

Lossy Compression Theory (Achievability)

John MacLaren Walsh, Ph.D.

ECET 602, Winter Quarter, 2013

1 References

- *Elements of Information Theory*, T. M. Cover and J. A. Thomas, John Wiley and Sons, 1991.
- *Information Theory and Reliable Communication*, R. G. Gallager, John Wiley and Sons, 1968.
- *Information Theory: Coding Theorems For Discrete Memoryless Systems.*, I. Csiszar, J. Koerner, Academic Press, 1981.

These proof technique in these notes follows extremely closely that in Chapter 13 in the first reference (Chapter 10 in the second edition).

2 Overview

This lecture covers:

- achievability of the information rate distortion function via a random coding argument employing distortion typical sequences.

3 Ingredients

Consider two sequences $\mathbf{X}^N = (X_1, \dots, X_N)$, and $\hat{\mathbf{X}}^N = (\hat{X}_1, \dots, \hat{X}_N)$, with $\{(X_i, \hat{X}_i)\}$ i.i.d. from a joint distribution $\mathbf{p}_{\mathbf{X}, \hat{\mathbf{X}}}$. Recall that the motivation for the set of jointly typical sequences was based on the convergence in probability of

$$\begin{aligned} -\frac{1}{N} \log(\mathbf{p}(\mathbf{X}^N)) &\rightarrow \mathcal{H}(\mathbf{X}) \\ -\frac{1}{N} \log(\mathbf{p}(\hat{\mathbf{X}}^N)) &\rightarrow \mathcal{H}(\hat{\mathbf{X}}) \\ -\frac{1}{N} \log(\mathbf{p}(\hat{\mathbf{X}}^N, \mathbf{X}^N)) &\rightarrow \mathcal{H}(\hat{\mathbf{X}}, \mathbf{X}) \end{aligned}$$

for which we provided an easy proof using the weak law of large numbers. The jointly typical set was then defined to be the set of deterministic realizations x^n, \hat{x}^n for $\hat{\mathbf{X}}^N, \mathbf{X}^N$, for which each of these three convergences had happened within a common ϵ .

In order to prove the achievability for rate distortion, we observed another convergence in probability, namely if $\hat{\mathbf{X}}^N, \mathbf{X}^N$ are specified as above, then additionally, we must have

$$\frac{1}{N} \sum_{i=1}^N d(X_i, \hat{X}_i) \rightarrow \mathbb{E} \left[d(\mathbf{X}, \hat{\mathbf{X}}) \right] \quad (1)$$

in probability as $N \rightarrow \infty$. This is clear through the WLLN because the i.i.d. nature of X_i, \hat{X}_i implies that $d(X_i, \hat{X}_i)$ must also be i.i.d.

We can then consider a subset of the jointly typical sequences, called the *distortion jointly typical sequences* $\mathcal{A}_{\epsilon, d}^N$ for which the convergence (1) has also occurred within the same ϵ .

Naturally, then, the convergences in probability we have mentioned make that the probability $\Pr \left[(\mathbf{X}^N, \hat{\mathbf{X}}^N) \in \mathcal{A}_{\epsilon, d}^N \right] \rightarrow 1$ as $N \rightarrow \infty$. Another interesting property is for $(\hat{x}^N, x^N) \in \mathcal{A}_{\epsilon, d}^N$ that we can bound $\mathbf{p}(\hat{x}^N)$ using $\mathbf{p}(\hat{x}^N | x^N)$ according to

$$\mathbf{p}(\hat{x}^N) \geq \mathbf{p}(\hat{x}^N | x^N) 2^{-N(\mathcal{I}(\mathbf{X}; \hat{\mathbf{X}}) + 3\epsilon)} \quad (2)$$

We proved this with the chain of inequalities

$$\mathbf{p}(\hat{x}^N | x^N) = \mathbf{p}(\hat{x}^N) \frac{\mathbf{p}(x^N, \hat{x}^N)}{\mathbf{p}(x^N)\mathbf{p}(\hat{x}^N)} \leq \mathbf{p}(\hat{x}^N) \frac{2^{-N(\mathcal{H}(\mathbf{X}, \hat{\mathbf{X}}) - \epsilon)}}{2^{-N(\mathcal{H}(\hat{\mathbf{X}}) + \epsilon)} 2^{-N(\mathcal{H}(\mathbf{X}) + \epsilon)}} = \mathbf{p}(\hat{x}^N) 2^{N(\mathcal{I}(\mathbf{X}; \hat{\mathbf{X}}) + 3\epsilon)}$$

which proves the proposition.

We presently use these two new tools to prove the achievability of rate distortion codes.

4 Rate Distortion Theory: Achievability

Here, we construct a lossy source code with rate $R = R^I(D)$ and distortion arbitrarily close to D as the block length $N \rightarrow \infty$. Consider a source \mathbf{X}_i i.i.d. according to \mathbf{p}_X , with a bounded distortion metric $d \leq d_{max}$ and information rate distortion function

$$R^I(D) := \max_{\mathbf{p}(\hat{\mathbf{X}}|\mathbf{X}) : \sum \mathbf{p}(\hat{\mathbf{X}}|\mathbf{X})\mathbf{p}(\mathbf{X})d(\hat{\mathbf{X}}, \mathbf{X}) \leq D} \mathcal{I}(\hat{\mathbf{X}}; \mathbf{X}) \quad (3)$$

The first step in our construction is to choose a $\mathbf{p}_{\hat{\mathbf{X}}|\mathbf{X}}$ achieving the maximum in (3) for a distortion D . This induces (together with the distribution for the source symbols) a distribution for $\hat{\mathbf{X}}$ that is $\mathbf{p}_{\hat{\mathbf{X}}}(\hat{x}) = \sum_{x \in \mathcal{X}} \mathbf{p}_{\hat{\mathbf{X}}|\mathbf{X}}(\hat{x}|x) \mathbf{p}_X(x)$.

