
Compressing Tabular Data via Latent Variable Estimation

Andrea Montanari¹ Eric Weiner¹

Abstract

Data used for analytics and machine learning often take the form of tables with categorical entries. We introduce a family of lossless compression algorithms for such data that proceed in four steps: (i) Estimate latent variables associated to rows and columns; (ii) Partition the table in blocks according to the row/column latents; (iii) Apply a sequential (e.g. Lempel-Ziv) coder to each of the blocks; (iv) Append a compressed encoding of the latents. We evaluate this approach on several benchmark datasets, and study optimal compression in a probabilistic model for tabular data, whereby latent values are independent and table entries are conditionally independent given the latent values. We prove that the model has a well defined entropy rate and satisfies an asymptotic equipartition property. We also prove that classical compression schemes such as Lempel-Ziv and finite-state encoders do not achieve this rate. On the other hand, the latent estimation strategy outlined above achieves the optimal rate.

1. Introduction

Classical theory of lossless compression (Cover & Thomas, 2006; Salomon, 2004) assumes that data take the form of a random vector $\mathbf{X}^N = (X_1, X_2, \dots, X_N)$ of length N with entries in a finite alphabet \mathcal{X} . Under suitable ergodicity assumptions, the entropy per letter converges to a limit $h := \lim_{N \rightarrow \infty} H(\mathbf{X}^N)/N$ (Shannon-McMillan-Breiman theorem). Universal coding schemes (e.g. Lempel-Ziv coding) do not require knowledge of the distribution of \mathbf{X}^N , and can encode such a sequence without information loss using (asymptotically) h bits per symbol.

While this theory is mathematically satisfying, its modeling assumptions (stationarity, ergodicity) are unlikely to be satisfied in many applications. This has long been recognized by practitioners. The main objective of this paper

¹Project N, Mountain View, CA, United States. Correspondence to: Andrea Montanari <am@projectn.co>.

is to investigate this fact mathematically in the context of tabular data, characterize the gap to optimality of classical schemes, and describe an asymptotically optimal algorithm that overcomes their limitations.

We consider a data table with m rows and n columns and entries in \mathcal{X} , $\mathbf{X}^{m,n} \in \mathcal{X}^{m \times n}$ $\mathbf{X}^{m,n} := (X_{ij})_{i \leq m, j \leq n}$. The standard approach to such data is: (i) Serialize, e.g. in row-first order, to form a vector of length $N = mn$, $\mathbf{X}^N = (X_{11}, X_{12}, \dots, X_{1n}, X_{21}, \dots, X_{mn})$; (ii) Apply a standard compressor (e.g., Lempel-Ziv) to this vector.

We will show, both empirically and mathematically, that this standard approach can be suboptimal in the sense of not achieving the optimal compression rate. This happens even in the limit of large tables, as long as the number of columns and rows are polynomially related (i.e. $n^\varepsilon \leq m \leq n^M$ for some small constant ε and large constant M).

We advocate an alternative approach:

1. Estimate row/column latents $\mathbf{u}^m = (u_1, \dots, u_m) \in \mathcal{L}^m$, $\mathbf{v}^n = (v_1, \dots, v_n) \in \mathcal{L}^n$, with \mathcal{L} a finite alphabet.
2. Partition the table in blocks according to the row/column latents, Namely, for $u, v \in \mathcal{L}$, define

$$\mathbf{X}(u, v) = \text{vec}(X_{ij} : u_i = u, v_j = v). \quad (1.1)$$

where $\text{vec}(\mathbf{M})$ denote the serialization of matrix \mathbf{M} (either row-wise or column-wise).

3. Apply a base compressor (generically denoted by $Z_{\mathcal{X}} : \mathcal{X}^* \rightarrow \{0, 1\}^*$) to each block $\mathbf{X}(u, v)$

$$\mathbf{z}(u, v) = Z_{\mathcal{X}}(\mathbf{X}(u, v)), \quad \forall u, v \in \mathcal{L}. \quad (1.2)$$

4. Encode the row latents and column latents using a possibly different compressor $Z_{\mathcal{L}} : \mathcal{L}^* \rightarrow \{0, 1\}^*$, to get $\mathbf{z}_{\text{row}} = Z_{\mathcal{L}}(\mathbf{u})$, $\mathbf{z}_{\text{col}} = Z_{\mathcal{L}}(\mathbf{v})$. Finally output the concatenation (denoted by \oplus)

$$\text{Enc}(\mathbf{X}^{m,n}) = \text{header} \oplus \mathbf{z}_{\text{row}} \oplus \mathbf{z}_{\text{col}} \oplus \bigoplus_{u, v \in \mathcal{L}} \mathbf{z}(u, v). \quad (1.3)$$

Here header is a header that contains encodings of the lengths of subsequent segments.

Note that encoding the latents can in general lead to a sub-optimal compression rate. While this can be remedied with techniques such as bits-back coding, we observed in our applications that using such techniques yields limited improvement. Our analysis shows that the rate improvement afforded by bits-back coding is only significant in certain special regimes. We refer to Sections 5 and 6 for further discussion.

The above description leaves several design choices undefined, namely: (a) The latents estimation procedure at point 1; (b) The base compressor $Z_{\mathcal{X}}$ for the blocks $\mathbf{X}(u, v)$; (c) The base compressor $Z_{\mathcal{L}}$ for the latents.

We will provide details for a specific implementation in Section 2, alongside empirical evaluation in Section 3. Section 4 introduces a probabilistic model for the data $\mathbf{X}^{m,n}$, and Section 5 establishes our main theoretical results: standard compression schemes are suboptimal on this model, while the above latents-based approach is asymptotically optimal. Finally we discuss extensions in Section 6.

1.1. Related work

The use of latent variables is quite prevalent in compression methods based on machine learning and probabilistic modeling. Hinton and Zemel (1993) introduced the idea that stochastically generated codewords (e.g., random latents) can lead to minimum description lengths via bits back coding. This idea was explicitly applied to lossless compression using arithmetic coding in (Frey & Hinton, 1996), and ANS coding in (Townsend et al., 2019a;b).

Compression via low-rank approximation is closely-related to our latents-based approach and has been studied in the past. An incomplete list of contributions includes (Cheng et al., 2005) (numerical analysis), (Li & Li, 2010) (hyper-spectral imaging), (Yuan & Oja, 2005; Hou et al., 2015) (image processing), (Taylor, 2013) (quantum chemistry), (Phan et al., 2020) (compressing the gradient for distributed optimization), (Chen et al., 2021) (large language models compression).

The present paper contributes to this line of work, but departs from it in a number of ways. (i) We study lossless compression while earlier work is mainly centered on lossy compression. (ii) Most of the papers in this literature do not precisely quantify compression rate: they do not ‘count bits.’ (iii) We show empirically an improvement in terms of lossless compression rate over state of the art.

Another related area is network compression: simple graphs can be viewed as matrices with entries in $\{0, 1\}$. In the case of graph compression, one is interested only in such matrices up to graph isomorphisms. The idea of reordering the nodes of the network and exploiting similarity between nodes has been investigated in this context, see e.g. (Boldi & Vigna, 2004; Chierichetti et al., 2009; Lim et al., 2014;

Algorithm 1 Latent-based Tabular Compressor

Input: Data matrix $\mathbf{X}^{m,n} \in \mathcal{X}^{m \times n}$, range $k = |\mathcal{L}|$
Output: Compressed data $\text{Enc}(\mathbf{X}^{m,n}) \in \{0, 1\}^*$

Estimate latents $\mathbf{u}^m \in [k]^m$, $\mathbf{v}^n \in [k]^n$ using Algorithm 2, with inputs $\mathbf{X}^{m,n}$, k

for $u, v \in [k]$ **do**

$\mathbf{X}(u, v) = \text{vec}(X_{ij} : u_i = u, v_j = v)$

$\mathbf{z}(u, v) = Z_{\mathcal{X}}(\mathbf{X}(u, v))$

end for

Compute $\mathbf{z}_{\text{row}} = Z_{\mathcal{L}}(\mathbf{u})$, $\mathbf{z}_{\text{col}} = Z_{\mathcal{L}}(\mathbf{v})$

return concatenation of $\{\mathbf{z}(u, v) : u, v \in [k]\}$, $\mathbf{z}_{\text{row}}, \mathbf{z}_{\text{col}}$, metadata

(Besta & Hoefler, 2018) However, we are not aware of results analogous to ours in this literature.

