

Universal Ensembles for Sample-Wise Lossy Compression

Neri Merhav

The Viterbi Faculty of Electrical & Computer Engineering
Technion—Israel Institute of Technology
Haifa, Israel

ITA 2023, San Diego, CA, February 12–17, 2023

A Very Quick Overview

Universal lossless coding:

- Davisson ('73): maximin+minimax universality; mixtures.
- Rissanen ('84): a converse for 'most' parameter values.
- Weinberger, Merhav & Feder ('94): 'semi-deterministic' analogue.
- Merhav & Feder ('95): parametric class → general class.
- Many: extensions, improvements, relations to prediction, etc.

A Very Quick Overview (Cont'd)

Universal lossy d -semifaithful coding:

- Zhang, Yang & Wei ('97): non-universal redundancy $\geq \frac{\log n}{2n}$; achievable $\leq \frac{\log n}{n}$; universality - larger constant.
- Yu & Speed ('93): weak universality.
- Ornstein & Shields ('90): stat. erg. sources, Hamming distortion.
- Kontoyiannis ('00): a.s. results – CLT, LIL, no-cost universality.
- Kontoyiannis & Zhang ('02): $-\log \Pr\{D\text{-ball}\}$.
- Mahmood & Wagner ('22): minimax distortion-universality.

In This Work ...

We adopt the semi-deterministic paradigm of Weinberger, Merhav & Feder ('94) for lossy compression:

Redundancy rates relative to the 'memoryless' empirical RDF

- Random coding using a mixture (Kontoyiannis & Zhang - '02).
- Asympt. accurate evaluation of $\Pr\{D\text{-ball}\}$.
- Universality w.r.t. the distortion measure.
- Converse.

Sequences "with memory"

- Optimal length $= -\log P_{\text{LZ}}\{D\text{-ball}\}$
- The main contribution is in the converse.
- Discussion

Notation & Definitions

- Source sequence: $\mathbf{x} = (x_1, \dots, x_n) \in \mathcal{X}^n$, $|\mathcal{X}| = J$.
- Reproduction sequence: $\hat{\mathbf{x}} = (\hat{x}_1, \dots, \hat{x}_n) \in \hat{\mathcal{X}}^n$, $|\hat{\mathcal{X}}| = K$.
- Distortion measure: $d : \mathcal{X} \times \hat{\mathcal{X}} \rightarrow \mathbb{R}^+$; $d(\mathbf{x}, \hat{\mathbf{x}}) = \sum_i d(x_i, \hat{x}_i)$.
- Encoder: $\phi_n : \mathcal{X}^n \rightarrow \mathcal{G}_n \subset \{0, 1\}^*$.
- Decoder: $\psi_n : \mathcal{G}_n \rightarrow \mathcal{C}_n \subseteq \hat{\mathcal{X}}^n$.
- D -semifaithful code: $\forall \mathbf{x} \in \mathcal{X}^n$, $d(\mathbf{x}, \psi_n(\phi_n(\mathbf{x}))) \leq nD$.
- Code ensemble: independent random selection under

$$W(\hat{\mathbf{x}}) = (K - 1)! \cdot \int_{\mathcal{Q}} dQ \prod_{i=1}^n Q(\hat{x}_i).$$

- D -sphere: $\mathcal{S}(\mathbf{x}, D) = \{\hat{\mathbf{x}} : d(\mathbf{x}, \hat{\mathbf{x}}) \leq nD\}$.
- $\mathcal{T}_n(P) = \{\text{all } \mathbf{x} \in \mathcal{X}^n \text{ with empirical distribution } P\}$.

