
Lossless Source Coding

Gadiel Seroussi Marcelo Weinberger

Information Theory Research
Hewlett-Packard Laboratories
Palo Alto, California

Notation

- A = **discrete (usually finite) alphabet**
- $\alpha = |A|$ = **size of A (when finite)**
- $x_1^n = x^n = x_1x_2x_3 \dots x_n =$ **finite sequence over A**
- $x_1^\infty = x^\infty = x_1x_2x_3 \dots x_t \dots =$ **infinite sequence over A**
- $x_i^j = x_ix_{i+1} \dots x_j =$ **sub-sequence (i sometimes omitted if = 1)**
- $p_X(x) = \text{Prob}(X=x) =$ **probability mass function (PMF) on A**
(subscript X and argument x dropped if clear from context)
- $X \sim p(x) :$ **X obeys PMF $p(x)$**
- $E_p[F] =$ **expectation of F w.r.t. PMF p (subscript and $[\]$ may be dropped)**
- $\hat{p}_{x_1^n}(x) =$ **empirical distribution obtained from x_1^n**
- $\log x =$ **logarithm to base 2 of x , unless base otherwise specified**
- $\ln x =$ **natural logarithm of x**
- $H(X), H(p) =$ **entropy of a random variable X or PMF p , in bits; also**
- $H(p) = -p \log p - (1-p) \log (1-p), 0 \leq p \leq 1 :$ **binary entropy function**
- $D(p||q) =$ **relative entropy (information divergence) between PMFs p and q**

Lossless Source Coding

3. Universal Coding

Universal Modeling and Coding

q So far, it was assumed that a model of the data is available, and we aimed at compressing the data optimally w.r.t. the model

q By Kraft's inequality, the models we use can be expressed in terms of probability distributions:

For every UD code with length function $L(s)$, we have $\sum_{s \in A^n} 2^{-L(s)} \leq 1$
string of length n over $\xrightarrow{\quad}$
finite alphabet A

\Rightarrow a code defines a probability distribution $P(s) = 2^{-L(s)}$ over A^n

q Conversely, given a distribution $P(\cdot)$ (a model), there exists a UD code that assigns $\lceil -\log p(s) \rceil$ bits to s (Shannon code)

q Hence, $P(\cdot)$ serves as a model to encode s , and every code has an associated model

| a (probabilistic) model is a tool to “understand” and predict the behavior of the data

Universal Modeling and Coding (cont.)

q Given a model $P(\cdot)$ on n -tuples, arithmetic coding provides an effective mean to **sequentially** assign a code word of length close to $-\log P(s)$ to s

| we don't need to see the whole string $s = x_1 x_2 \dots x_n$: encode x_t using conditional probability

$$p(x_t | x_1 x_2 \dots x_{t-1}) = \frac{P(x^t)}{P(x^{t-1})} = \frac{\sum_{u \in A^{n-t}} P(x^t u)}{\sum_{v \in A^{n-t+1}} P(x^{t-1} v)}$$

| the model probabilities can vary arbitrarily and “adapt” to the data

| in a **strongly** sequential model, $P(x^t)$ is independent of n

q **CODING SYSTEM = MODEL + CODING UNIT**

| two separate problems: design a model and use it to encode

We will view data compression as a problem of assigning probabilities to data

Coding with Model Classes

- q **Universal data compression** deals with the optimal description of data in the absence of a given model
 - | in most practical applications, the model is not given to us
- q How do we make the concept of “optimality” meaningful?
 - | there is always a code that assigns just 1 bit to the data at hand!

The answer: Model classes

- q We want a “universal” code to perform as well as the best model in a **given** class \mathcal{C} for **any** string s , where the best competing model changes from string to string
 - | universality makes sense only w.r.t. a model class
- q A code with length function $L(x^n)$ is **pointwise universal** w.r.t. a class \mathcal{C} if $R_c(L, x^n) = \frac{1}{n} [L(x^n) - \min_{C \in \mathcal{C}} L_C(x^n)] \rightarrow 0$ when $n \rightarrow \infty$
pointwise 
redundancy

How to Choose a Model Class?

Universal coding tells us how to encode optimally w.r.t. to a class; it doesn't tell us how to choose a class!

q Some possible criteria:

- | complexity**
- | prior knowledge on the data**
- | some popular models were already presented**

q We will see that the bigger the class, the slower the best possible convergence rate of the redundancy to 0

- | in this sense, prior knowledge is of paramount importance: don't learn what you already know!**

Ultimately, the choice of model class is an art

Example: Bernoulli Models

$$A = \{0,1\}, \quad C = \{P_\theta, \theta \in \Lambda = [0,1]\}$$

 **i.i.d. distribution
with parameter $p(1) = \theta$**

$$\min_{C \in \mathcal{C}} L_C(x^n) = \min_{q \in \Lambda} \log \frac{1}{P_q(x^n)} = \log \frac{1}{P_{\hat{q}(x^n)}(x^n)} = nH(x^n)$$

 **ML-estimate of θ**

Therefore, our goal is to find a code such that

$$\frac{L(x^n)}{n} - H(x^n) \rightarrow 0$$

q Here is a trivial example of a universal code:

Use $\lceil \log(n+1) \rceil$ bits to encode n_1 , and then “tune” your Shannon code to the parameter $\theta = n_1/n$, which is precisely the ML-estimate

\Rightarrow this is equivalent to telling the decoder what the best model is!

Bernoulli Models: Enumerative Code

q The total code length for this code satisfies:

$$\frac{L(x^n)}{n} \leq -\binom{n_0}{n} \log \binom{n_0}{n} - \binom{n_1}{n} \log \binom{n_1}{n} + \frac{\log(n+1)+2}{n} = \hat{H}(x^n) + \underbrace{\frac{\log(n+1)+2}{n}}_{\rightarrow 0}$$

| wasteful: knowledge of n_1 already discards part of the sequences, to which we should not reserve a code word

q A slightly better code: enumerate all $\binom{n}{n_1}$ sequences with n_1 ones, and describe the sequence with an index

$$\frac{L'(x^n)}{n} = \frac{1}{n} \left\lceil \log \binom{n}{n_1} \right\rceil + \frac{\lceil \log(n+1) \rceil}{n} \leq \hat{H}(x^n) + \frac{\log(n+1)+2}{n}$$

Stirling: $\binom{n}{n_1} \leq \frac{2^{n\hat{H}(x^n)}}{\sqrt{2\pi(n_1 n_0)/n}} \Rightarrow$ for sequences such that n_1/n is bounded away from 0 and 1, $\frac{L'(x^n)}{n} \leq \hat{H}(x^n) + \frac{\log n}{2n} + O(1/n)$

Bernoulli Models: Mixture Code

q The enumerative code length is close to $-\log Q'(x^n)$, where

q Is $Q(\cdot)$ a probability assignment?

$$Q'(x^n) = \binom{n}{n_1}^{-1} \cdot \frac{1}{n+1}$$

$$\int_0^1 P_q(x^n) dq = \int_0^1 q^{n_1} (1-q)^{n_0} dq = \binom{n}{n_1}^{-1} \cdot \frac{1}{n+1} = Q'(x^n) \Rightarrow \sum_{x^n \in A^n} Q'(x^n) = 1$$

q Aha!! So we get the same result by **mixing** all the models in the class and using the mixture as a model!

q This **uniform** mixture is very appealing because it can be **sequentially** implemented, independent of n :

$$Q'(x^n) = \frac{n_0!n_1!}{(n+1)!} = \prod_{t=0}^{n-1} q'(x_{t+1} | x^t) \quad \text{where} \quad q'(0 | x^t) = \frac{n_0(x^t) + 1}{t + 2}$$

| Laplace's rule of succession!

| this results in a **"plug-in" approach**: estimate θ and plug it in

Bernoulli Models: Mixture Code (cont.)

q Maybe we can make Stirling work for us with a different mixture, that puts more weight in the problematic regions where n_1/n approaches 0 and 1?

q Consider Dirichlet's density $w(q) = \frac{1}{\Gamma(\frac{1}{2})^2 \sqrt{q(1-q)}} \Rightarrow$

$$Q''(x^n) = \int_0^1 P_q(x^n) dw(q) = \frac{1}{\Gamma(\frac{1}{2})^2} \int_0^1 q^{n_1 - \frac{1}{2}} (1-q)^{n_0 - \frac{1}{2}} dq = \frac{\Gamma(n_0 + \frac{1}{2})\Gamma(n_1 + \frac{1}{2})}{n! \Gamma(\frac{1}{2})^2}$$

\Rightarrow by Stirling for the Gamma function, for **all** sequences x^n

$$\frac{L''(x^n)}{n} = -\frac{1}{n} \log Q''(x^n) \leq \hat{H}(x^n) + \frac{\log n}{2n} + O(1/n) \quad \text{Can we do any better?}$$

q This mixture also has a plug-in interpretation:

$$Q''(x^n) = \prod_{t=0}^{n-1} q''(x_{t+1} | x^t) \quad \text{where} \quad q''(0 | x^t) = \frac{n_0(x^t) + \frac{1}{2}}{t+1}$$

Summary of Codes

q Two-part code

- | describe best parameter and then code based on it: works whenever the number of possible optimal codes in the class is **sub-exponential**
- | the most natural approach (akin to Kolmogorov complexity), but not sequential
- | not efficient in the example: can be improved by describing the parameter more coarsely \Rightarrow trade-off: spend less bits on the parameter and use an **approximation** of best parameter
- | alternatively, don't allocate code words to impossible sequences

q Mixture code

- | can be implemented sequentially
- | a suitable **prior** on the parameter gave the smallest redundancy

q Plug-in codes

- | use a biased estimate of the conditional probability in the model class based on the data seen so far
- | can often be interpreted as mixtures

q But what's the best we can do?