To the best of our knowledge, our work is the first to prove that classical lossless compression techniques do not achieve the ideal compression rate under a probabilistic model for tabular data. We characterize this ideal rate as well as the one achieved by classical compressors, and prove that latents estimation can be used to close this gap.

1.2. Notations

We generally use boldface for vectors and uppercase boldface for matrices, without making any typographic distinction between numbers and random variables. When useful, we indicate by superscripts the dimensions of a matrix or a vector: \mathbf{u}^m is a vector of length m , and $\mathbf{X}^{m,n}$ is a matrix of dimensions $m \times n$. For a string v and $a \leq b$, we use $v_a^b = (v_a, \dots, v_b)$ to denote the substring of v .

If X, Y are random variables on a common probability space $(\Omega, \mathcal{F}, \mathbb{P})$, we denote by $H(X)$, $H(Y)$ their entropies, $H(X, Y)$ their joint entropy, $H(X|Y)$ the conditional entropy of X given Y . We will overload this notation: if p is a discrete probability distribution, we denote by $H(p)$ its entropy. Unless stated otherwise, all entropies will be measured in bits. For $\varepsilon \in [0, 1]$, $h(\varepsilon) := -\varepsilon \log_2 \varepsilon - (1 - \varepsilon) \log_2 (1 - \varepsilon)$.

2. Implementation

The overall structure of the compression algorithm was already described in the introduction. Algorithm 1 summarizes it. In this section we provide further details about the two basic components it relies on: the base compressor Z_{\cdot} , and the latents estimation algorithm. In both cases, the choice is in no-way unique and we only describe what we used in our implementation.

2.1. Base compressors

We implemented the following two options for the base compressors $Z_{\mathcal{X}}$ (for data blocks) and $Z_{\mathcal{L}}$ (for latents).

Dictionary-based compression (Lempel-Ziv, LZ). For this we used Zstandard (ZSTD) Python bindings to the C implementation using the library `zstd`, with level 12. While ZSTD can use run-length encoding schemes or literal encoding schemes, we verified that in this case ZSTD always use its LZ algorithm.

The LZ algorithm in ZSTD is somewhat more sophisticated than the plain LZ algorithm used in our proofs. In particular it includes (Collet & Kucherawy, 2018) Huffman coding of literals 0-255 and entropy coding of the LZ stream. Experiments with other (simpler) LZ implementations yielded similar results. We focus on ZSTD because of its broad adoption in industry.

Frequency-based entropy coding (ANS). For each data portion (i.e each block $\mathbf{X}(u, v)$ and each of the row latents u and column latents v) compute empirical frequencies of the corresponding symbols. Namely for all $u, v \in \mathcal{L}$, $x \in \mathcal{X}$, we compute

$$\hat{Q}(x|u, v) := \frac{1}{N(u, v)} \sum_{i:u_i=u} \sum_{j:v_j=v} \mathbf{1}_{x_{ij}=x},$$

$$\hat{q}_r(u) := \frac{1}{m} \sum_{i=1}^m \mathbf{1}_{u_i=u}, \quad \hat{q}_c(v) := \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{v_i=v},$$

where $N(u, v)$ is the number of $i \leq m, j \leq n$ such that $u_i = u, v_j = v$. We then apply ANS coding (Duda, 2009) to each block $\mathbf{X}(u, v)$ modeling its entries as independent with distribution $\hat{Q}(\cdot | u, v)$, and to the latents u^m, v^n using the distributions $\hat{q}_r(\cdot), \hat{q}_c(\cdot)$. We separately encode these counts as long integers.

Since our main objective was to study the impact of learning latents, we did not try to optimize these base compressors.

2.2. Latent estimation

We implemented latents estimation using a spectral clustering algorithm outlined in Algorithm 2.

In words, the algorithm encodes the data matrix $\mathbf{X}^{m,n}$ as an $m \times n$ real-valued matrix $\mathbf{M}^{m,n} \in \mathbb{R}^{m \times n}$ using a map $\psi : \mathcal{X} \rightarrow \mathbb{R}$. It then computes the top $k-1$ left and right singular vectors of \mathbf{M} , and stores them as matrices $\mathbf{A} \in \mathbb{R}^{m \times (k-1)}, \mathbf{B} \in \mathbb{R}^{n \times (k-1)}$. The rows $\mathbf{a}_i, \mathbf{b}_j \in \mathbb{R}^{k-1}$ of these matrices are used as embedding of the rows and columns indices in $k-1$ dimensions. Finally, we run KMeans on these vectors to construct k clusters of rows/columns.

A few remarks are in order. The algorithm encodes the data matrix $\mathbf{X}^{m,n}$ as an $m \times n$ real-valued matrix $\mathbf{M}^{m,n} \in$

Algorithm 2 Spectral latents estimation

Input: Data matrix $\mathbf{X}^{m,n} \in \mathcal{X}^{m \times n}$
 latents range $k = |\mathcal{L}|$; map $\psi : \mathcal{X} \rightarrow \mathbb{R}$
Output: Factors $\mathbf{u}^m \in \mathcal{L}^m, \mathbf{v}^n \in \mathcal{L}^n$

Compute top $(k-1)$ singular vectors of $\mathbf{M}^{m,n} = \psi(\mathbf{X}^{m,n}), (\tilde{\mathbf{a}}_i)_{i \leq k-1}, (\tilde{\mathbf{b}}_i)_{i \leq k-1}$
 Stack singular vectors in matrices $\mathbf{A} = [\tilde{\mathbf{a}}_1 | \dots | \tilde{\mathbf{a}}_{k-1}] \in \mathbb{R}^{m \times (k-1)}, \mathbf{B} = [\tilde{\mathbf{b}}_1 | \dots | \tilde{\mathbf{b}}_{k-1}] \in \mathbb{R}^{n \times (k-1)}$;
 Let $(\mathbf{a}_i)_{i \leq m}, \mathbf{a}_i \in \mathbb{R}^{k-1}$ be the rows of \mathbf{A} ; $(\mathbf{b}_i)_{i \leq n}, \mathbf{b}_i \in \mathbb{R}^{k-1}$ the rows of \mathbf{B}
 Apply KMeans to $(\mathbf{a}_i)_{i \leq m}$; store the cluster labels as vector \mathbf{u}^m
 Apply KMeans to $(\mathbf{b}_i)_{i \leq n}$; store the cluster labels as vector \mathbf{v}^n
return $\mathbf{u}^m, \mathbf{v}^n$

$\mathbb{R}^{m \times n}$ using a map $\psi : \mathcal{X} \rightarrow \mathbb{R}$. In our experiments we did not optimize this map and encoded the elements of \mathcal{X} as $0, 1, \dots, |\mathcal{X}| - 1$ arbitrarily, cf. also Section 5.3

The singular vector calculation turns out to be the most time consuming part of the algorithm. Computing approximate singular vectors via power iteration requires in this case of the order of $\log(m \wedge n)$ matrix vector multiplications for each of k vectors. This amounts to $mnk \log(m \wedge n)$ operations, which is larger than the time needed to compress the blocks or to run KMeans. A substantial speed-up is obtained via row subsampling, cf. Section 6

For the clustering step we use the scikit-learn implementation via `sklearn.cluster.KMeans`, with random initialization.