A Key Lemma - Assessing $W[\mathcal{S}(\mathbf{x}, D)]$

Let $\mathbf{x} \in \mathcal{T}_n(P)$ and define

$$F(s, Q) \triangleq - \sum_x P(x) \ln \left[\sum_{\hat{x}} Q(\hat{x}) e^{-sd(x, \hat{x})} \right] - sD.$$

Then, it is well known that

$$R_d(D, P) = \sup_{s \geq 0} \min_Q F(s, Q) = \min_Q \sup_{s \geq 0} F(s, Q).$$

Let (s^*, Q^*) be the saddle-point that achieves $R_d(D, P)$ and define

$$V(P, d) = \left| \det \left\{ \left. \text{Hess}F(s^* + j\omega, Q) \right|_{(0, Q^*)} \right\} \right|, \quad j = \sqrt{-1}.$$

A Key Lemma - Assessing $W[\mathcal{S}(\mathbf{x}, D)]$ (Cont'd)

Suppose that $\{d(j, k), 1 \leq j \leq J, 1 \leq k \leq K\}$ are **commensurable** and let Δ be their largest common divisor, and define

$$T_n(P, d) = (K-1)! \cdot (2\pi)^{K/2-1} \cdot \frac{\Delta \exp\{-s^*[(nD)\bmod\Delta]\}}{(1 - e^{-s^*\Delta})\sqrt{V(P, d)}},$$

If $\{d(j, k), 1 \leq j \leq J, 1 \leq k \leq K\}$ are incommensurable, take $\Delta \rightarrow 0$:

$$T_n(P, d) = \frac{(K-1)! \cdot (2\pi)^{K/2-1}}{s^* \sqrt{V(P, d)}}.$$

Lemma:

$$W[\mathcal{S}(\mathbf{x}, D)] = \frac{T_n(P, d)}{n^{K/2}} \cdot \exp\{-nR_d(D, P)\} \cdot [1 - \epsilon_{P, d}(n)].$$

The exact pre-exponent is essential for an exact characterization of the code-length redundancy in the sequel.

Main Analysis Tool - the Saddlepoint Method

Representing the unit step function $U(t)$ as the inverse Laplace transform of $1/z$:

$$U(t) = \frac{1}{2\pi j} \int_{c-j\infty}^{c+j\infty} \frac{e^{zt}}{z} dz, \quad c > 0,$$

we have:

$$\begin{aligned} W[\mathcal{S}(\mathbf{x}, D)] &= (K-1)! \sum_{\{\hat{\mathbf{x}}: d(\mathbf{x}, \hat{\mathbf{x}}) \leq nD\}} \int_{\mathcal{Q}} Q(\hat{\mathbf{x}}) dQ \\ &= (K-1)! \sum_{\hat{\mathbf{x}} \in \hat{\mathcal{X}}^n} U\left(nD - \sum_{i=1}^n d(x_i, \hat{x}_i)\right) \int_{\mathcal{Q}} Q(\hat{\mathbf{x}}) dQ \\ &= \frac{(K-1)!}{2\pi j} \int_{c-j\infty}^{c+j\infty} \int_{\mathcal{Q}} \frac{e^{-nF(z, Q)}}{z} \cdot dQ dz, \end{aligned}$$

and we select $c = s^*$ to pass thru all saddle-points.

A Universal Coding Scheme

- Generate \mathcal{C}_n with A^n independent random codewords ($A > K$), $\hat{\mathbf{X}}_i \sim W$, $i = 1, 2, \dots, A^n$.
- Reveal the codebook to both parties.
- Given \mathbf{x} and d , find $I_d(\mathbf{x}) = \min\{i : \hat{\mathbf{X}}_i \in \mathcal{S}(\mathbf{x}, D)\}$.
- Encode $I_d(\mathbf{x})$ using a Shannon code w.r.t. the distribution $u[i] \propto 1/i$, $i = 1, 2, \dots, A^n$.
- The decoder decodes $I_d(\mathbf{x})$ and outputs the $I_d(\mathbf{x})$ -th reproduction vector from \mathcal{C}_n .

Note that the codebook is the same for every (bounded) d – distortion-universality.