Normalized Maximum-Likelihood Code

- q **Goal:** find a code that attains the best worst-case pointwise redundancy

$$R_C = \min_L \max_{x^n \in A^n} R_C(L, x^n) = \frac{1}{n} \min_L \max_{x^n \in A^n} [L(x^n) - \min_{C \in \mathcal{C}} L_C(x^n)]$$

- q **Consider the code defined by the probability assignment**

$$Q(x^n) = \frac{2^{-\min_{C \in \mathcal{C}} L_C(x^n)}}{\sum_{x^n \in A^n} 2^{-\min_{C \in \mathcal{C}} L_C(x^n)}} \quad \text{D} \quad R_C(L, x^n) = \frac{1}{n} \log \left[\sum_{x^n \in A^n} 2^{-\min_{C \in \mathcal{C}} L_C(x^n)} \right]$$

This quantity depends only on the class!

- q **Since the redundancy of this code is the same for all x^n , it must be optimal in the minimax sense!**
- q **Drawbacks:** - sequential probability assignment depends on horizon n
- hard to compute

Example: NML Code for the Bernoulli Class

q We evaluate the minimax redundancy R_C for the Bernoulli class:

$$R_C = \frac{1}{n} \log \left[\sum_{x^n \in A^n} 2^{-H(x^n)} \right] = \frac{1}{n} \log \left[\sum_{i=0}^n \binom{n}{i} 2^{-h(\frac{i}{n})} \right] = \frac{1}{2n} \log \frac{np}{2} + o(1/n)$$

typical of sufficient statistics:

$$P_q(x^n) = p(x^n | s(x^n)) p_q(s(x^n))$$

Hint: Stirling +

$$\int_0^1 \frac{dq}{\sqrt{q(1-q)}} = \Gamma(\frac{1}{2})^2 = \pi$$

⇒ the pointwise redundancy cannot vanish faster than $(\log n)/2n$
(Dirichlet mixture or NML code)

Parametric Model Classes

q A useful limitation to the model class is to assume $\mathcal{C} = \{P_\theta, \theta \in \Theta_d\}$

a parameter space
of dimension d

q Examples:

| Bernoulli: $d = 1$, general i.i.d. model: $d = \alpha - 1$ ($\alpha = |A|$)

| FSM model with k states: $d = k(\alpha - 1)$

| memoryless geometric distribution on the integers $i \geq 0$: $P(i) = \theta^i (1 - \theta)$
 $d = 1$

q The dimension of the parameter space (number of one-dimensional parameters) plays a fundamental role in modeling problems

| with more parameters we can better fit the model to the data (e.g., a Markov model of higher order \Rightarrow the entropy cannot increase)

| but on the other hand, the class is richer and the redundancy is higher: e.g., in a two-part code it takes more bits to describe the best parameter

we will quantify this tradeoff of the model selection step

Minimax Redundancy for Parametric Classes

q For a parametric class, $\min_{C \in \mathcal{C}} L_C(x^n) = -\log P_{\hat{q}(x^n)}(x^n)$

ML-estimate
of θ based on x^n

q A code corresponding to a distribution $Q(x^n)$, has a redundancy

$\log \frac{P_{\hat{q}(x^n)}(x^n)}{Q(x^n)}$ and the minimax redundancy (with the NML code) is

$$R_C = \frac{1}{n} \log \left[\sum_{x^n \in A^n} P_{\hat{q}(x^n)}(x^n) \right]$$

| the more sensitive the class to the parameter, the larger the redundancy

Assumptions on the Parametric Model Classes

- q We will assume that each model in the class satisfies the marginality condition for a random process, namely

$$\sum_{x_{n+1} \in A} P_q(x^{n+1}) = P_q(x^n)$$

- | this means that each model can be implemented in a strongly sequential manner
- q We will usually assume that the model class is “nice”, including:
 - | Θ_d is an open bounded subset of \mathcal{R}^d which includes the ML estimates
 - | the Fisher information matrix $\mathbf{I}(\theta)$ is “nice” $I_{ij}(\mathbf{q}) = -\lim_{n \rightarrow \infty} \frac{1}{n} E \left[\frac{\partial^2 \ln P_q(x^n)}{\partial q_i \partial q_j} \right]$
 - u for i.i.d., this is the classical Fisher information
 - | the maximum-likelihood estimator satisfies the CLT: the distribution of $\sqrt{n}(\hat{\mathbf{q}}(x^n) - \mathbf{q})$ converges to $N(0, \mathbf{I}^{-1}(\theta))$

Redundancy of NML Code

q **Theorem:**
$$R_C = \frac{d}{2n} \log \frac{n}{2p} + \frac{1}{n} \log \int_{\Theta_d} \sqrt{|I(q)|} dq + o(1/n)$$

As expected, R_C grows with the number of parameters

q **Idea of the proof:**

- | partition the parameter space into small hypercubes with sides $r = o(1/\sqrt{n})$
- | represent each hypercube $V_r(\theta^r)$ by a parameter θ^r belonging to it
- | associate to $V_r(\theta^r)$ the mass
$$P_r(q^r) = \sum_{q(x^n) \in V_r(q^r)} P_{q^r}(x^n)$$
- | with a Taylor expansion of $P_\theta(x^n)$ (as a function of θ) around its maximum $\hat{q}(x^n)$ we get the exact asymptotics of $R_C \rightarrow \sum_{q^r} P_r(q^r)$
- | for each hypercube, approximate $P_r(\theta^r)$ using CLT \Rightarrow integral of a normal distribution

Interpretation of NML as a Two-part Code

q Encode x^n in two parts:

- | **Part I:** encode the parameter θ^r for the hypercube $V_r(\theta^r)$ into which $\hat{q}'(x^n)$ falls, with a code tuned to the distribution $P_r(q^r) / \sum_{q^r} P_r(q^r)$
- | **Part II:** encode x^n using the model \hat{q}'^r ; **conditioned** on the fact that the ML estimate belongs to $V_r(\hat{q}'^r) \Rightarrow$ distribution $\frac{P_{\hat{q}'^r}(x^n)}{P_r(\hat{q}'^r)}$

q In earlier two-part codes, the parameters were represented with a **fixed** precision $O(1/\sqrt{n})$ and then x^n was encoded based on the **approximate** ML parameter

- | evaluating the approximation cost with Taylor, this precision was shown to be optimal
- | intuitively, $1/\sqrt{n}$ is the magnitude of the estimation error in $\hat{q}'(x^n)$ and therefore there is no need to encode the estimator with better precision

$$\frac{d}{2} \log n = \text{MODEL COST}$$

Redundancy of the Mixture code

$$Q_w(x^n) = \int_{q \in \Theta_d} P_q(x^n) dw(q)$$

q As in the Bernoulli example, for any i.i.d. **exponential family**, and for **Markov** models, mixture code also gives the same “magic” model cost when the distribution $w(\theta)$ is proportional to $\sqrt{|I(q)|}$ (Jeffrey’s prior)

- | in the Bernoulli case, this is precisely Dirichlet’s distribution: $I(q) = \frac{1}{q(1-q)}$
- | the tool for solving the integral is **Laplace’s integration method**

q The advantage of mixtures is that they yield sequential **horizon-free**

codes:
$$\sum_{x_t \in A} Q_w(x^t) = \int_{q \in \Theta_d} \sum_{x_t \in A} P_q(x^t) dw(q) = \int_{q \in \Theta_d} P_q(x^{t-1}) dw(q) = Q_w(x^{t-1})$$

$$q_w(x_t | x^{t-1}) = \frac{Q_w(x^t)}{Q_w(x^{t-1})} = \int_{q \in \Theta_d} P_q(x_t | x^{t-1}) \frac{P_q(x^{t-1})}{\int_{q \in \Theta_d} P_q(x^{t-1}) dw(q)} dw(q)$$

Important Example: FSM model classes

q Given FSM S with k states (**fixed** initial state)

parameters = conditional probabilities per state $p(x/s), x \in A, s \in S$

ML estimates:
$$\hat{p}_{x^n}(x|s) = \frac{n_{x^n}(x|s)}{n_{x^n}(s)}$$

target:
$$\min_{q \in \Theta_{k(a-1)}} \log \frac{1}{P_q(x^n)} = \log \frac{1}{P_{\hat{q}(x^n)}(x^n)} = nH(x^n | S)$$

mixture distribution:

$$Q_w(x^n) = \int_{q \in \Theta_{k(a-1)}} \prod_{s \in S, x \in A} p(x|s)^{n_{x^n}(x|s)} dw(q) = \prod_{s \in S, x \in A} \int_0^1 q^{n_{x^n}(x|s)} dw(q)$$

⇒ equivalent to doing i.i.d. mixture for every state sub-sequence

⇒ optimal universal probability assignment:

$$q_{x^t}(x|s) = \frac{n_{x^t}(x|s) + 1/2}{n_{x^t}(s) + a/2}$$

Krichevski-Trofimov (KT) estimator \rightarrow

FSM model: Numerical Example

q FSM = First order Markov, initial state 0

$$x^{20} = \mathbf{0\ 0\ 1\ 1\ 1\ 1\ 0\ 1\ 1\ 1\ 1\ 1\ 1\ 0\ 0\ 0\ 0\ 1\ 0\ 0}$$

$$\frac{1}{2} \ \frac{3}{4} \ \frac{1}{6} \ \frac{1}{2} \ \frac{3}{4} \ \frac{5}{6} \ \frac{1}{8} \ \frac{3}{8} \ \frac{7}{10} \ \frac{3}{4} \ \frac{11}{14} \ \frac{13}{16} \ \frac{5}{6} \ \frac{3}{20} \ \frac{1}{2} \ \frac{7}{12} \ \frac{9}{14} \ \frac{5}{16} \ \frac{5}{22} \ \frac{13}{18}$$

total probability:

$$Q(x^n) = \prod_{s \in S} \frac{\Gamma(n(0|s) + \frac{1}{2})\Gamma(n(1|s) + \frac{1}{2})}{n(s)!\Gamma(\frac{1}{2})^2} = \frac{\Gamma(6.5)\Gamma(3.5)\Gamma(3.5)\Gamma(8.5)}{9!1!\Gamma(\frac{1}{2})^4}$$

$$-\log Q(x^{20}) \approx 17.61 \quad \Rightarrow \quad \text{code length: } \mathbf{18 \text{ bits}}$$

$$n\hat{H}(x^{20} | S) = 9h(1/3) + 11h(3/11) \approx 17.56 \quad \mathbf{\text{bits}}$$