3. Empirical evaluation

We evaluated our approach on tabular datasets with different origins. Our objective is to assess the impact of using latents in reordering columns and rows, so we will not attempt to achieve the best possible data reduction rate (DRR) on each dataset, but rather to compare compression with latents and without in as-uniform-as-possible fashion.

Since our focus is on categorical variables, we preprocess the data to fit in this setting as described in Section A.2. This preprocessing step might involve dropping some of the columns of the original table. We denote the number of columns after preprocessing by n .

We point out two simple improvements we introduce in the implementation: (i) We use different sizes for rows latent alphabet and column latent alphabet $|\mathcal{L}_r| \neq |\mathcal{L}_c|$; (ii) We choose $|\mathcal{L}_r|, |\mathcal{L}_c|$ by optimizing the compressed size .

3.1. Datasets

More details on these data can be found in Appendix A.1:

Taxicab. A table with $m = 62,495$, $n = 18$ (NYC.gov, 2022). LZ: $|\mathcal{L}_r| = 9$, $|\mathcal{L}_c| = 15$. ANS: $|\mathcal{L}_r| = 5$, $|\mathcal{L}_c| = 14$.

Network. Four social networks from (Leskovec & Krevl, 2014) with $m = n \in \{333, 747, 786, 1187\}$. LZ and ANS: $|\mathcal{L}_r| = 5$, $|\mathcal{L}_c| = 5$.

Card transactions. A table with $m = 24,386,900$ and $n = 12$ (Altman, 2019). LZ and ANS: $|\mathcal{L}_r| = 3$, $|\mathcal{L}_c| = n$.

Business price index. A table with $m = 72,750$ and $n = 10$ (stats.govt.nz, 2022). LZ: $|\mathcal{L}_r| = 6$, $|\mathcal{L}_c| = 7$. ANS: $|\mathcal{L}_r| = 2$, $|\mathcal{L}_c| = 6$.

Forest. A table from the UCI data repository with $m = 581,011$, $n = 55$ (Dua & Graff, 2017). LZ and ANS: $|\mathcal{L}_r| = 6$, $|\mathcal{L}_c| = 17$.

US Census. Another table from (Dua & Graff, 2017) with $m = 2,458,285$ and $n = 68$. LZ and ANS: $|\mathcal{L}_r| = 9$, $|\mathcal{L}_c| = 68$.

Jokes. A collaborative filtering dataset with $m = 23,983$ rows and $n = 101$ (Goldberg et al., 2001; Goldberg et al.). LZ: $|\mathcal{L}_r| = 2$, $|\mathcal{L}_c| = 101$. ANS: $|\mathcal{L}_r| = 8$, $|\mathcal{L}_c| = 8$.

3.2. Results

Given a lossless encoder $\phi : \mathcal{X}^{m \times n} \rightarrow \{0, 1\}^*$, we define its compression rate and data reduction rate (DRR) as

$$\begin{aligned} R_\phi(\mathbf{X}^{m,n}) &:= \frac{\text{len}(\phi(\mathbf{X}^{m,n}))}{mn \log_2 |\mathcal{X}|}, \\ \text{DRR}_\phi(\mathbf{X}^{m,n}) &:= 1 - R_\phi(\mathbf{X}^{m,n}). \end{aligned} \quad (3.1)$$

(Larger DRR means better compression.)

The DRR of each algorithm is reported in Table 1. For the table of results, **LZ** refers to row-major order ZSTD, **LZ (c)** refers to column-major order ZSTD. We run KMeans on the data 5 times, with random initializations finding the DRR each time and reporting the average.

We make the following observations on the empirical results of Table 1. First, Latent + ANS encoder achieves systematically the best DRR. Second, the use of latent in several cases yields a DRR improvement of 5% (of the uncompressed size) or more. Third, as intuitively natural, this improvement appears to be larger for data with a large number of columns (e.g. the network data).

The analysis of the next section provides further support for these findings.

4. A probabilistic model

In order to better understand the limitations of classical approaches, and the optimality of latent-based compression, we introduce a probabilistic model for the table $\mathbf{X}^{m,n} \in \mathcal{X}^{m \times n}$. We assume the true latents $(u_i)_{i \leq m}$, $(v_j)_{j \leq n}$ to be

independent random variables with

$$\mathbb{P}(u_i = u) = q_r(u), \quad \mathbb{P}(v_j = v) = q_c(v). \quad (4.1)$$

We assume that the entries $(X_{ij})_{i \leq m, j \leq n}$ are conditionally independent given $\mathbf{u}^m = (u_i)_{i \leq m}$, $\mathbf{v}^n = (v_j)_{j \leq n}$, with

$$\mathbb{P}(X_{ij} = x | \mathbf{u}^m, \mathbf{v}^n) = Q(x | u_i, v_j). \quad (4.2)$$

The distributions q_r, q_c , and conditional distribution Q are parameters of the model (a total of $2(|\mathcal{L}| - 1) + |\mathcal{L}|^2(|\mathcal{X}| - 1)$ real parameters). We will write $(\mathbf{X}^{m,n}, \mathbf{u}^m, \mathbf{v}^n) \sim \mathcal{T}(Q, q_r, q_c; m, n)$ to indicate that the triple $(\mathbf{X}^{m,n}, \mathbf{u}^m, \mathbf{v}^n)$ is distributed according to the model.

Remark 4.1. Some of our statements will be non-asymptotic, in which case $m, n, \mathcal{X}, \mathcal{L}, Q, q_r, q_c$ are fixed. Others will be of asymptotic. In the latter case, we have in mind a sequence of problems indexed by n . In principle, we could write $m_n, \mathcal{X}_n, \mathcal{L}_n, Q_n, q_{r,n}, q_{c,n}$ to emphasize the fact that these quantities depend on n . However, we will typically omit these subscripts.

Example 4.2 (Symmetric Binary Model). As a toy example, we will use the following Symmetric Binary Model (SBM) which parallels the symmetric stochastic block model for community detection (Holland et al., 1983). We take $\mathcal{L} = [k] := \{1, \dots, k\}$, $\mathcal{X} = \{0, 1\}$, $q_r = q_c = \text{Unif}([k])$ (the uniform distribution over $[k]$) and

$$Q(1|u, v) = p_1 \text{ if } u = v, \quad Q(1|u, v) = p_0 \text{ if } u \neq v. \quad (4.3)$$

We will write $(\mathbf{X}^{m,n}, \mathbf{u}^m, \mathbf{v}^n) \sim \mathcal{T}_{\text{SBM}}(p_0, p_1, k; m, n)$ when this distribution is used.

Figure 1 reports the results of simulations within the SBM, for ZSTD and ANS base compressors. In this case $m = n = 1000$, $k = 3$, and we average DRR values over 4 realizations. Appendix B reports additional simulations under the same model for $k \in \{5, 7\}$: the results are very similar to the ones of Figure 1. As expected, the use of latents is irrelevant along the line $p_1 \approx p_0$ (in this case, the latents do not impact the distribution of X_{ij}). However, it becomes important when p_1 and p_0 are significantly different.

The figures also report contour lines of the theoretical predictions for the asymptotic DRR of various compression algorithms (cf. Example 5.4). The agreement is excellent.

5. Theoretical analysis

In this section we present our theoretical results on compression rates under the model $\mathcal{T}(Q, q_r, q_c, k; m, n)$ introduced above. We first characterize the optimal compression rate in Section 5.1, then prove that standard compression methods fail to attain this goal in Section 5.2, and finally show that

Table 1: Data reduction rate (DRR) achieved by classical and latent-based compressors on real tabular data.