Coding Theorem

$\forall \epsilon > 0$, \exists a sequence of codebooks, $\{\mathcal{C}_n\}_{n \geq 1}$, and $\{\psi_n\}$, such that
 $\forall d \in \bigcup_{k \geq 1} \{0, d_{\max}/k, \dots, d_{\max}\}^{JK}$, $\exists \{\phi_n\}$, such that $\forall P \in \bigcup_{k \geq 1} \mathcal{P}_k$,
 $n \in \mathcal{N} \triangleq \{\hat{n} : d \in \mathcal{D}_{\hat{n}}, P \in \mathcal{P}_{\hat{n}}\}$ and $\mathbf{x} \in \mathcal{T}_n(P)$:

(a)

$$\begin{aligned} L_d(\mathbf{x}) &\leq nR_d(D, P) + \left(\frac{K}{2} + 2 + \epsilon \right) \cdot \ln n + \\ &\quad \beta_{P,d}(n) + \log(\log A + 1) + O(J^n e^{-n^{1+\epsilon}}). \end{aligned}$$

(b) The code is d -semifaithful: $d(\mathbf{x}, \psi_n(\phi_n(\mathbf{x}))) \leq nD$.

\mathcal{C}_n and ψ_n do not depend on P and d , but ϕ_n does.

Mahmood & Wagner ('22): 3 schemes with $\log n$ -coefficients: $2JK + J + 3$, $J(K + 1)$ and $J^2K^2 + J - 2$.

Converse Theorem

Let P and d be given. $\forall \epsilon > 0$ and sufficiently large n , \forall codebook that covers $\mathcal{T}_n(P)$ and every one-to-one variable-length code applied to that codebook, the following lower bound applies to a fraction of at least $(1 - 2n^{-\epsilon})$ of the codewords that cover $\mathcal{T}_n(P)$:

$$L_d(\hat{x}) \geq nR_d(D, P) + \left(\frac{1}{2} - \epsilon\right) \log n + c - c' \log(\log n),$$

where c and c' are constants that depend on P .

Converse Theorem (Cont'd)

The proof is based on a sphere-covering argument:

$$\log |\mathcal{T}_n(P)| \geq nH(P) - \frac{J-1}{2} \log n + c(P)$$

and

$$\begin{aligned} & \ln \left| \mathcal{T}_n(P) \bigcap \{ \mathbf{x} : d(\mathbf{x}, \hat{\mathbf{x}}) \leq nD \} \right| \\ & \leq \max_{\{P_{\hat{X}|X} : \mathbf{E}\{d(X, \hat{X})\} \leq D\}} H(X|\hat{X}) - \frac{J}{2} \log n + c' \log(\log n), \quad P_X = P \end{aligned}$$

and so,

$$|\mathcal{C}_n| \geq \exp_2 \left\{ nR_d(D, P) + \frac{\log n}{2} + \dots \right\}.$$

Most codewords cannot have code-length much less than $\log |\mathcal{C}_n|$.

Beyond the Memoryless Structure

Consider the universal distribution

$$U(\hat{\mathbf{x}}) = \frac{2^{-LZ(\hat{\mathbf{x}})}}{\sum_{\hat{\mathbf{x}}'} 2^{-LZ(\hat{\mathbf{x}}')}}$$

and let

$$U[\mathcal{S}(\mathbf{x}, D)] = \sum_{\hat{\mathbf{x}} \in \mathcal{S}(\mathbf{x}, D)} U(\hat{\mathbf{x}}).$$

Converse theorem: Let ℓ divide n and let $\mathcal{T}_n(\hat{P}^\ell)$ be any ℓ -th order type of source sequences. Let d be a distortion function that depends on $(\mathbf{x}, \hat{\mathbf{x}})$ only via $\hat{P}_{\mathbf{x}\hat{\mathbf{x}}}^1$. Then, $\forall d$ -semifaithful variable-length block code, and $\forall \epsilon > 0$, the following lower bound applies to a fraction of at least $(1 - 2n^{-\epsilon})$ of the codewords, $\{\phi_n(\mathbf{x}), \mathbf{x} \in \mathcal{T}_n(\hat{P}^\ell)\}$:

$$L(\phi_n(\mathbf{x})) \geq -\log(U[\mathcal{S}(\mathbf{x}, D)]) - n\Delta_n(\ell) - \epsilon \log n,$$

where $\lim_{n \rightarrow \infty} \Delta_n(\ell) = 1/\ell$.