FSM model classes and the Lempel-Ziv algorithm

- q The LZ algorithm is universal for **ANY** class of FSM models (of any size) \Rightarrow in particular, for all Markov models

$$\limsup_{n \rightarrow \infty} \frac{1}{n} L_{LZ}(x^n) \leq \limsup_{n \rightarrow \infty} H(x^n | S_k)$$

Markov machine
of any order k

and for an infinite sequence x^∞

$$\limsup_{n \rightarrow \infty} \frac{1}{n} L_{LZ}(x^n) \leq \lim_{k \rightarrow \infty} \limsup_{n \rightarrow \infty} H(x^n | S_k) = H(x^\infty)$$

limit exists

Markov
compressibility

Expected Pointwise Redundancy

- q So far, we have considered the code length for **individual sequences**, without taking any average
- q We have found that for parametric classes, the best we can do is to achieve a **worst-case** pointwise redundancy of about $(d/2)\log n$ bits
- q This means guaranteed performance, but maybe there are only a few such “unlucky” strings?
- q How to weight the redundancies in order to compute an average?
Since we are assuming that the class $\mathcal{C} = \{P_\theta, \theta \in \Theta_d\}$ is a good model for the data, it makes sense to assume that the data was drawn from a source with some distribution P_θ in \mathcal{C}
- q **Expected pointwise redundancy** of a code

$$E_q [R_C(L, x^n)] = \frac{1}{n} E_q [L(x^n) + \log P_{q(x^n)}(x^n)]$$

Expected Redundancy

- q Maybe there exist codes that are “good” (with expected pointwise redundancy smaller than $(d/2)\log n$) for a significant fraction of models in the class?
- q We will be even “larger”: we consider the **expected redundancy**

$$\bar{R}_C(L, q) = \frac{1}{n} E_q [L(x^n)] - H_n(q) = D_n(P_q \parallel Q)$$

normalized divergence ↑ ↑ $Q(x^n) = 2^{-L(x^n)}$
 for distr. defined on n -tuples

This is more “liberal” because

$$H_n(q) = \frac{1}{n} E_q [-\log P_q(x^n)] \geq \frac{1}{n} E_q [-\log P_{q(x^n)}(x^n)]$$

and therefore

$$\bar{R}_C(L, q) \leq E_q [R_C(L, x^n)]$$

Lower Bound on Expected Redundancy

q The pointwise universal codes we saw are - a fortiori - average universal for all parameters $\theta \Rightarrow$ this means that the corresponding distribution Q is “close” to **all the models in the class** in the sense that $D_n(P_q \parallel Q) \rightarrow 0$

q Still, the answer is that $(d/2)\log n$ is about the best we can do for **most** models in the class: it's the **inevitable** cost of universality

q Theorem:

Assume that either CLT holds for ML estimator of parameters in Θ_d or $\Pr\{\sqrt{n}(\hat{q}_i(x^n) - q_i) \geq \log n\} \leq d(n) \rightarrow 0$

Then for all Q and all $\varepsilon > 0$,

$$-n^{-1} E_q [\log Q(x^n)] \geq H_n(q) + k \frac{\log n}{2n} (1 - \varepsilon)$$

for all points θ in Θ_d except in a set whose volume $\rightarrow 0$ as $n \rightarrow \infty$

Lower Bound (cont.)

- q This lower bound parallels Shannon's coding theorem: when we consider a model class instead of a single distribution, a **model cost** gets added to the entropy
- q The bound cannot hold for **all** models in the family, but it holds for **most**
- q One interpretation of the lower bound: if the parameters can be estimated well, they are “distinguishable” (P_θ is sensitive to θ), so the class cannot be coded without a model cost

Conclusion:

the number of parameters affects the achievable convergence rate of a universal code length to the entropy

Variations on the Lower Bound

q Assuming, in addition, that $\sum_n d(n) < \infty$, then the **cumulative volume** of “bad” parameters for **all** sufficiently large $n \rightarrow 0 \Rightarrow$ the “bad” parameters form a set of Lebesgue volume 0

q **Strong Redundancy-Capacity Theorem:**
Under very mild conditions on the model class,

$$-n^{-1} E_q [\log Q(x^n)] \geq H_n(q) + (1 - \epsilon) C_n$$

where $C_n = \sup_w \left[H_n(Q_w) - \int_{q \in \Theta_d} H_n(q) dw(q) \right] \geq 0$ and Q_w is a mixture,

for all θ except for a set B for which $w^*(B) \leq e 2^{-nC_n \epsilon}$, where w^* is the density that achieves the supremum

q For parametric classes, C_n indeed behaves as $(d/2n) \log n$ and $w^*(B) \rightarrow 0$

Optimal Codes for Average Redundancy

q How do we achieve $\inf_L \sup_q \{E_q [L(x^n)] - H_n(q)\}$?

We use the mixture Q_w for which the **weighted** expected redundancy

$$\int_{q \in \Theta_d} [E_q [-\log Q_w(x^n)] - H_n(q)] dw(q) \text{ is maximum } \Rightarrow C_n$$

\Rightarrow an appropriate mixture is the best one can do to minimize the expected redundancy for the worst-case parameter (close to Jeffrey's prior in the case of exponential families/Markov models)

q In many situations, the minimum expected redundancy (or at least its main term $(d/2)\log n$) can be achieved by a **plug-in** code of the form

$$Q(x^n) = \prod_{t=1}^n P_{\tilde{q}_{t-1}}(x_t)$$

where \tilde{q}_{t-1} is an estimate of θ based on x^{t-1}

- | example: Bernoulli/FSM cases (KT estimator; even Laplace's estimator works in the average sense)

Summary of Codes and Bounds

q Minimax pointwise redundancy: NML code

- | for “nice” parametric families with d parameters,

$$R_c = \frac{d}{2n} \log \frac{n}{2p} + \frac{1}{n} \log \int_{\Theta_d} \sqrt{|I(q)|} dq + o(1/n)$$

- | horizon-dependent, mixture code with suitable prior is a good horizon-free approximation
- | in fact, for i.i.d. models with $\alpha-1$ parameters, it can be shown that

$$\lim_{n \rightarrow \infty} q_{\text{NML}}(x_t | x_1 x_2 \dots x_{t-1}) = \lim_{n \rightarrow \infty} \frac{\sum_{u \in A^{n-t}} Q_{\text{NML}}(x^t u)}{\sum_{v \in A^{n-t+1}} Q_{\text{NML}}(x^{t-1} v)} = \frac{n_{x^t}(x | s) + 1/2}{n_{x^t}(s) + a/2}$$



KT mixture!

Summary (cont.)

- q **Minimax average redundancy: mixture codes with prior such that the weighted expected redundancy is maximum**

$$C_n = \sup_w \left[H_n(Q_w) - \int_{q \in \Theta_d} H_n(q) dw(q) \right]$$

- | for i.i.d. exponential families and Markov models, close to Jeffrey's prior

$$C_n = \frac{d}{2n} \log \frac{n}{2pe} + \frac{1}{n} \log \int_{\Theta_d} \sqrt{|I(q)|} dq + o(1/n)$$

- | not only we cannot do better than C_n for one parameter (minimax), but for **most** parameters

- | for “nice” parametric families,

$$C_n \approx \frac{d}{2} \frac{\log n}{n}$$

Twice-universal coding

q Consider a model class to be the **union** of **nested** classes of growing dimensionality: $\Theta_1 \cup \Theta_2 \cup \dots \cup \Theta_d \cup \dots$

where $\Theta_1 \subset \Theta_2 \subset \dots \subset \Theta_d \subset \dots$

| example: Markov of different order, or FSM

q We seek a code length close to the minimum over the classes of the universal code length for each class:

$$\frac{1}{n} L(x^n) \approx \min_d \min_{q \in \Theta_d} \left\{ H_n(q) + \frac{d}{2n} \log n \right\}$$

| here, we also try to optimize the **model size**

| trade-off: first part diminishes with d , second part grows with d

| we answer questions such as: should we model the data as i.i.d. or as Markov of order 1?

| since the second level of universality is over a discrete set, we expect it not to affect the main redundancy term

Mixture Approach to Double-Universality

q Consider the series $1 = \sum_i l_i$ where $l_i > 0$

The probability assignment $Q(x^n) = \sum_d l_d Q_d(x^n)$

universal sequential code for Θ_d

is sequential: $Q(x^n) = \sum_{x_{n+1}} Q(x^{n+1})$

It is universal as desired: for all d

$$-\frac{1}{n} \log Q(x^n) \leq -\frac{1}{n} \log Q_d(x^n) - \frac{\log l_d}{n} < \min_{q \in \Theta_d} \log \frac{1}{P_q(x^n)} + \frac{d}{2n} \log n + O(1/n)$$

q However, it involves an infinite sum

- | this can be avoided by modifying Q_d so that $Q_d(x_t) = a^{-1}$ for all $t \leq d$: this way, all models with $d \geq n$ assign the same probability a^{-n} and the sum becomes finite, without affecting universality and sequentiality
- | anyway, this is obviously not practical but shows achievability of the goal