Data	Size	LZ	LZ (c)	ANS	Latent + LZ	Latent + ANS
Taxicab	380 KB	0.41	0.44	0.43	0.48	0.54
FB Network 1	13.6 KB	0.63	0.63	0.76	0.58	0.78
FB Network 2	68.1 KB	0.44	0.44	0.57	0.64	0.75
FB Network 3	75.4 KB	0.59	0.59	0.75	0.69	0.80
GP Network 1	172 KB	0.46	0.46	0.65	0.58	0.70
Forest (s)	6.10 MB	0.29	0.38	0.47	0.41	0.49
Card Transactions (s)	123 MB	0.03	0.21	0.29	0.20	0.30
Business price index (s)	153 KB	-0.03	0.20	0.28	0.25	0.32
US Census	43.9 MB	0.38	0.31	0.47	0.52	0.62
Jokes	515 KB	-0.21	-0.15	0.07	-0.03	0.14

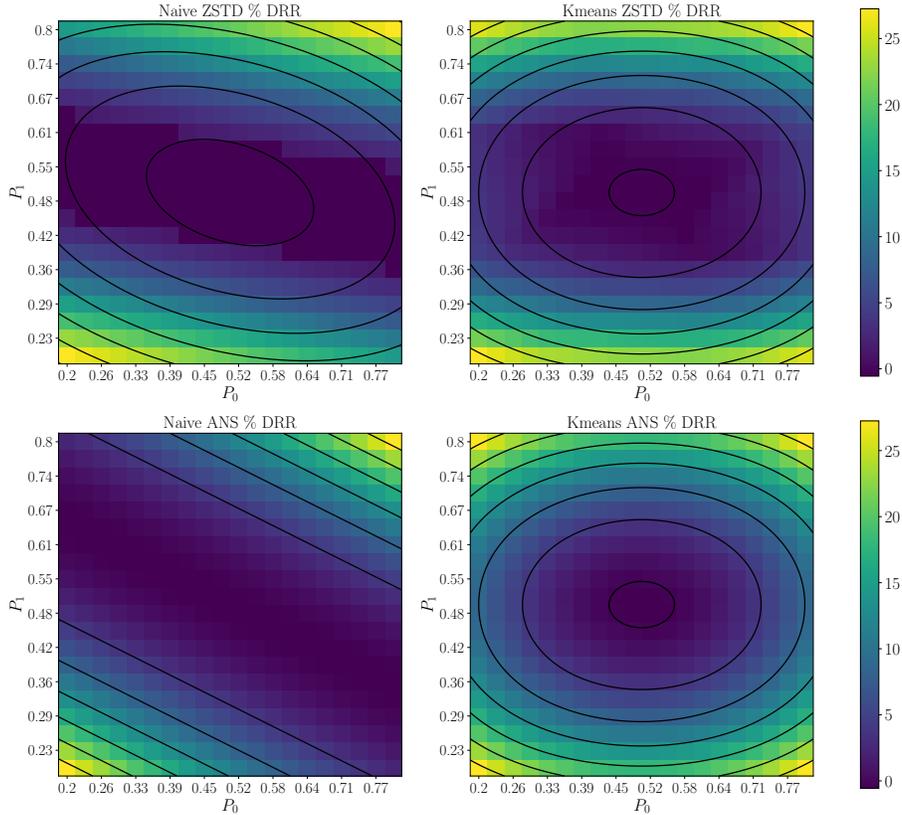


Figure 1: Comparing data reduction rate of naive coding and latent-based coding for synthetically generated data. Top: ZSTD base compressor. Bottom: ANS base compressor. Contour lines correspond to the compression rate predicted by the theorems of Section 5 (coinciding with optimal rate for latent-based encoders).

latent-based compression does in Section 5.3. Proofs are deferred to Appendices D, E, F, G.

Throughout, we denote by (X, U, V) a triple with joint distribution $\mathbb{P}(X = x, U = u, V = v) = Q(x|u, v)q_U(u)q_V(v)$ (this is the same as the joint distribution of (X_{ij}, u_i, v_j) for fixed i, j).

5.1. Ideal compression

Our first lemma provides upper and lower bounds on the entropy per symbol $H(\mathbf{X}^{m,n})/mn$.

Lemma 5.1. *Defining $H_{m,n}^+(X|U, V) := H(X|U, V) + \frac{1}{n}H(U) + \frac{1}{m}H(V)$, we have*

$$H(X|U, V) \leq \frac{1}{mn}H(\mathbf{X}^{m,n}) \leq H_{m,n}^+(X|U, V). \quad (5.1)$$

Further, for any estimators $\hat{u} : \mathcal{X}^{m \times n} \rightarrow \mathcal{L}^m$, $\hat{v} : \mathcal{X}^{m \times n} \rightarrow \mathcal{L}^n$, let $\text{Err}_U := \min_{\pi \in \mathfrak{S}_{\mathcal{L}}} \sum_{i=1}^m \mathbf{1}_{\hat{u}_i \neq \pi(u_i)}/m$, $\text{Err}_V := \min_{\pi \in \mathfrak{S}_{\mathcal{L}}} \sum_{i=1}^n \mathbf{1}_{\hat{v}_i \neq \pi(v_i)}/n$ (min over permutations of \mathcal{L}), letting $\varepsilon_U := \mathbb{E} \text{Err}_U$, $\varepsilon_V := \mathbb{E} \text{Err}_V$, we have

$$H_{m,n}^+(X|U, V) - \delta_{m,n} \leq \frac{1}{mn} H(\mathbf{X}^{m,n}) \leq H_{m,n}^+(X|U, V). \quad (5.2)$$

where $\delta_{m,n} := \delta(\varepsilon_U)/n + \delta(\varepsilon_V)/m$ and $\delta(\varepsilon) := h(\varepsilon) + \varepsilon \log(|\mathcal{L}| - 1)$.

Corollary 5.2. *There exists a lossless compressor ϕ whose rate (cf. Eq. (3.1)) is*

$$\mathbb{E} R_{\phi}(\mathbf{X}^{m,n}) \leq \frac{1}{\log_2 |\mathcal{X}|} \left\{ H_{m,n}^+(X|U, V) + \frac{1}{mn} \right\}. \quad (5.3)$$

Further, for any lossless compressor ϕ , $\mathbb{E} R_{\phi}(\mathbf{X}^{m,n}) \geq H_{m,n}^+(X|U, V) - \delta_{m,n} - 2 \log_2(mn)/mn$.

Remark 5.1. The simpler bound (5.1) implies that the entropy per entry is $H(X|U, V) + O(1/(m \wedge n))$. The operational interpretation of this result is that we should be able to achieve the same compression rate per symbol as if the latents were given to us.

The additional terms $\frac{1}{n} H(U) + \frac{1}{m} H(V)$ in Eq. (5.2) account for the additional memory required for the latents. The lower bound in Eq. (5.2) implies that, if the latents can be accurately estimated from the data $\mathbf{X}^{m,n}$ (that is if $\varepsilon_U, \varepsilon_V$ are small), then this overhead is essentially unavoidable.

The nearly ideal compression rate in Eq. (5.3) can be achieved by Huffman or arithmetic coding, and requires knowledge of the probability distribution of $\mathbf{X}^{m,n}$. Under these schemes, the length of the codeword associated to $\mathbf{X}^{m,n}$ is within constant number of bits from $-\log_2 P(\mathbf{X}^{m,n})$, where $P(\mathbf{X}_0) := \mathbb{P}(\mathbf{X}^{m,n} = \mathbf{X}_0)$ is the probability mass function of the random table $\mathbf{X}^{m,n}$ (Cover & Thomas, 2006; Salomon, 2004). The next lemma implies that the length concentrates tightly around the entropy.

