_1^n) = \prod_{s \in T} Q_{KT}^{|T(x_1^n, s)|}(T(x_1^n, s))$$

- **Theorem :** (Willems, Shtar'kov, Tjalkens '93) there is a constant  $C$  such that :

$$-\log Q_T^n(x_1^n) \leq \inf_{\theta \in \mathcal{S}_{|A|}^T} -\log \mathbb{P}_\theta(x_1^n | x_{-\infty}^{-1}) + |T| \frac{|A| - 1}{2} \log \left( \frac{n}{|T|} \right) + C|T|$$

- first-order optimal
- What if we don't know which model  $T$  to use ?

$\implies$  model selection

# Minimum Description Length Principle

- Occam's razor (Jorma Rissanen '78) :

Choose the model that gives the shortest description of data

- Information theory :

objective codelength = codelength of a minimax coder.

- Estimator associated with minimax mixtures  $(\pi_T)_T$  :

$$\hat{T}_{BIC} = \arg \min_T - \log \int_{\theta \in \Theta_T} P_{\theta}^n(x_1^n) \pi_T(d\theta).$$

- Looking directly at the bounds :

$$\hat{T}_{KT} = \arg \min_T \inf_{\theta \in \Theta_T} - \log P_{\theta}^n(x_1^n) + \frac{|T|(|A| - 1)}{2} \log n.$$

**Theorems** (Csiszár& Shields '96, Csiszár& Talata '04, G. '05, Galves & Leonardi '07, Leonardi '09) :

- the BIC estimator is strongly consistent
- the Krichevski-Trofimov estimator is not consistent

Many Rates of convergence, minimal penalties, . . .

The Krichevski-Trofimov estimator fails to identify Bernoulli-1/2 iid sources !

One can derive from this bounds on  $\sup_T \sup_{\theta \in \Theta_T} KL \left( P_{\theta}^n, Q_{\hat{T}}^n \right)$





## Lossless Source Coding = density estimation with log-loss

Source Coding and Universal Coding

Bayesian coding

## Context-tree Models and Double Mixture

Rissanen's model

Coding in Context Tree models

The Context-Tree Weighting Method

## Context Tree Weighting on Renewal Processes

CTW on Renewal Processes

Perspectives

- $\nu$  = measure on the set  $\mathcal{T}$  of all context trees defined by :

$$\nu(T) = 2^{-2|T|+1}$$

is a **probability distribution**.

- **Context Tree Weighting** coding distribution :

$$Q_{\text{CTW}}^n(x_1^n) = \sum_T \nu(T) Q_T^n(x_1^n)$$

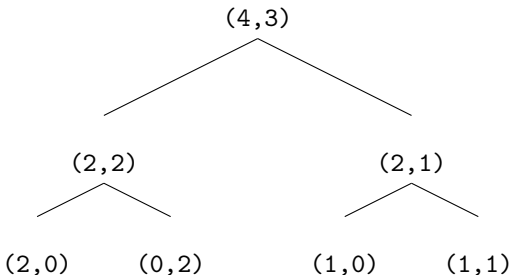
- **Theorem** (Shtarkov and al '93)) : oracle inequality :

$$\begin{aligned} -\log Q_{\text{CTW}}^n(x_1^n | x_{-\infty}^0) &\leq \inf_{T \in \mathcal{T}} \inf_{\theta \in \Theta_T} p_{\theta}(x_1^n | x_{-\infty}^0) \\ &+ \frac{|A| - 1}{2} |T| \log^+ \frac{n}{T} + |T|(2 + \log m) + m - 2. \end{aligned}$$



# CTW : Efficient Algorithm

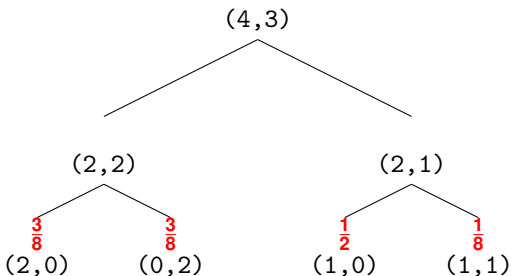
**Idea :** in each node, compute the arithmetical mean of a *self-prob* and a *sub-prob*.