Plug-in Approach to Double-Universality

- q The idea: estimate the best model class Θ_d based on x^{t-1} , and plug it to encode x_t with Q_d

$$Q(x^n) = \prod_{t=1}^n Q_{d(x^{t-1})}(x_t | x^{t-1})$$

- | the decoder can reproduce the process without overhead

- q In general, this approach does not work pointwise, but for some model classes it was shown to work **on the average**

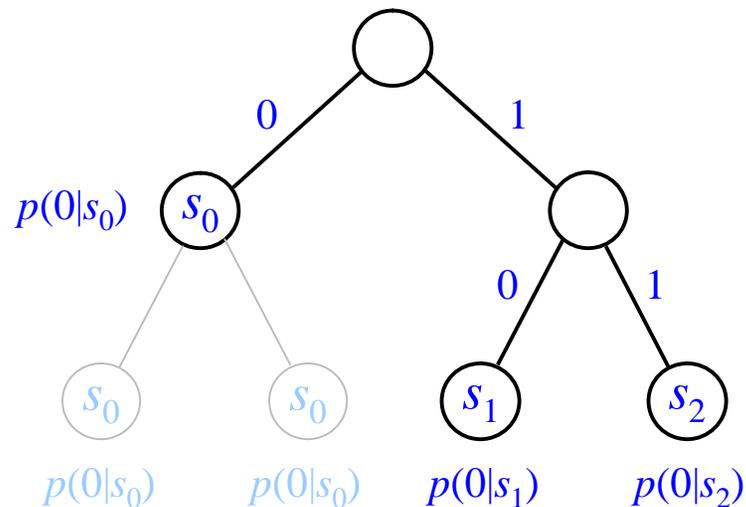
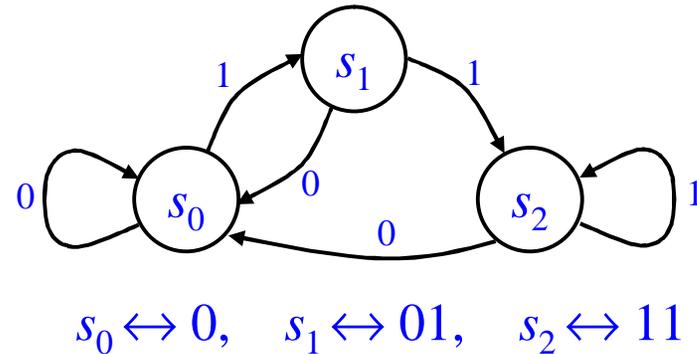
Lossless Source Coding

5. Universal Algorithms for Tree Models

Double-Universality for Tree Models

q Tree sources (FSMX)

- | finite memory $\leq k$ (Markov)
- | # of past symbols needed to determine the state might be $< k$ for some states



- | by merging nodes from the full Markov tree, we get a model with a **smaller number of free parameters**
- | the set of tree sources with unbalanced trees has **measure zero** in the space of Markov sources of any given order: **otherwise, double-universality would contradict the lower bound!!**
- | yet, tree models have proven very useful in practice, partly because there exist efficient **twice-universal, sequential** schemes in the class of tree models of **any size**

Two-pass Context Algorithm

q First pass: gather all the context statistics for x^n in a tree

| $O(n)$ complexity with suffix tree methods

q After the first pass, associate to each node a cost $L(s) + \frac{a}{a-1}$

KT code length for symbols occurring at s : $\frac{a}{a-1}$
 includes **penalty** for overparametrization

and “prune” the tree to find a subtree T with total minimum cost for the leaves (dynamic programming algorithm)

q The total cost for T is:

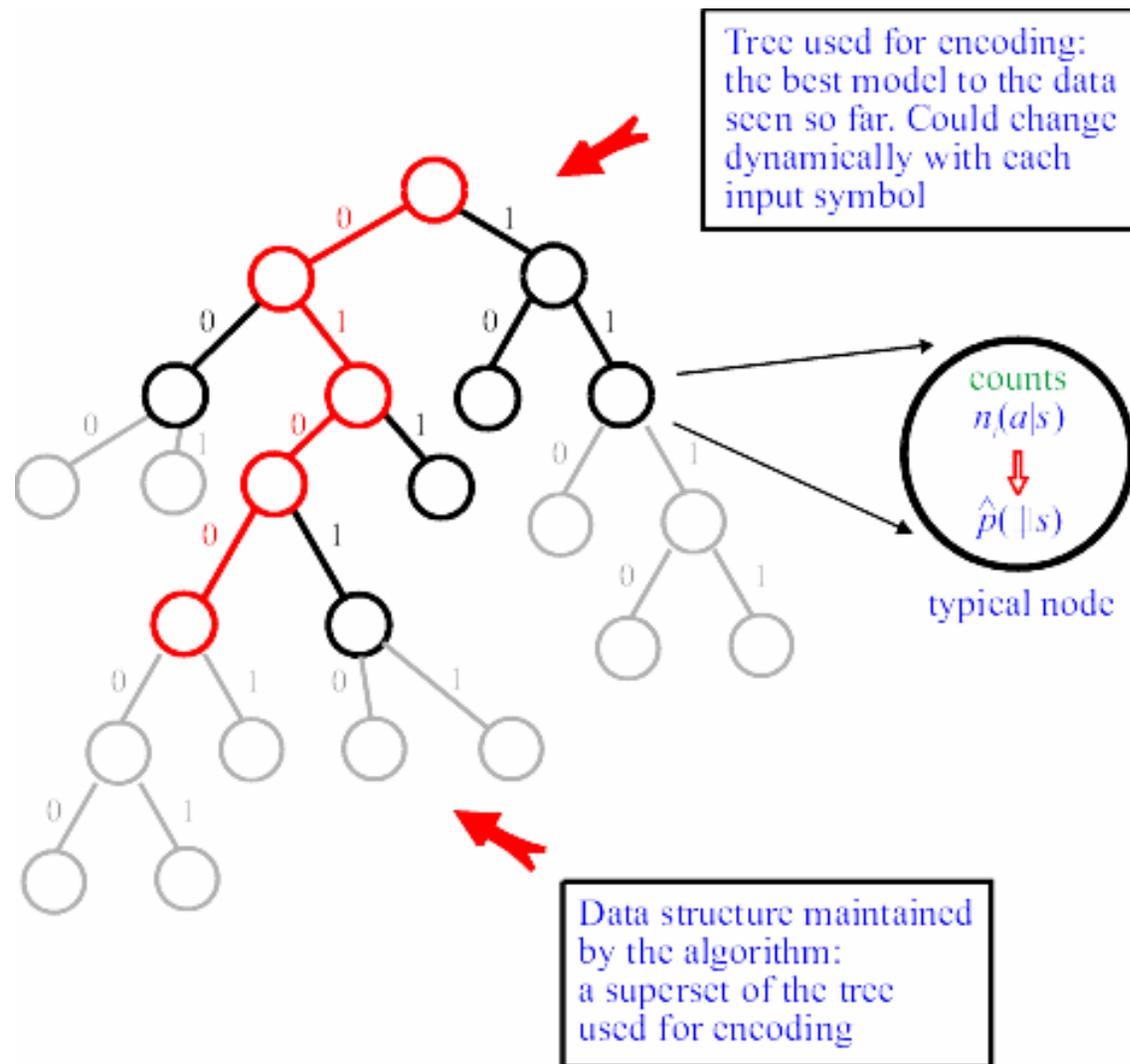
$$L_T(x^n) + \frac{\overset{\text{\# leaves in } T}{l_T} a}{a-1} = L_T(x^n) + \overset{\text{\# nodes in } T}{n_T} + \frac{1}{a-1}$$

Two-pass Context Algorithm (cont.)

q **Second pass: describe T to the decoder using n_T bits and encode x^n conditioned on T with KT \Rightarrow by definition, we have the best tree**

$$\frac{1}{n} L(x^n) = \frac{1}{n} L_T(x^n) + \frac{n_T}{n} < H(x^n | T) + \frac{(a-1)n_T}{2n} \log n + O(1/n)$$

Pruning of Context Trees: Example



Drawbacks of Two-pass Approach

q Non-sequential

q The scheme has asymptotically minimum **pointwise** redundancy, but is redundant (given the tree, not all x^n are possible)

$$Q(x^n) = 2^{-L(x^n)} = 2^{-L_T(x^n)} 2^{-n_T} = 2^{-L_T(x^n)} \prod_T < \sum_T \prod_T 2^{-L_T(x^n)}$$

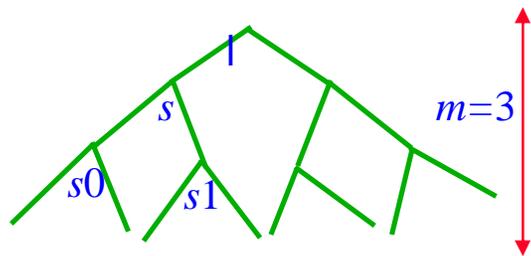
⇒ mixture with weights λ_T is better!

q Can the mixture be efficiently done?