Lemma 5.3 (Asymptotic Equipartition Property). *For $\mathbf{X}_0 \in \mathcal{X}^{m \times n}$, let $P(\mathbf{X}_0) = P_{Q, q_r, q_c; m, n}(\mathbf{X}_0)$ the probability of $\mathbf{X}^{m,n} = \mathbf{X}_0$ under model $\mathbf{X}^{m,n} \sim \mathcal{T}(Q, q_r, q_c; m, n)$. Assume there exists a constant $c > 0$ such that $\min_{x \in \mathcal{X}} \min_{u, v \in \mathcal{L}} Q(x|u, v) \geq c$. Then there exists a constant C (depending on c) such that the following happens.*

For $\mathbf{X}^{m,n} \sim \mathcal{T}(Q, q_r, q_c; m, n)$ and any $t \geq 0$ with probability at least $1 - 2e^{-t}$:

$$|-\log P(\mathbf{X}^{m,n}) - H(\mathbf{X}^{m,n})| \leq C \sqrt{mn(m+n)} t. \quad (5.4)$$

For the sake of simplicity, in the last statement we assume a uniform lower bound on $Q(x|u, v)$. While such a lower

bound holds without loss of generality when Q is independent of m, n (symbols with zero probability can be dropped), it might not hold in the n -dependent case. Appendix D gives a more general statement.

5.2. Failure of classical compression schemes

We analyze two types of codes: finite-state encoders and Lempel-Ziv codes. Both operate on the serialized data $\mathbf{X}^N = \text{vec}(\mathbf{X}^{m,n})$, $N = mn$, obtained by scanning the table in row-first order (obviously column-first yields symmetric results).

5.2.1. FINITE STATE ENCODERS

A finite state (FS) encoder takes the form of a triple (Σ, f, g) with Σ a finite set of cardinality $M = |\Sigma|$ and $f : \mathcal{X} \times \Sigma \rightarrow \{0, 1\}^*$, $g : \mathcal{X} \times \Sigma \rightarrow \Sigma$.

We assume that Σ contains a special ‘initialization’ symbol s_{init} . Starting from state $s_0 = s_{\text{init}}$, the encoder scans the input \mathbf{X}^N sequentially. Assume after the first ℓ input symbols it is in state s_{ℓ} , and produced encoding $\mathbf{z}_1^{k(\ell)}$. Given input symbol $X_{\ell+1}$, it appends $f(X_{\ell+1}, s_{\ell})$ to the codeword, and updates its state to $s_{\ell+1} = g(X_{\ell+1}, s_{\ell})$.

With an abuse of notation, denote by $f_{\ell}(\mathbf{X}^{\ell}, s_{\text{init}}) \in \{0, 1\}^*$ the binary sequence obtained by applying the finite state encoder to $\mathbf{X}^{\ell} = (X_1, \dots, X_{\ell})$. We say that the FS encoder is information lossless if for any $\ell \in \mathbb{N}$, $\mathbf{X}^{\ell} \mapsto f_{\ell}(\mathbf{X}^{\ell}, s_{\text{init}})$ is injective.

Theorem 5.4. *Let $\mathbf{X} = \mathbf{X}^{m,n} \sim \mathcal{T}(Q, q_r, q_c; m, n)$ and $\phi := (\Sigma, f, g)$ be an information lossless finite state encoder. Define the corresponding compression rate $R_{\phi}(\mathbf{X})$, as per Eq. (3.1). Assuming $m > 10$, $|\Sigma| \geq |\mathcal{X}|$, and $\log_2 |\Sigma| \leq n \log_2 |\mathcal{X}|/9$,*

$$\mathbb{E} R_{\phi}(\mathbf{X}) \geq \frac{H(X|U)}{\log_2 |\mathcal{X}|} - 10 \sqrt{\frac{\log |\Sigma|}{n \log |\mathcal{X}|}} \cdot \log(n \log |\Sigma|). \quad (5.5)$$

The asymmetry between U and V in the last statement (and below) arises because we assume that the table is serialized in row-major order. Of course the roles of U and V are exchanged if we use column major.

Remark 5.2. The leading term of the above lower bound is $H(X|U)/\log_2 |\mathcal{X}|$. Since conditioning reduces entropy, this is strictly larger than the ideal rate which is roughly $H(X|U, V)/\log_2 |\mathcal{X}|$, cf. Eq. (5.3).

The next term is negligible provided $\log |\Sigma| \ll n \log |\mathcal{X}|$. This condition is easy to interpret: it amounts to say that the finite state machine does not have enough states to memorize a row of the table $\mathbf{X}^{m,n}$.

The gap between $H(X|U, V)$ (appearing in the ideal rate of Lemma 5.1) and $H(X|U)$ (appearing in the last statement

and in the analysis of LZ encoders) can be illustrated by a toy example. Assume u_i, v_j are uniform in $\{0, 1\}$ and $X_{ij} = u_i + v_j \pmod{2}$. Then $H(X|U, V) = 0$ while $H(X|U) = H(X|V) = \log 2$. It is easy to compress very well by identifying the latents (just store the latents), but if we scan a single row (or column), we will only see a random sequence of bits.

5.2.2. LEMPEL-ZIV

The pseudocode of the Lempel-Ziv algorithm that we will analyze is given in Appendix F.

In words, after the first k characters of the input have been parsed, the encoder finds the longest string $X_k^{k+\ell-1}$ which appears in the past. It then encodes a pointer to the position of the earlier appearance of the string T_k , and its length L_k . If a symbol X_k never appeared in the past, we use a special encoding, cf. Appendix F.

We encode the pointer T_k in plain binary using $\lceil \log_2(N + |\mathcal{X}|) \rceil$ bits (note that $T_k \in \{-|\mathcal{X}| + 1, \dots, 1, \dots, N\}$), and L_k using an instantaneous prefix-free code, e.g. Elias δ -coding, taking $2\lceil \log_2 L_k \rceil + 1$ bits.

Assumption 5.5. There exist a constant $c_0 > 0$ such that

$$\max_{x \in \mathcal{X}} \max_{u, v \in \mathcal{L}} Q(x|u, v) \leq 1 - c_0.$$

Further $Q, q_r, c_0, \mathcal{X}, \mathcal{L}$ are fixed and $m, n \rightarrow \infty$ with $m = n^{\alpha+o(1)}$, i.e.

$$\lim_{n \rightarrow \infty} \frac{\log m}{\log n} = \alpha \in (0, \infty). \quad (5.6)$$

Theorem 5.6. Under Assumption 5.5, the asymptotic Lempel-Ziv rate is

$$\lim_{m, n \rightarrow \infty} \mathbb{E} R_{\text{LZ}}(\mathbf{X}^{m, n}) = R_{\text{LZ}}^\infty := \sum_{u \in \mathcal{L}} \frac{q_r(u) R_{\text{LZ}}^\infty(u)}{\log_2 |\mathcal{X}|}, \quad (5.7)$$

$$R_{\text{LZ}}^\infty(u) := H(X|U = u) \wedge \left(\frac{1 + \alpha}{\alpha} \right) H(X|U = u, V).$$

Remark 5.3. The asymptotics of the Lempel-Ziv rate is given by the minimum of two expressions, which correspond to different behaviors of the encoder. For $u \in \mathcal{L}$, define $\alpha_*(u) := H(X|U = u, V) / (H(X|U = u) - H(X|U = u, V))$ (with $\alpha_*(u) = \infty$ if $H(X|U = u) = H(X|U = u, V)$). Then:

If $\alpha < \alpha_*(u)$, then we are a ‘skinny table’ regime. The algorithm mostly deduplicates segments in rows with latent u by using strings in different rows but aligned in the same columns. If $\alpha > \alpha_*(u)$, then we are a ‘fat table’ regime. The algorithm mostly deduplicates segments on rows with latent u by using rows and columns that are not the same as the current segment.