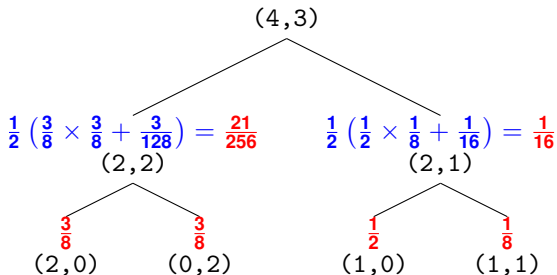


# CTW : Efficient Algorithm

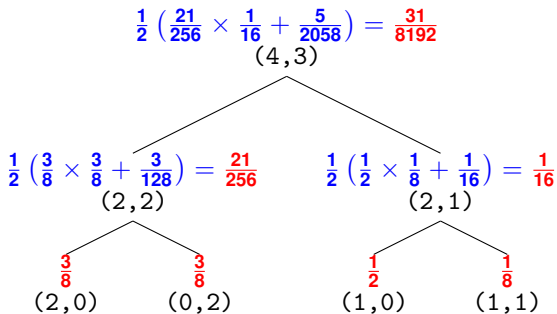
**Idea :** in each node, compute the arithmetical mean of a *self-prob* and a *sub-prob*.



**Idea :** in each node, compute the arithmetical mean of a *self-prob* and a *sub-prob*.



**Idea :** in each node, compute the arithmetical mean of a *self-prob* and a *sub-prob*.





## Lossless Source Coding = density estimation with log-loss

Source Coding and Universal Coding

Bayesian coding

## Context-tree Models and Double Mixture

Rissanen's model

Coding in Context Tree models

The Context-Tree Weighting Method

## Context Tree Weighting on Renewal Processes

CTW on Renewal Processes

Perspectives

# Renewal Processes

- Let  $\mathcal{R}$  be the (non-parametric) class of renewal processes on the binary alphabet  $\{0, 1\}$  :

$$X = \dots 1 \underbrace{00 \dots 1}_{N_i} \underbrace{00 \dots 1}_{N_{i+1}} \dots, \quad N_i \stackrel{iid}{\sim} \mu \in \mathfrak{M}_1(\mathbb{N}_+)$$

- Theorem** : (Csiszár and Shields '96) There exist two constants  $c$  and  $C$  such that

$$c\sqrt{n} \leq \inf_{Q^n \in \mathcal{R}} \sup_{P^n} KL(P^n, Q^n) = C\sqrt{n}$$

$\implies$  first example of an *intermediate complexity* class.

- Problem : non-constructive
- Does a *general purpose* coder behave well on Renewal

Processes ?



# Redundancy of the CTW method on Renewal Processes

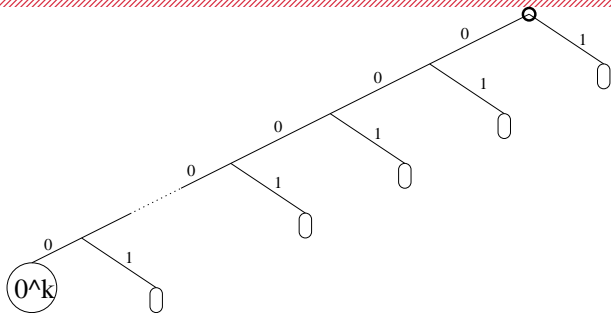
**Theorem :** (G. '04)

There exist constants  $c$  and  $C$  such that for all  $\mathbb{P} \in \mathcal{R}$  :

$$c\sqrt{n}\log n \leq KL(\mathbb{P}^n, Q_{\text{CTW}}^n) \leq C\sqrt{n}\log n$$

- CTW is efficient on long-memory processes
- *Adaptivity* : if the renewal distribution is bounded, CTW achieves regret  $O(\log n)$  (contrary to ad-hoc coders).
- **Requires deep contexts** in the double mixture ( $\implies$  the tree should not be cut off at depth  $\log n$ ).
- Kind of **non-parametric estimation** : need for a balance between approximation and estimation.