Mixture Approach: Context Tree Weighting

- q **CTW**: An efficient implementation of the mixture for the class of binary tree models of maximum length bounded by m
 - | there exist extensions for non-binary and unbounded, but not as elegant

- q Probability assignment associated with a node s of the tree



$$Q_t^w(s) = \begin{cases} \frac{Q_t^{\text{KT}}(s) + Q_t^w(s0)Q_t^w(s1)}{2} \\ Q_t^{\text{KT}}(s) & \text{if } \text{depth}(s) = m \end{cases}$$

$$q_{\text{CTW}}(x_t | x_1 x_2 \dots x_{t-1}) = \frac{Q_{\text{CTW}}(x^t)}{Q_{\text{CTW}}(x^{t-1})} = \frac{Q_t^w(l)}{Q_{t-1}^w(l)}$$

- q **Theorem**: Let T have l_T leaves ($\Rightarrow n_T = 2^{l_T} - 1$), m_T leaves at level m

$$Q_{\text{CTW}}(x^t) = \sum_T 2^{-2l_T + 1 + m_T} Q_T^{\text{KT}}(x^t) \quad \Rightarrow \quad L_{\text{CTW}}(x^n) \leq L_T^{\text{KT}}(x^t) + n_T$$

↑
↑

KT prob. assignment
conditioned on T
for **any** T

CTW: "Proof" by Example

q Let $m = 2$

$$Q_t^w(s) = \begin{cases} \frac{Q_t^{\text{KT}}(s) + Q_t^w(s0)Q_t^w(s1)}{2} \\ Q_t^{\text{KT}}(s) & \text{if depth}(s) = m \end{cases}$$

$$\begin{aligned}
 Q^w(1) &= \frac{Q^{\text{KT}}(1) + Q^w(0)Q^w(1)}{2} \\
 &= \frac{Q^{\text{KT}}(1)}{2} + \frac{[Q^{\text{KT}}(0) + Q^w(00)Q^w(01)][Q^{\text{KT}}(1) + Q^w(10)Q^w(11)]}{8} \\
 &= \frac{Q^{\text{KT}}(1)}{2} + \frac{[Q^{\text{KT}}(0) + Q^{\text{KT}}(00)Q^{\text{KT}}(01)][Q^{\text{KT}}(1) + Q^{\text{KT}}(10)Q^{\text{KT}}(11)]}{8} \\
 &= \frac{Q^{\text{KT}}(1)}{2} + \frac{Q^{\text{KT}}(0)Q^{\text{KT}}(1)}{8} + \frac{Q^{\text{KT}}(0)Q^{\text{KT}}(10)Q^{\text{KT}}(11)}{8} \left. \vphantom{\frac{Q^{\text{KT}}(0)Q^{\text{KT}}(1)}{8}} \right\} \\
 &\quad \text{root} \quad \wedge \text{P } n_T - m_T = 3 \quad \wedge \text{P } n_T - m_T = 3 \\
 &+ \frac{Q^{\text{KT}}(1)Q^{\text{KT}}(01)Q^{\text{KT}}(00)}{8} + \frac{Q^{\text{KT}}(00)Q^{\text{KT}}(01)Q^{\text{KT}}(10)Q^{\text{KT}}(11)}{8} \left. \vphantom{\frac{Q^{\text{KT}}(00)Q^{\text{KT}}(01)Q^{\text{KT}}(10)Q^{\text{KT}}(11)}{8}} \right\} \\
 &\quad \wedge \text{P } n_T - m_T = 3 \quad \wedge \text{P } n_T - m_T = 3
 \end{aligned}$$

CTW: Practical Issues

- q **Finite precision:** - as node probabilities get smaller they cannot be stored \Rightarrow may need to re-compute from counts
 - probability assignment is a ratio of vanishing numbers
- q In general, m needs to be assumed large, and complexity raises
- q **Binary only**
- q **However, there are good implementations of CTW which provide the best compression ratios on text and binary file compression**
- q **From a practical viewpoint, a plug-in approach is more appealing**
 - | however, universality is only shown on an average sense

Plug-in Approach

- q In a plug-in approach, the idea is to encode x_{t+1} by selecting a context from x^t that strikes the “right balance” between **conditional entropy** and **model cost**
 - | a long context reduces entropy by introducing more conditioning
 - | but the longer the context, the smaller occurrence counts it has and the statistics that we learn from them are less reliable
 - | for each time t , it is only necessary to select a context in the path $x_t x_{t-1} x_{t-2} \dots$, not the whole tree
- q A plug-in twice-universal scheme consists of a **context selection rule**, and a coding scheme based on the statistics stored at the selected context (e.g., KT estimator)
- q The context selection rule is also a tool for model selection for purposes other than data compression

Context Algorithm

- q The algorithm consists of three interleaved stages:
 - | growing of a tree that captures, essentially, all occurrences of each symbol at each node
 - | a context selection rule that selects for x_{t+1} a context $s_t(x^t)$ from the tree T_t grown by x^t
 - | a KT sequential probability assignment for x_{t+1} based on the counts stored at $s_t(x^t)$
- q For each new symbol, we first encode, then update $T_t \rightarrow T_{t+1}$ (think of the decoder!)
 - | the update consists of incrementing the occurrence counts for all the nodes in the path $x_t x_{t-1} x_{t-2} \dots$
 - | when we get to a leaf that was already visited, we extend the tree one level and initialize the count of x_{t+1} to **1** (the others remain **0**) \Rightarrow only a few initial occurrences are missing, so we can basically assume that for all nodes s and all symbols a , the count $n_{x^t}(a | s)$ is available
 - | nodes in other paths don't need to be visited

Context Selection Rule

q **Basic principle:**

The node which would have assigned the shortest code length for its symbol occurrences in the past string should be selected

q **Most intuitive choice:** find minimum cost tree T_t for all times t , and do KT-coding for x_{t+1} conditioned on the context $x_t x_{t-1} x_{t-2} \dots$ in T_t

- | unlike two-pass, the cost should **not** include tree description: the decoder already has it !
- | very complex and was not shown to be asymptotically optimal, even on the average

q **Another possibility:** for each node $sb \in T_t$, $b \in A$ define

$$\Delta_t(sb) = L_t(sb) - L'_t(s)$$

KT code length for symbols occurring at sb based on counts gathered at s

KT code length for symbols occurring at sb

Choose the deepest node sb in T_t such that $\Delta_t(sb) < 0$

Universality of Context Algorithm

q **Easier to analyze:**

$$\Delta_t(sb) = \sum_{a \in A} n_{x^t}(a|sb) \log \frac{\hat{P}_{x^t}(a|sb)}{P_{x^t}(a|s)}$$

the selected tree T_t is the smallest complete super-tree of
 { **the deepest nodes** w in T_t **s.t.** $\Delta_t(w) \geq C \log(t+1)$ }

q **Theorem:** For any minimal complete tree T with k leaves defining a tree source $P_T(x^n)$ with probabilities bounded away from 0, if $C > 2(\alpha+1)$ then the probability assignment Q of Context Algorithm satisfies

$$\frac{1}{n} E_T \left[\log \frac{P_T(x^n)}{Q(x^n)} \right] \leq k(\alpha - 1) \frac{\log n}{2n} + O(1/n)$$

and, moreover,

$$\frac{1}{n} \log \frac{P_T(x^n)}{Q(x^n)} \leq k(\alpha - 1) \frac{\log n}{2n} + O(1/n)$$

with P_T -probability 1

Idea of the Proof

q The proof is based on the fact that

$$\sum_{t=1}^{\infty} P_T \{x^t \mid \mathbb{T}_t \neq T\} \log t < \infty$$

and so the contribution of the “bad” sequences is $O(1/n)$

q Two (non-disjoint) classes of errors:

- | **overestimation**: the selected tree contains an internal node which is a leaf of the true tree; taken care of by the penalty term
- | **underestimation**: a leaf of the selected tree is an internal node of the true tree; requires large deviation techniques

PPM algorithm

- q A context selection rule very popular in the data compression community: choose the shortest context that occurred more than a certain number of times
 - | rationale: the context gathered enough statistics in order to be reliable
 - | the rule is totally *ad hoc*: if the best tree model is short, it will tend to overestimate (think of data generated by an i.i.d. source!)
- q A family of algorithms based on variations of this selection rule is called PPM (**Prediction by Partial Matching**)
 - | it is a very popular scheme for text compression and yields some of the best compression ratios
 - | however, algorithms with a stronger theoretical basis tend to do better

Application to Statistics: the MDL Principle

- q Model selection is probably the most important problem in statistical inference
- q Minimum Description Length (MDL) principle of statistical inference:
choose the model class that provides the shortest code length for the model and the data in terms of the model \Rightarrow universal coding theory provides the yardstick to measure this code length
- q Rationale: models serve as tools to describe regularities in the data
 - | we should use the simplest explanation for the data, but not too simple
- q Strong consistency shown in various settings, solves problem of model order selection avoiding the use of artificial penalty terms
- q **Bayesian** interpretation through mixture codes
 - | however, NML code cannot be explained as a mixture
- q Other interpretations: maximum-entropy principle