Example 5.4 (Symmetric Binary Model, dense regime). Under the Symmetric Binary Model $\mathcal{T}_{\text{SBM}}(p_0, p_1, k; m, n)$ of Example 4.2, we can compute the optimal compression rate of Corollary 5.2, the finite state compression rate of Theorem 5.4, the Lempel-Ziv rate of Theorem 5.6.

If p_0, p_1 are of order one, and $m = n^{\alpha+o_n(1)}$ as $m, n \rightarrow \infty$, letting $\bar{p} := ((k-1)/k)p_0 + (1/k)p_1$, $\bar{h}(p_0, p_1) := ((k-1)/k)h(p_0) + (1/k)h(p_1)$, we obtain:

$$\mathbb{E} R_{\text{opt}}(\mathbf{X}) = \left(1 - \frac{1}{k}\right) h(p_0) + \frac{1}{k} h(p_1) + o_n(1),$$

$$\mathbb{E} R_{\text{fin. st.}}(\mathbf{X}) \geq h(\bar{p}) + o_n(1),$$

$$\mathbb{E} R_{\text{LZ}}(\mathbf{X}) = h(\bar{p}) \wedge \left(\frac{1 + \alpha}{\alpha} \right) \bar{h}(p_0, p_1) + o_n(1).$$

These theoretical predictions are used to trace the contour lines in Figure 1. (ANS coding is implemented as a finite state code here.)

5.3. Practical latent-based compression

Achieving the ideal rate of Corollary 5.2 via arithmetic or Huffman coding requires to compute the probability $P(\mathbf{X}^{m, n})$, which is intractable. We will next show that we can achieve a compression rate that is close to the ideal rate via latents estimation.

We begin by considering general latents estimators $\hat{u} : \mathcal{X}^{m \times n} \rightarrow \mathcal{L}^m$, $\hat{v} : \mathcal{X}^{m \times n} \rightarrow \mathcal{L}^n$. We measure their accuracy by the error (cf. Lemma 5.1)

$$\text{Err}_U(\mathbf{X}; \hat{u}) := \frac{1}{m} \min_{\pi \in \mathfrak{S}_{\mathcal{L}}} \sum_{i=1}^m \left\{ \mathbf{1}_{\hat{u}_i(\mathbf{X}) \neq \pi(u_i)} \right\}$$

and the analogous $\text{Err}_V(\mathbf{X}; \hat{v})$. Here the minimization is over the set $\mathfrak{S}_{\mathcal{L}}$ of permutations of the latents alphabet \mathcal{L} .

We can use any estimators \hat{u}, \hat{v} to reorder rows and columns and compress the table $\mathbf{X}^{m, n}$ according to the algorithm described in the introduction. We denote by $R_{\text{lat}}(\mathbf{X})$ the compression rate achieved by such a procedure.

Our first result implies that, if the latent estimators are consistent (namely, they recover the true latents with high probability, up to permutations), then the resulting rate is close to the ideal one.

Lemma 5.7. Assume data distributed according to model $\mathbf{X}^{m, n} \sim \mathcal{T}(Q, q_r, q_c; m, n)$, with $m, n \geq \log_2 |\mathcal{L}|$. Further assume there exists $c_0 > 0$ such that $q_r(u), q_c(v) \geq c_0$ for all $u, v \in \mathcal{L}$. Let $R_{\text{lat}}(\mathbf{X})$ be the rate achieved by the latent-based scheme with latents estimators \hat{u}, \hat{v} , and base encoders $Z_{\mathcal{X}} = Z_{\mathcal{L}} = Z$. Then

$$\begin{aligned} \mathbb{E} R_{\text{lat}}(\mathbf{X}) &\leq \frac{H(\mathbf{X}^{m, n})}{mn \log_2 |\mathcal{X}|} + 2P_{\text{err}}(m, n) + \frac{4 \log(mn)}{mn} \\ &\quad + |\mathcal{L}|^2 \Delta_Z(c \cdot mn; \mathcal{L}) + 2\Delta_Z(m \wedge n; \{q_r, q_c\}). \end{aligned} \quad (5.8)$$

Here $P_{\text{err}}(m, n) := \mathbb{P}(\text{Err}_U(\mathbf{X}^{m,n}; \hat{\mathbf{u}}) > 0) + \mathbb{P}(\text{Err}_V(\mathbf{X}^{m,n}; \hat{\mathbf{v}}) > 0)$, $\Delta_Z(N; \mathcal{P}_*)$ is the worst-case redundancy of encoder Z over i.i.d. sources with distributions in \mathcal{P}_* (see comments below), $\mathcal{Q} := \{Q(\cdot | u, v)\}_{u,v \in \mathcal{L}}$.

The redundancies of Lempel-Ziv, frequency-based arithmetic coding and ANS coding can be upper bounded as (in the last bound Q, q_r, q_c need to be independent of N)

$$\Delta_{\text{LZ}}(N; \mathcal{P}_*) \leq 40c_*(\mathcal{P}_*) \left(\frac{\log \log N}{\log N} \right)^{1/2}, \quad (5.9)$$

$$\Delta_{\text{AC}}(N; \mathcal{P}_*) \leq \frac{2|\mathcal{X}|}{\log |\mathcal{X}|} \cdot \frac{\log N}{N}, \quad (5.10)$$

$$\Delta_{\text{ANS}}(N; \mathcal{P}_*) \leq \frac{2|\mathcal{X}| \log N + C|\mathcal{X}|}{N}. \quad (5.11)$$

Here Eq. (5.9) holds for $N \geq \exp\{\sup_{q \in \mathcal{P}_*} (4 \log(2/H(q)))^2\}$, and $c_*(\mathcal{P}_*) := \sup_{q \in \mathcal{P}_*} \sum_{x \in \mathcal{X}} (\log q(x))^2 / |\mathcal{X}|$.

The proof of this lemma is given in Appendix G.1. The main content of the lemma is in the general bound (5.8) which is proven in Appendix G.1.1.

Remark 5.5. We define the worst case redundancy $\Delta_Z(N_0; \mathcal{P}_*) := \max_{N \geq N_0} \hat{\Delta}_Z(N; \mathcal{P}_*)$, where

$$\hat{\Delta}_Z(N; \mathcal{P}_*) := \max_{q \in \mathcal{P}_*} \frac{\mathbb{E}_q \text{len}(Z(\mathbf{Y}^N)) - H(\mathbf{Y}^N)}{N \log_2 k}, \quad (5.12)$$

where $\mathcal{P}_* \subseteq \mathcal{P}([k]) := \{(p_i)_{i \leq k} \in \mathbb{R}^k : p_i \geq 0 \forall i \text{ and } \sum_{i \leq k} p_i = 1\}$ is a set of probability distributions over $[k]$ and \mathbf{Y}^N is a vector with i.i.d. entries $Y_i \sim q$.

While Eqs. (5.9)—(5.11) are closely related to well known facts, there are nevertheless differences with respect to statements in the literature. We address them in Section G.1.2. Perhaps the most noteworthy difference is in the bound (5.9) for the LZ algorithm. Existing results, e.g. Theorem 2 in (Savari, 1998), assume a single, N -independent, distribution q and are asymptotic in nature. Equation (5.9) is a non-asymptotic statement and applies to a collection of distributions \mathcal{P}_* that could depend on N .