# Why it works



$\mathbb{P}^n$  is approximated by  $Q_T^n$  with  $T$  as above. Two losses :

- estimating the transition parameters :  $k \log(n)/2$
- approximation due to boundedness :  $n \log(k)/k$

$\implies$  for  $k = c\sqrt{n}$ , both are  $O(\sqrt{n} \log(n))$

## Extension : Markovian Renewal Processes

- A *Markovian Renewal Processes* is such that successive distances between occurrences of symbol 1 form a Markov Chain on  $\mathbb{N}_+$  :

$$X = \dots 1 \underbrace{00 \dots 1}_{N_i} \underbrace{00 \dots 1}_{N_{i+1}} \dots, \quad (N_i)_i \text{ Markov chain}$$

- **Theorem** : (Csiszár and Shields '96) There exist two constants  $c$  and  $C$  such that

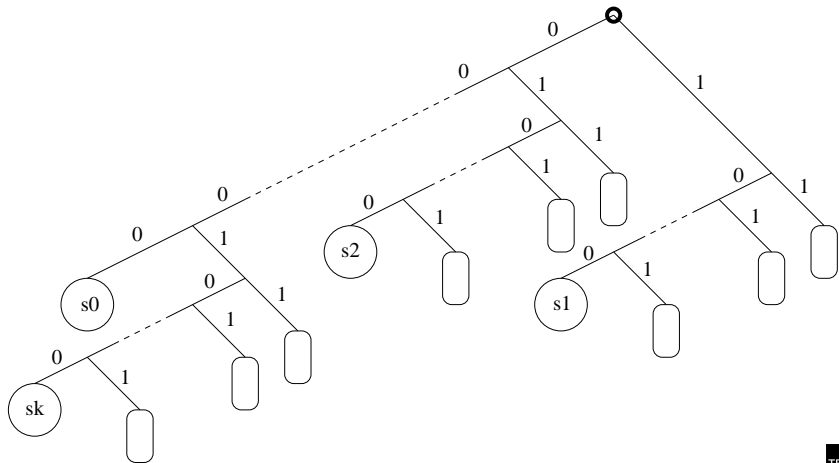
$$cn^{2/3} \leq \inf_{Q^n \in \mathcal{MR}} \sup_{P^n} KL(\mathbb{P}^n, Q^n) \leq Cn^{2/3}$$

- **Theorem** : (G. '04) CTW is almost adaptive on the class  $\mathcal{MR}$  of Markovian Renewal processes :  $\exists C \in \mathbb{R}, \forall P \in \mathcal{MR},$

$$KL(P^n, Q_{CTW}^n) \leq Cn^{2/3} \log n$$

# Markovian Renewal Processes are Context Tree Sources

And so are r-order Markov Renewal Processes





## Lossless Source Coding = density estimation with log-loss

Source Coding and Universal Coding

Bayesian coding

## Context-tree Models and Double Mixture

Rissanen's model

Coding in Context Tree models

The Context-Tree Weighting Method

## Context Tree Weighting on Renewal Processes

CTW on Renewal Processes

Perspectives



# Perspectives : oracle inequalities on more general classes ?

- CTW is probably more useful than expected : it obtains good results on infinite memory processes
- We want to obtain oracle inequalities for different non-markovian classes
- But the setting is more difficult than usually
- And approximation theory is difficult !  
Ex : let  $\mathbb{P}_H$  be a HMM with 3 hidden states, what is

$$\inf_{|T| \leq k} \inf_{\theta \in \mathcal{S}_{|A|}^T} KL(P_H^n, P_\theta^n) \quad ???$$



the end...

---

Thank you for your attention !