Main Ideas of the Proof

Relating sphere-covering and $U[\mathcal{S}(\mathbf{x}, D)]$ in a few steps.

First, observe that

$$\begin{aligned} N(D) &\stackrel{\triangle}{=} \sum_{\mathbf{x}, \hat{\mathbf{x}}} \mathbb{I}\{\mathbf{x} \in \mathcal{T}_n(P^\ell), \hat{\mathbf{x}} \in \mathcal{T}_n(Q^\ell), d(\mathbf{x}, \hat{\mathbf{x}}) \leq nD\} \\ &= |\mathcal{T}_n(P^\ell)| \cdot \left| T_n(Q^\ell) \cap \mathcal{S}(\mathbf{x}, D) \right| \\ &= |\mathcal{T}_n(Q^\ell)| \cdot \left| T_n(P^\ell) \cap \hat{\mathcal{S}}(\hat{\mathbf{x}}, D) \right|, \quad \hat{\mathcal{S}}(\hat{\mathbf{x}}, D) \stackrel{\triangle}{=} \{\mathbf{x} : d(\mathbf{x}, \hat{\mathbf{x}}) \leq nD\} \end{aligned}$$

and so,

$$\frac{|T_n(P^\ell)|}{\left| T_n(P^\ell) \cap \hat{\mathcal{S}}(\hat{\mathbf{x}}, D) \right|} = \frac{|T_n(Q^\ell)|}{\left| T_n(Q^\ell) \cap \mathcal{S}(\mathbf{x}, D) \right|}$$

LHS = sphere-covering ratio;

RHS = $1/U_Q[\mathcal{S}(\mathbf{x}, D)] \geq 1/U[\mathcal{S}(\mathbf{x}, D)] \rightarrow$ use U for random coding!

Direct Theorem

Let $d : \mathcal{X}^n \times \hat{\mathcal{X}}^n \rightarrow \mathbb{R}^+$ be an arbitrary distortion function. Then, $\forall \epsilon > 0$, \exists sequence of d -semifaithful, variable-length block codes of block length n , such that $\forall \mathbf{x} \in \mathcal{X}^n$, the code length for \mathbf{x} is upper bounded by

$$L(\mathbf{x}) \leq -\log(U[\mathcal{S}(\mathbf{x}, D)]) + (2 + \epsilon) \log n + c + \delta_n,$$

where $c > 0$ is a constant and $\delta_n = O(nJ^n e^{-n^{1+\epsilon}})$.

The proof is very similar to that of the previous direct theorem.

Discussion

♠ Related to the Kontoyiannis-Zhang converse:

$$\forall \mathbf{x}, \mathcal{C}_n \exists Q : L(\mathbf{x}) \geq -\log Q[\mathcal{S}(\mathbf{x}, D)].$$

♠ $-\log(U[\mathcal{S}(\mathbf{x}, D)]) \sim \min_L \{L - \log |\{\hat{\mathbf{x}} : LZ(\hat{\mathbf{x}}) = L\} \cap \mathcal{S}(\mathbf{x}, D)|\}$, analogous to $\min_{P_{\hat{X}}} [H(\hat{X}) - \max\{H(\hat{X}|X) : \mathbf{E}d(X, \hat{X}) \leq D\}]$.

♠ Easy to see that the proposed scheme is better than $\min_{\hat{\mathbf{x}} \in \mathcal{S}(\mathbf{x}, D)} LZ(\hat{\mathbf{x}})$.
Complexity of both schemes depend on D .

♠ Universality w.r.t. a wide (continuous, parametric) class of distortion measures can also be proved. Here, the class distortion measures is quite arbitrary.