Lemma 5.7 can be used in conjunction with any latent estimation algorithm, as we next demonstrate by considering the spectral algorithm of Section 2.2. Recall that the algorithm makes use of a map $\psi : \mathcal{X} \rightarrow \mathbb{R}$. For $(X, U, V) \sim Q(\cdot | \cdot, \cdot) q_r(\cdot) q_c(\cdot)$, we define $\bar{\psi}(u, v) := \mathbb{E}[\psi(X) | U = u, V = v]$, $\Psi := (\bar{\psi}(u, v))_{u,v \in \mathcal{L}}$ and the parameters:

$$\mu_n := \sigma_{\min}(\Psi), \quad \nu_n := \max_{u,v \in \mathcal{L}} |\bar{\psi}(u, v)|, \quad (5.13)$$

$$\sigma_n^2 := \max_{u,v \in \mathcal{L}} \text{Var}(\psi(X) | U = u, V = v). \quad (5.14)$$

We further will assume, without loss of generality $\max_{x \in \mathcal{X}} |\psi(x)| \leq 1$.

Finally, we need to formally specify the version of the KMeans primitive in the spectral clustering algorithm. In fact, we establish correctness for a simpler thresholding procedure. Considering to be definite the row latents, and for a given threshold $\theta > 0$, we construct a graph $G_\theta = ([m], E_\theta)$ by letting (for distinct $i, j \in [m]$)

$$\frac{\|\mathbf{a}_i - \mathbf{a}_j\|_2}{(\|\mathbf{a}_i\|_2 + \|\mathbf{a}_j\|_2)/2} \leq \theta \Leftrightarrow (i, j) \in E_\theta. \quad (5.15)$$

The algorithm then output the connected components of G_θ .

Theorem 5.8. Assume data $\mathbf{X}^{m,n} \sim \mathcal{T}(Q, q_r, q_c; m, n)$, with $m, n \geq \log_2 |\mathcal{L}|$ and $\min_{u \in \mathcal{L}} (q_r(u) \wedge q_c(u)) \geq c_0$ for a constant $c_0 > 0$. Let $R_{\text{lat}}(\mathbf{X})$ be the rate achieved by the latent-based scheme with spectral latents estimators $\hat{\mathbf{u}}, \hat{\mathbf{v}}$, base compressors $Z_{\mathcal{X}} = Z_{\mathcal{L}} = Z$, and thresholding algorithm as described above. Then, assuming $\sigma_n \geq c\sqrt{(\log n)/n}$, $\nu_n/\sigma_n \leq c\sqrt{\log n}$, $\mu_n \geq C(\sigma_n \sqrt{(\log n)/m} \vee (\log n)/\sqrt{mn})$, $\theta \leq \sqrt{c_0}/100$, we have

$$\mathbb{E} R_{\text{lat}}(\mathbf{X}) \leq \frac{H(\mathbf{X}^{m,n})}{mn \log_2 |\mathcal{X}|} + \frac{10 \log(mn)}{mn} + |\mathcal{L}|^2 \Delta_Z(c \cdot mn; \mathcal{Q}) + 2\Delta_Z(m \wedge n; \{q_r, q_c\}).$$

We focus on the simpler thresholding algorithm of Eq. (5.15) instead of KMeans in order to avoid technical complications that are not the main focus of this paper. We expect it to be relatively easy to generalize this result, e.g. using the results of (Makarychev et al., 2020) for KMeans++.

Example 5.6. Consider the Symmetric Binary Model $\mathcal{T}_{\text{SBM}}(p_0, p_1, k; m, n)$ of Example 4.2, with $p_0 = p_{0,n}$, $p_1 = p_{1,n}$ potentially dependent on n . Since in this case $\mathcal{X} = \{0, 1\}$ the choice of the map ψ has little impact and we set $\psi(x) = x$. We assume, to simplify formulas, $|p_{1,n} - p_{0,n}| \leq k p_{0,n}, p_{1,n} \vee p_{0,n} \leq 9/10$. It is easy to compute $\mu_n = |p_{1,n} - p_{0,n}|, \nu_n = p_{0,n} \vee p_{1,n}, \sigma_n^2 \asymp p_{1,n} \vee p_{0,n}$: Theorem 5.8 implies nearly optimal compression rate under the following conditions on the model parameters:

$$p_{1,n} \vee p_{0,n} \gtrsim \sqrt{\frac{\log n}{n}}, \quad |p_{1,n} - p_{0,n}| \gtrsim \frac{\log n}{\sqrt{mn}},$$

$$\frac{|p_{1,n} - p_{0,n}|}{p_{1,n} \vee p_{0,n}} \gtrsim \sqrt{\frac{\log n}{n}}.$$

Here \gtrsim hides factors depending on k, c_0 . The last of these condition amounts to requiring that the signal-to-noise ratio is large enough to consistently reconstruct the latents. In the special case of square symmetric matrices ($m = n$), sharp constants in these bounds can be derived from (Abbe, 2017).

Method	DRR ($k = 4$)	($k = 6$)	($k = 10$)
Oracle	0.27	0.27	0.27
No Latent	0.06	0.11	0.17
$k_{\text{alg}} = k - 2$	0.14	0.21	0.25
$k_{\text{alg}} = k - 1$	0.21	0.25	0.26
$k_{\text{alg}} = k$	0.27	0.27	0.26
$k_{\text{alg}} = k + 1$	0.27	0.27	0.26
$k_{\text{alg}} = k + 2$	0.25	0.27	0.23

Method	Taxi	Census	Forest	Jokes
ZSTD 22	0.53 sec	146 sec	17.2 sec	2.8 sec
Latents + LZ	1.2 sec	96 sec	5.2 sec	3.8 sec
Subsampling	0.7 sec	30 sec	3.7 sec	1.5 sec

Table 2: Left: DRR of the latent-based compressor in the SBM with miss-specified latent space size k_{alg} . Right: average runtime (seconds) of latent-based compressors on several datasets, compared with state of the art.

6. Discussion and extensions

We proved that classical lossless compression schemes, that serialize the data and then apply a finite state encoder or a Lempel-Ziv encoder to the resulting sequence are sub-optimal when applied to tabular data. Namely, we introduced a simple model for tabular data, and made the following novel contributions:

1. We characterized the optimal compression rate under this model.
2. We rigorously quantified the gap in compression rate suffered by classical compressors.
3. We showed that a compression scheme that estimates the latents performs well in practice and provably achieves optimal rate on our model.

The present work naturally suggests several questions and directions for future work.

Model miss-specification. The numerical simulations of Section 4 were carried out within the Symmetric Binary Model, under the assumption that the correct cardinality of the latents space, k , is known at the encoder. While the rigorous guarantees of Section 5.3 are more general, it is important to understand the effect of model misspecification.

Table 2 presents initial evidence of the robustness of the proposed approach. We generate tables $\mathbf{X}^{m,n} \sim \mathcal{T}_{\text{SBM}}(p_0, p_1, k; m, n)$, with $p_0 = 0.2$, $p_1 = 0.8$, and dimensions $m = n = 1000$. We vary $k \in \{4, 6, 8\}$ and run the latent-based compressor in a miss-specified setting by using a latent space size $k_{\text{alg}} \neq k$. We observe that the resulting DRR is fairly robust to miss-specification in the the number of latents.

It is also worth pointing out that parameters of the latents model, such as k_{alg} , can be chosen at compression time to optimize DRR. However this implies severe performance losses, and therefore further investigation into model misspecification is warranted.

Computational efficiency. While we did not attempt to develop a highly optimized compressor, we believe that the

present approach is amenable to such a development. The largest overhead over standard compressors is in the clustering step, and in particular computing the singular value decomposition (SVD) of the data. This can be accelerated using methods from randomized linear algebra. We implemented row subsampling SVD (Drineas et al., 2006), and observed essentially no loss in DRR by using 10% of the rows as compared to full SVD.