Lossless Source Coding

6. Sequential Decision Problems

The Sequential Decision Problem

q The framework

- | **observations:** $x^n = x_1 x_2 \dots x_n, x_t \in A$
- | **corresponding actions:** $b^n = b_1 b_2 \dots b_n, b_t \in B$
- | **instantaneous losses** $l(b_t, x_t)$ **accumulate over time:**

$$L(x^n) = \sum_{t=1}^n l(b_t, x_t)$$

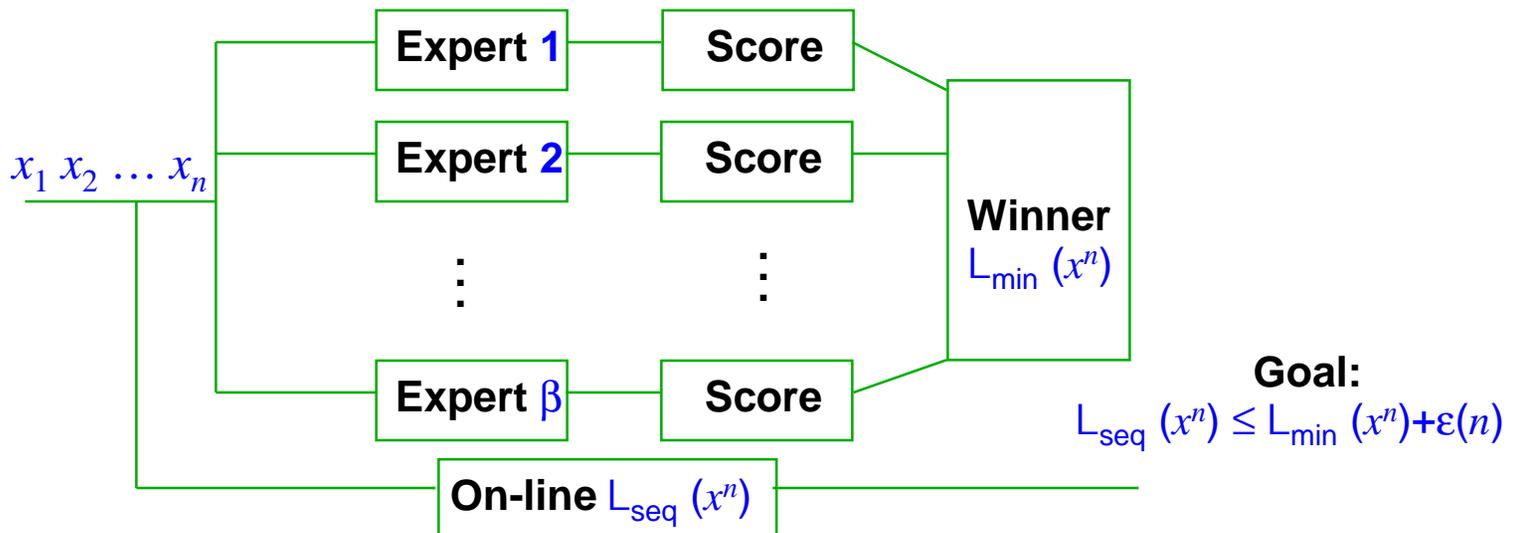
q On-line (sequential) strategy

- | $\{b_t\}$, action b_t is decided **before** observing x_t
- | possibly **randomized:** $\{ p_t(b_t | x^{t-1}, b^{t-1}) \}$

q The goal: as $n \rightarrow \infty$, approximate performance of best strategy in a given **reference class**, for **arbitrary** x^n (**deterministic** setting)

- | **excess loss w.r.t. reference:** **regret** or **redundancy**
- | the reference class (or **expert set**) may reflect limited resources

Prediction with expert advice



- q Most general setting: **reference class = set of generic “experts”**
- q Basic principle for on-line strategy

select an expert’s prediction randomly, with probability dependent on its accumulated loss

Sequential Decision Problem: Examples

q Binary prediction with Hamming loss

- | $x_1 x_2 \dots x_n$ is a binary sequence ($|A|=2$)
- | the action: predict either $b_t = 0$ or 1 (deterministic), or assign a probability p_t to 1 (randomized strategy) ($|B|=2$)
- | the loss:
 - deterministic strategy: $l(b_t, x_t) = 0$ if $b_t = x_t$ and 1 otherwise
 - ▢ accumulated loss = total # of prediction errors
 - randomized strategy: $\mathbf{E}[l(b_t, x_t)] = |x_t - p_t|$

q A less trivial example: lossless data compression

- | $x_1 x_2 \dots x_n$ is the data to encode, finite alphabet A
- | the (**deterministic**) action b_t is a **probability distribution** assigned to x_t
 $b_t = \{p_t(\cdot | x_1 x_2 \dots x_{t-1})\}$ (B continuous and vectorial!)
- | the loss: $l(b_t, x_t) = -\log p_t(x_t | x_1 x_2 \dots x_{t-1})$ (given the assigned distribution, an encoder can generate a code word of length $l(b_t, x_t)$)
- | the accumulated loss (total code length) is $-\log$ of the probability assigned to $x_1 x_2 \dots x_n$

Exponential-Weighting Algorithm

q The most general scheme for on-line expert selection

| here, we will assume A, B finite and loss bounded by l_{\max}

q $L_F(x^t) =$ loss of expert $F \in F$ accumulated through time t

At time $t+1$ choose the action suggested by expert F with probability

$$P_{t+1}(F) = \frac{e^{-hL_F(x^t)}}{\sum_{F' \in F} e^{-hL_{F'}(x^t)}}$$


 some given positive constant
 number of experts

Then,

$$L_{\text{ew}}(x^n) \leq \min_{F \in F} L_F(x^n) + \frac{\ln b}{h} + \frac{nhl_{\max}^2}{8}$$

▮ with $h = \frac{1}{l_{\max}} \sqrt{(8 \ln b)/n}$, normalized excess loss over best expert is

$$l_{\max} \sqrt{(\ln b)/(2n)} \quad (\text{horizon-dependent scheme})$$

Horizon-Free Exponential Weighting

q **Horizon-free variant:** divide the data into blocks of exponentially growing size, and apply the horizon-dependent algorithm to each block \triangleright it is easy to see that the overall loss increases only by a constant factor for all n

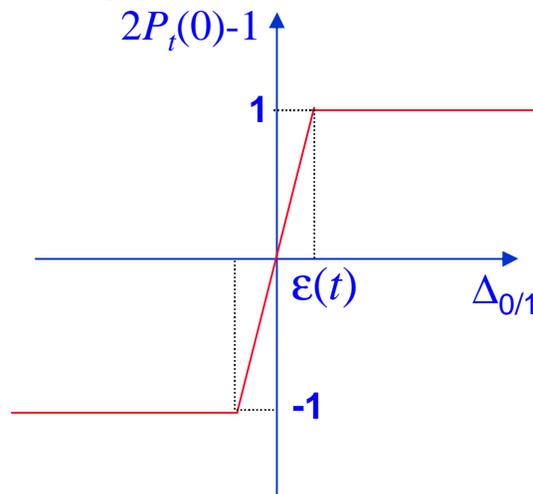
Example: Binary Prediction with Constant Experts

- Two experts: one always says “0”, the other always says “1”
 - analogous to memoryless model in data compression

$L_{\min}(x^n) = \min(n_0, n_1)$

$$P_{t+1}(0) = \frac{1}{1 + e^{-h[n_0(x^t) - n_1(x^t)]}}$$

Horizon-free approximation:



Choose the current winner except if up to $\varepsilon(t)$ from tie, with $\varepsilon(t) = O(\sqrt{1/n})$, in which case randomize \square
similar redundancy, different constant

- Minimum worst-case redundancy: draw $x_{t+1} x_{t+2} \dots x_n$ at random and predict the winner in the overall sequence of length $n-1$

FS reference strategies (L-Z framework)

q A **given** FSM is driven by the observations $\{x_t\}$

S = set of states f = next-state function

s_1 = fixed initial state $s_{t+1} = f(s_t, x_t)$

Reference strategy is allowed to vary following the FSM: $b_t = g(s_t)$

| example: one-state machine = constant strategy

q Best g for the specific x^n :

$$g(s) = \arg \min_{b \in B} \mathbb{E}_x [l(b, x) | s]$$

expectation w.r.t.  conditional empirical distribution

▷ normalized regret vanishes by applying single- state strategy to sub-sequences at each state

q Take best FSM for x^n , consider normalized loss for $n \rightarrow \infty$, and $|S| \rightarrow \infty$

q For log loss: FS “**compressibility**” of an infinite individual sequence

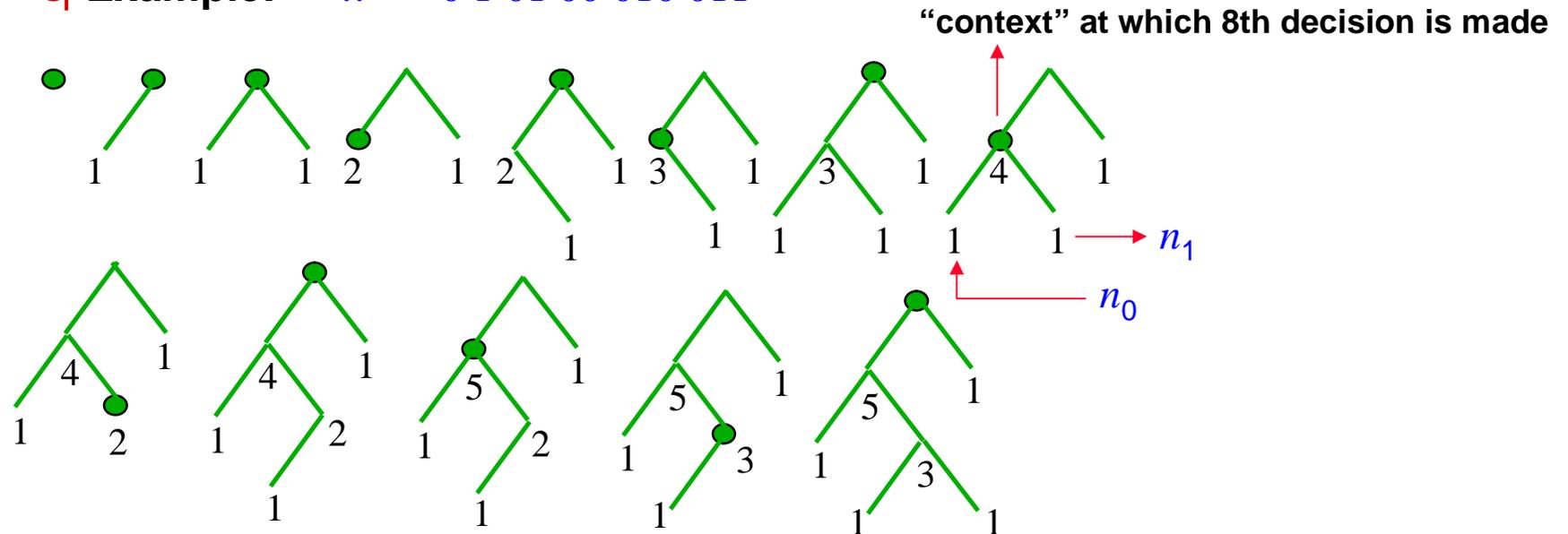
| **efficiently** achievable: LZ data compression scheme