The next step in our construction is to generate a $2^{NR} \times N$ matrix \mathbf{C} of codewords, with the w th row denoted by $\hat{\mathbf{X}}^N(w)$, $w \in \{1, \dots, 2^{NR}\}$, by sampling each element of the matrix i.i.d. from the distribution $\mathbf{p}_{\hat{\mathbf{X}}}(\hat{x})$.

Encoding of an observation $\mathbf{X}^N = x^N$ is accomplished by sending w if the w th row of the codebook $\hat{\mathbf{X}}^N(w)$ is distortion jointly typical with x^N . If there is more than one such w we transmit the smallest, and if there is no such w , we transmit $w = 1$.

Decoding is accomplished at the decoder by mapping w to $\hat{\mathbf{X}}^N(w)$ (via knowledge of the codebook).

We presently prove that this code has expected (over both word to encode \mathbf{X}^N and codebook \mathbf{C}) distortion $D + \delta$ with δ vanishing as $N \rightarrow \infty$. To do this, we consider the following two possible situations

- There exists a w such that $\hat{\mathbf{X}}^N(w)$ and x^N are in the jointly distortion typical set $\mathcal{A}_{\epsilon, d}^N$. This set of sequences x^N and codebooks \mathbf{C} can not have more probability than one, and furthermore, by definition of the jointly distortion typical set the distortion $d(x^n, \hat{\mathbf{X}}^N(w)) < D + \epsilon$, thus these sequences contribute no more than $D + \epsilon$ to the overall distortion.
- The sequence x^N and codebook \mathbf{C} are such that there does not exist such a w with the w th codeword jointly distortion typical with x^N . Let the collective probability of these sequences and codebook pairs be denoted by P_e , then these sequences can contribute no more than $P_e d_{max}$ to the overall expected distortion.

Thus, having established that the expected distortion is bounded above by $D + \epsilon + P_e d_{max}$, all that remains to show is that P_e can be made arbitrarily small as $N \rightarrow \infty$. To do this, we obtain an upper bound on P_e . First off all, note that

$$P_e = \sum_{\mathbf{C}} \mathbf{p}(\mathbf{C}) \sum_{x^N | \nexists \hat{x}^N \in \mathbf{C} \text{ s.t. } (x^N, \hat{x}^N) \in \mathcal{A}_{\epsilon, d}^N} \mathbf{p}_{\mathbf{X}^N}(x^N)$$

we can rewrite this as

$$P_e = \sum_{x^N \in \mathcal{X}^N} \mathbf{p}_{\mathbf{X}^N}(x^N) \sum_{\mathbf{C} | \nexists \hat{x}^N \in \mathbf{C} \text{ s.t. } (x^N, \hat{x}^N) \in \mathcal{A}_{\epsilon, d}^N} \mathbf{p}(\mathbf{C})$$

Now, considering the fact that the rows of the code matrix (i.e. the codewords) are selected independently of one another in an identically distributed, we can write the result of the latter sum as the product over

the 2^{NR} different code words of the probability that one selected code word is not jointly typical with the given x^N . Thus, we have

$$P_e = \sum_{x^N \in \mathcal{X}^N} \mathbf{p}_{X^N}(x^N) \left(1 - \sum_{\hat{x}^N | (x^N, \hat{x}^N) \in \mathcal{A}_{\epsilon, d}^N} \mathbf{p}_{\hat{X}^N}(\hat{x}^N) \right)^{2^{NR}}$$

We can bound the probability in the parenthesis using (2) to get

$$P_e \leq \sum_{x^N \in \mathcal{X}^N} \mathbf{p}_{X^N}(x^N) \left(1 - \sum_{\hat{x}^N | (x^N, \hat{x}^N) \in \mathcal{A}_{\epsilon, d}^N} \mathbf{p}(\hat{x}^N | x^N) 2^{-N(\mathcal{I}(X; \hat{X}) + 3\epsilon)} \right)^{2^{NR}}$$

We then use the deterministic bound that for $0 \leq x, y \leq 1, N > 0$ $(1 - xy)^N \leq 1 - x + e^{-yN}$ to get

$$P_e \leq \sum_{x^N \in \mathcal{X}^N} \mathbf{p}_{X^N}(x^N) \left(1 - \sum_{\hat{x}^N | (x^N, \hat{x}^N) \in \mathcal{A}_{\epsilon, d}^N} \mathbf{p}(\hat{x}^N | x^N) + e^{-2^{-N(\mathcal{I}(X; \hat{X}) + 3\epsilon) + NR}} \right) \quad (4)$$

$$= 1 - \sum_{(x^N, \hat{x}^N) \in \mathcal{A}_{\epsilon, d}^N} \mathbf{p}(\hat{x}^N, x^N) + e^{-2^{-N(\mathcal{I}(X; \hat{X}) + 3\epsilon) + NR}} \quad (5)$$

The first two terms are then recognized as one minus the probability of the distortion typical set, which we know to be $\leq \epsilon$. Finally, remembering that we selected $\mathbf{p}(\hat{X}|X)$ to achieve the maximum in $R^I(D)$ for D , we observe that the last term may be rewritten to yield

$$P_e \leq \epsilon + e^{-2^{N(R - R^I(D) - 3\epsilon)}}$$

which goes to zero as $N \rightarrow \infty$ if $R \geq R^I(D) + 3\epsilon$. We have thus established achievability of the rate distortion function.