FS Predictability

q Similarly, FS “predictability”

l sequential “LZ-like” decision scheme performs essentially as well as the **best FSM (of any size!) for the sequence**

q Example: $x^{12} = 0\ 1\ 01\ 00\ 010\ 011$



Average Loss

q When the goal is to minimize the **average** number of prediction errors, the redundancy can be made much smaller than the worst pointwise case

| this is different from data compression (log loss)!

q For example, for binary prediction and two constant experts, “0” and “1”, the redundancy is given by $n^{-1}E_q[\text{\#errors}] - \min(\theta, 1-\theta)$ and it is $O(1/n)$ for a majority predictor **without randomization**

Proof: Assuming (without loss of generality) $p(1) = \theta < 0.5$,

$$P_t(\text{error}) = qP_q\{x_t = 0\} + (1-q)P_q\{x_t = 1\} = q + (1-2q)P_q\{x_t = 1\} \Rightarrow$$

$$E_q[\text{\#errors}] = \sum_{t=1}^n P_t(\text{error}) = nq + (1-2q) \sum_{t=1}^n P_q\{x_t = 1\} \Rightarrow$$

$$\frac{1}{n} E_q[\text{\#errors}] - q < \frac{(1-2q)}{n} \sum_{t=1}^{\infty} P_q\{n_{x^t}(1) \geq n_{x^t}(0)\} = O(1/n)$$

 vanishes exponentially fast (**Chernoff**)

Lossless Source Coding

7. Lossless image compression

Lossless Image Compression (the real thing...)



q Some applications of lossless image compression:

- | Images meant for further analysis and processing (as opposed to just human perception)

 - u Medical, space

- | Images where loss might have legal implications

 - u Medical

- | Images obtained at great cost

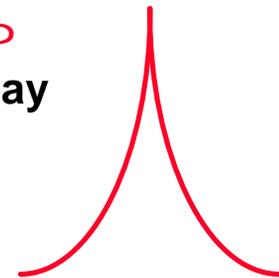
- | Applications with intensive editing and repeated compression/decompression cycles

- | Applications where desired quality of rendered image is unknown at time of acquisition

q A new international standard (ISO/IEC: “JPEG Committee): **JPEG-LS**

Universality vs. Prior Knowledge

- q Application of universal algorithms for tree models directly to real images yields poor results
 - | some structural symmetries typical of images are not captured by the model
 - | a universal model has an associated “learning cost:” why learn something we already know?
- q Modeling approach: limit model class by use of “**prior knowledge**”
 - | for example, images tend to be a combination of smooth regions and edges
 - | **predictive coding** was successfully used for years:
it encodes the difference between a pixel and a **predicted value** of it
 - | prediction errors tend to follow a Laplacian distribution \triangleright
AR model + Laplacian, where both the center and the decay are **context dependent**
 - | Prediction = **fixed** prediction + **adaptive** correction



Models for Images

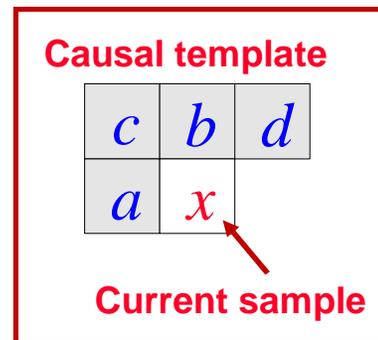
Continuous tone images

Gray scale: a 2D array of **pixel intensity values** (integers) in a given range $[0..(a-1)]$ (often $a=256$)

123	255	8	15	...
0	128	200	217	...
⋮				

Color: a 2D array of **vectors** (e.g. triplets) whose coordinates represent intensity in a given **color space** (e.g., RGB, YUV); similar principles

In practice, contexts are formed out of a finite subset of the past sequence

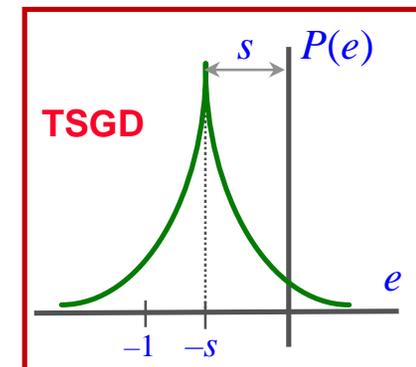


Conditional probability model for prediction errors: **two-sided geometric distribution (TSGD)**

$$P(e) = c_0 q^{|e+s|}, \quad q \in (0,1), \quad s \in [0,1)$$

“discrete Laplacian”

shift s constrained to $[0,1)$ by integer-valued **adaptive correction (bias cancellation)** on the fixed predictor



Complexity Constraints

- q Are sophisticated models worth the price in complexity?
 - | Algorithm Context and CTW are linear time algorithms for tree sources of limited depth, but quite expensive in practice
 - | even arithmetic coding is not something that a practitioner will easily buy in many applications!
- q Is high complexity required to approach the best possible compression?
- q The idea in JPEG-LS: apply judicious modeling to **reduce complexity**, rather than to improve compression

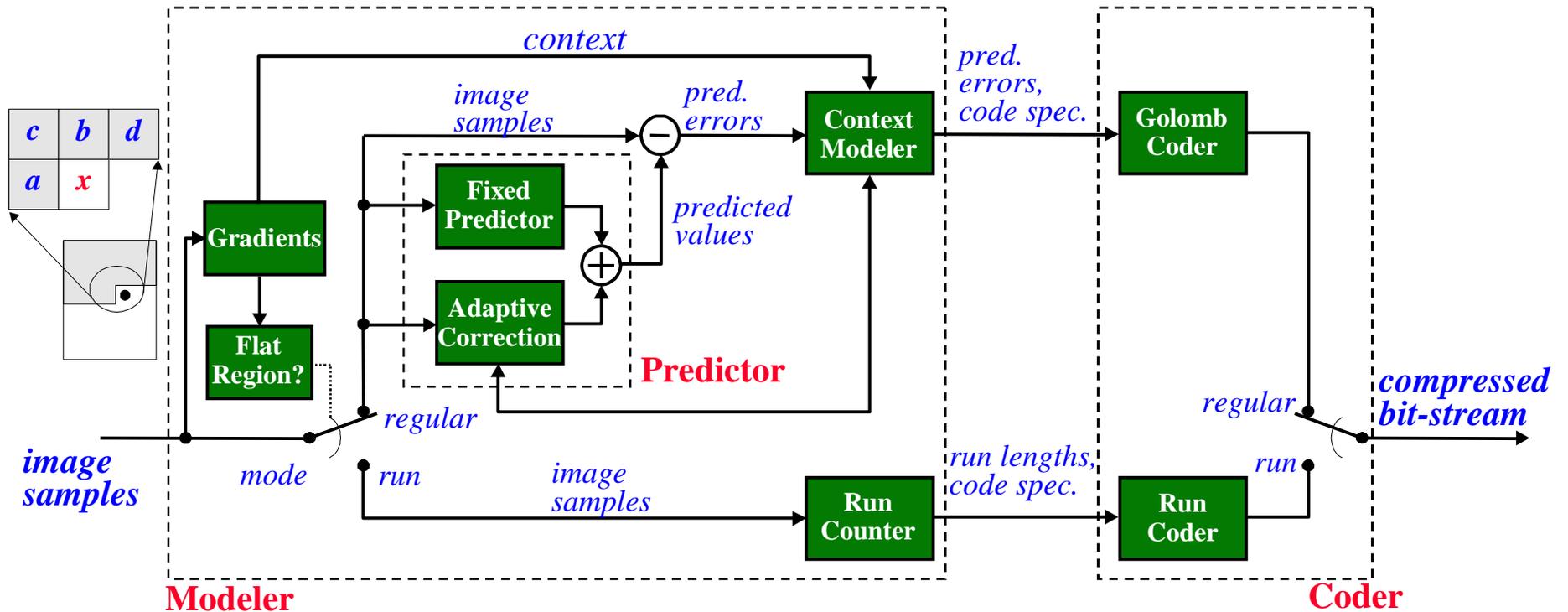
the modeling/coding separation paradigm is less neat without complex models or arithmetic coding

The LOCO-I algorithm

- q JPEG-LS is based on the **LOCO-I** algorithm:
Low COmplexity LOSSless COmpression of Images

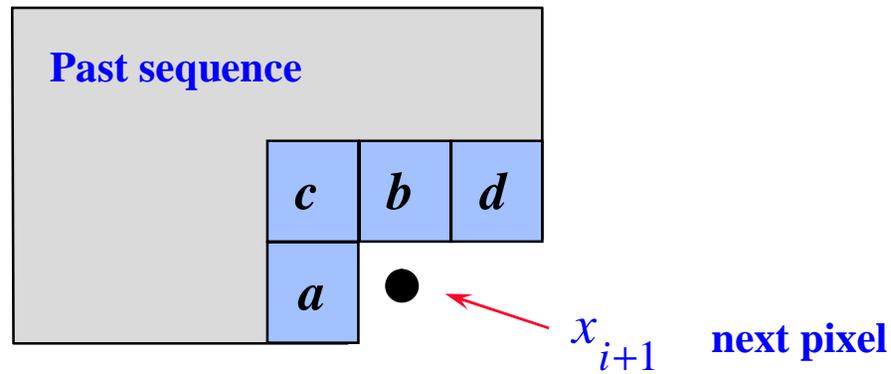
- q **Basic components:**
 - | **Fixed + Adaptive prediction**
 - | **Conditioning contexts based on quantized gradients**
 - | **Two-parameter conditional probability model (TSGD)**
 - | **Low complexity adaptive coding** matched to the model (variants of Golomb codes)
 - | **Run length coding** in flat areas to address drawback of symbol-by-symbol coding

JPEG-LS (LOCO-I Algorithm): Block Diagram



Fixed Predictor

q **Causal template** for prediction and conditioning



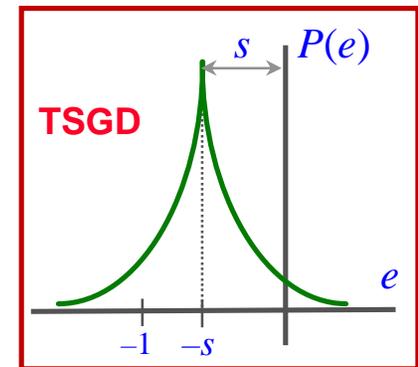
predicted value $\rightarrow \hat{x}_{i+1} = \left. \begin{cases} \min(a, b) & \text{if } c \geq \max(a, b) \\ \max(a, b) & \text{if } c \leq \min(a, b) \\ a + b - c & \text{otherwise} \end{cases} \right\} \text{median of } a, b, \text{ and } a+b-c$

q **Nonlinear**, has some “**edge detection**” capability:

- | Predicts b in “vertical edge”
- | Predicts a in “horizontal edge”
- | Predicts $a+b-c$ in “smooth region”

Parameter Reduction and Adaptivity

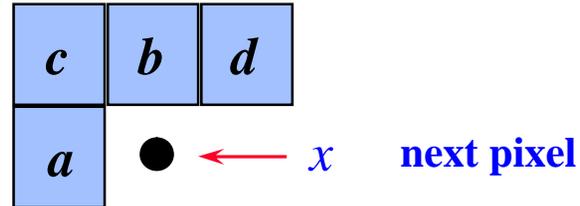
- q The goal in selecting the number of parameters: capture **high order dependencies** without excessive **model cost**
- q **Adaptive coding** is needed, but arithmetic coding ruled out (if possible...) due to complexity constraints
- q **A solution that addresses both issues:**
 - Model prediction residuals with a TSGD
 - | only **two parameters** per context
 - u “sharpness” (rate of decay, variance, etc.)
 - u shift (often non-zero in a context-dependent scheme)
 - | allows for large number of contexts (365 in JPEG-LS)
 - | suited to low complexity adaptive coding



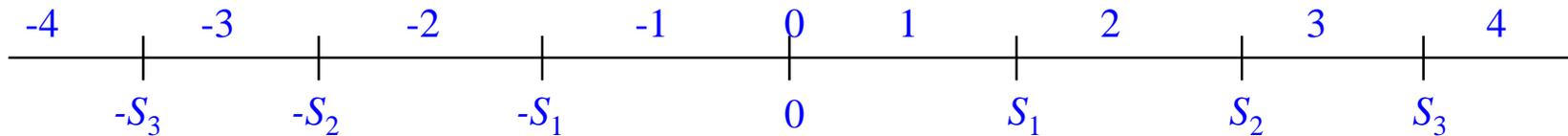
$$P(e) = c_0 q^{|e+s|}, \quad q \in (0,1), \quad s \in [0,1)$$

Context Determination

Causal template:



- q Look at the **gradients** $g_1 = (d-b)$, $g_2 = (b-c)$, $g_3 = (c-a)$,
 - | gradients capture the level of activity (smoothness, edginess) surrounding a pixel
 - | g_1, g_2, g_3 quantized into **9** regions determined by **3** thresholds S_1, S_2, S_3
 - | maximum information on x_{i+1} suggests equiprobable regions



- q **Symmetric contexts merged:**

$$P(e / [q_1, q_2, q_3]) \leftrightarrow P(-e / [-q_1, -q_2, -q_3])$$

- q **A fixed number of contexts: $(9^3 + 1)/2 = 365$**

Coding of TSGD's: Golomb Codes

q Optimal prefix codes for TSGDs are built out of the **Golomb codes** for nonnegative integers

q Given a **code parameter** m and an integer j ,

$$j \text{ } \textcircled{R} \text{ } \{ j \bmod m \text{ (in binary)}, \lfloor j/m \rfloor \text{ (in unary)} \}$$

| **example:** $j = 19, m = 4$: $19 \bmod 4 = 3$ $\lfloor 19/4 \rfloor = 4$

$$\underline{11} \quad \underline{00001}$$

| Golomb codes are optimal for **geometric** distributions

| JPEG-LS uses the subfamily of Golomb **power-of-2 (PO2)** codes: $m = 2^k$

$$G_{2^k} : n \rightarrow \left\{ \begin{array}{l} (n \bmod 2^k), \lfloor n / 2^k \rfloor \\ \text{(binary)} \qquad \text{(unary)} \end{array} \right\} \quad \begin{array}{l} n \geq 0, \\ k = \text{code parameter} \end{array}$$

| Encoding is very simple and **explicit**: no tables

| Very suited for **adaptive** coding: adapt k

Adaptive Coding of TSGD's in JPEG-LS

- q Optimal prefix codes for TSGD's are approximated in JPEG-LS by applying the Golomb-PO2 subfamily to a **mapped** error value:

$$\varepsilon \rightarrow M(\varepsilon) \text{ or } \varepsilon \rightarrow M(-1-\varepsilon)$$

$$0, -1, +1, -2, +2, \dots \text{ } \textcircled{R} \text{ } 0, 1, 2, 3, 4, \dots \text{ (Rice mapping), or}$$

$$-1, 0, -2, +1, -3, \dots \text{ } \textcircled{R} \text{ } 0, 1, 2, 3, 4, \dots$$

- q Adaptive code selection (parameter k , mapping)

- l approximation of optimal strategy based on **ML estimation** for TSGD parameters θ, s through **sufficient statistics**

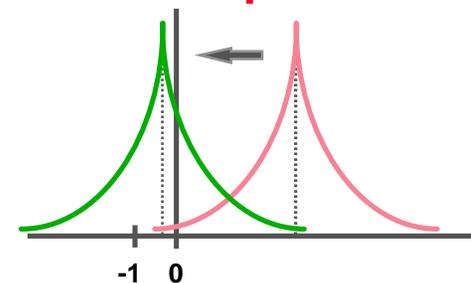
A = accumulated sum of **error magnitudes**

N_- = number of **negative** samples

- q Assumption $s \hat{\tau} [0, 1)$ satisfied through the use of **adaptive correction** of the predictor, using also

B = accumulated sum of **error values**

N = **total** number of samples



Adaptive Coding (cont.)

q For the Golomb-PO2 code with parameter k applied on remapped prediction errors

- | compute the expected code length explicitly as a function of θ , s , and k
- | replace (the **unknown**) θ and s by their ML estimates as a function of sufficient statistics

q Optimal decision regions for k result in

- | Let $\varphi = (\sqrt{5} + 1) / 2 \approx 1.618$ (golden ratio)
- | If $A - N_- \leq N \varphi$, use N_- to
 - u choose $k = 0$ or $k = 1$
 - u choose $s \geq 1/2$ or $s < 1/2$ if $k = 0$ (irrelevant if $k > 0$)
- | If $A - N_- > N \varphi$, choose k such that

$$\frac{1}{\varphi^{2^{-k+2}} - 1} < \frac{A - N_-}{N} \leq \frac{1}{\varphi^{2^{-k+1}} - 1} \quad k > 1$$

Approximation of Decision Regions

q **Observe** $j^{2^{-k}} - 1 \approx 2^{-k} \ln j$
 $\ln j \approx 0.48 \approx 0.5$

p $\frac{1}{j^{2^{-k}} - 1} + 0.5$ is always close to a power of 2

q **Since** $N_- / N \approx 0.5$, k can be approximated by $k \cong \lceil \log_2 (A/N) \rceil$



**Estimate k using the trivial loop
for ($k=0$; ($N \ll k$) < A ; $k++$);**

- | If $k=0$ and $s < -1/2$, encode $M(-\varepsilon-1)$, otherwise $M(\varepsilon)$
- | Only 4 variables per context are needed

Embedded Run-length Coding

- q Aimed at overcoming the basic limitation of 1 bit/pixel inherent to pixel-wise prefix codes, most damaging in low-entropy regions
- q Creates **super-symbols** representing **runs** of the same pixel value in the “flat region” $a = b = c = d$ \rightarrow special context
 $[q_1, q_2, q_3] = [0, 0, 0]$
- q A run of a is counted and the count is encoded using **block-MELCODE**, a fast adaptation technique for Golomb-type codes

LOCO-I in One Page

loop:

- q Get context pixels a, b, c, d , next pixel x
- q Compute gradients $d-b, b-c, c-a$ and quantize $\triangleright [q_1, q_2, q_3, sign]$
- q $[q_1, q_2, q_3] = 0$? YES: Process **run state** NO: Continue
- q $x_{pred} = predict(a, b, c)$
- q Update correction value for context. Correct x_{pred}
- q $e = x - x_{pred}$. If $sign < 0$ then $e = -e$
- q Estimate k for the context
- q Remap $e \text{ } \textcircled{R} \text{ } M(e)$ or $e \text{ } \textcircled{R} \text{ } M(-1-e)$
- q Encode M with Golomb-PO2(k)
- q Back to **loop**

run state:

- q Count run of a until $x \neq a \triangleright$ run length L
- q Encode L using block-MELCODE
- q Update MELCODE state

Compression/Complexity trade-off