Table 2 reports running time experiments to compare the original latent based compressor, its subsampling version and ZSTD at its strongest compression level. Runtimes were averaged over 5 runs on a Macbook Pro single-threaded with a 2 GHz 4-core Intel i5 chip.

We observe that row-subsampling yields significant performance improvements, and runtimes that compare well with the industry state of the art.

Bits back coding. As mentioned several times, encoding the latents is sub-optimal, unless these can be estimated accurately in the sense of $\mathbb{E} \text{Err}_U(\mathbf{X}; \hat{\mathbf{u}}) \rightarrow 0$, $\mathbb{E} \text{Err}_V(\mathbf{X}; \hat{\mathbf{v}}) \rightarrow 0$, cf. Lemma 5.1. If this is not the case, then optimal rates can be achieved using bits-back coding.

Continuous latents. Of course, using discrete latents is somewhat un-natural, and it would be interesting to consider continuous ones, in which case bits-back coding is required.

Multi-way tables (tensors). A natural extension of the current work is to order- s multiway tables or (equivalently for our purposes) tensors $\mathbf{X} \in \mathcal{X}^{n_1 \times \dots \times n_s}$. The basic scheme would essentially be the same: estimate s vector of latents $v_i \in \mathcal{L}^{n_i}$, $i \leq s$ and use them to partition rows/columns.

Acknowledgements

This work was carried out while Andrea Montanari was on leave from Stanford and a Chief Scientist at Ndata Inc dba Project N. The present research is unrelated to AMs Stanford activity.

References

- Abbe, E. Community detection and stochastic block models: recent developments. *The Journal of Machine Learning Research*, 18(1):6446–6531, 2017.
- Abbe, E., Fan, J., Wang, K., and Zhong, Y. Entrywise eigenvector analysis of random matrices with low expected rank. *Annals of statistics*, 48(3):1452, 2020.
- Altman, E. R. Synthesizing credit card transactions, 2019. <https://arxiv.org/abs/1910.03033>.
- Besta, M. and Hoefler, T. Survey and taxonomy of lossless graph compression and space-efficient graph representations. *arXiv preprint arXiv:1806.01799*, 2018.
- Boldi, P. and Vigna, S. The webgraph framework i: compression techniques. In *Proceedings of the 13th international conference on World Wide Web*, pp. 595–602, 2004.
- Chen, P., Yu, H.-F., Dhillon, I., and Hsieh, C.-J. Drone: Data-aware low-rank compression for large nlp models. *Advances in neural information processing systems*, 34: 29321–29334, 2021.
- Cheng, H., Gimbutas, Z., Martinsson, P.-G., and Rokhlin, V. On the compression of low rank matrices. *SIAM Journal on Scientific Computing*, 26(4):1389–1404, 2005.
- Chierichetti, F., Kumar, R., Lattanzi, S., Mitzenmacher, M., Panconesi, A., and Raghavan, P. On compressing social networks. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pp. 219–228, 2009.
- Collet, Y. and Kucherawy, M. Zstandard compression and the application/zstd media type. Technical report, 2018.
- Cover, T. M. and Thomas, J. A. *Elements of information theory*. Wiley, 2006.
- Drineas, P., Kannan, R., and Mahoney, M. W. Fast monte carlo algorithms for matrices ii: Computing a low-rank approximation to a matrix. *SIAM Journal on computing*, 36(1):158–183, 2006.
- Dua, D. and Graff, C. UCI machine learning repository, 2017. <http://archive.ics.uci.edu/ml>.
- Duda, J. Asymmetric numeral systems. *arXiv:0902.0271*, 2009.
- Duda, J. Asymmetric numeral systems: entropy coding combining speed of huffman coding with compression rate of arithmetic coding. *arXiv preprint arXiv:1311.2540*, 2013.
- Frey, B. J. and Hinton, G. E. Free energy coding. In *Proceedings of Data Compression Conference-DCC’96*, pp. 73–81. IEEE, 1996.
- Goldberg, K., Roeder, T., Gupta, D., and Perkins, C. Jester Datasets for Recommender Systems and Collaborative Filtering Research. <https://eigentaste.berkeley.edu/dataset/>.
- Goldberg, K., Roeder, T., Gupta, D., and Perkins, C. Eigentaste: A constant time collaborative filtering algorithm. *Information Retrieval*, 4(2):133–151, 2001. <https://doi.org/10.1023/A:1011419012209>.
- Hinton, G. E. and Zemel, R. Autoencoders, minimum description length and helmholtz free energy. *Advances in neural information processing systems*, 6, 1993.
- Holland, P. W., Laskey, K. B., and Leinhardt, S. Stochastic blockmodels: First steps. *Social networks*, 5(2):109–137, 1983.
- Hou, J., Chau, L.-P., Magnenat-Thalmann, N., and He, Y. Sparse low-rank matrix approximation for data compression. *IEEE Transactions on Circuits and Systems for Video Technology*, 27(5):1043–1054, 2015.
- IBM. Stats NZ Business Price indexes: March 2022 quarter. <https://www.stats.govt.nz/information-releases/business-price-indexes-march-2022-quarter/>, March 2022.
- Kosolobov, D. The efficiency of the ans entropy encoding. *arXiv preprint arXiv:2201.02514*, 2022.
- Leskovec, J. and Krevl, A. SNAP Datasets: Stanford large network dataset collection, June 2014. <http://snap.stanford.edu/data>.
- Li, N. and Li, B. Tensor completion for on-board compression of hyperspectral images. In *2010 IEEE International Conference on Image Processing*, pp. 517–520. IEEE, 2010.
- Lim, Y., Kang, U., and Faloutsos, C. Slashburn: Graph compression and mining beyond caveman communities. *IEEE Transactions on Knowledge and Data Engineering*, 26(12):3077–3089, 2014.
- Makarychev, K., Reddy, A., and Shan, L. Improved guarantees for k-means++ and k-means++ parallel. *Advances in Neural Information Processing Systems*, 33:16142–16152, 2020.
- NYC.gov. TLC Trip Record Data: NYC Taxi & Limousine Commission, January 2022.
- Phan, A.-H., Sobolev, K., Sozykin, K., Ermilov, D., Gusak, J., Tichavský, P., Glukhov, V., Oseledets, I., and Cichocki, A. Stable low-rank tensor decomposition for compression

of convolutional neural network. In *European Conference on Computer Vision*, pp. 522–539. Springer, 2020.

Salomon, D. *Data compression: the complete reference*. Springer Science & Business Media, 2004.

Savari, S. A. Redundancy of the lempel-ziv string matching code. *IEEE Transactions on Information Theory*, 44(2): 787–791, 1998.

stats.govt.nz. Stats NZ Business Price indexes: March 2022 quarter. <https://www.stats.govt.nz/information-releases/business-price-indexes-march-2022-quarter/>, March 2022.

Taylor, P. R. Lossless compression of wave function information using matrix factorization: A gzip for quantum chemistry. *The Journal of Chemical Physics*, 139(7): 074113, 2013.

Townsend, J., Bird, T., and Barber, D. Practical lossless compression with latent variables using bits back coding. In *7th International Conference on Learning Representations, ICLR 2019*, volume 7. International Conference on Learning Representations (ICLR), 2019a.

Townsend, J., Bird, T., Kunze, J., and Barber, D. Hilloc: lossless image compression with hierarchical latent variable models. In *International Conference on Learning Representations*, 2019b.

Yuan, Z. and Oja, E. Projective nonnegative matrix factorization for image compression and feature extraction. In *Scandinavian Conference on Image Analysis*, pp. 333–342. Springer, 2005.

A. Details on the empirical evaluation

A.1. Datasets

We used the following datasets:

- *Taxicab*. A table with $m = 62,495$ rows, $n_0 = 20$ columns comprising data for taxi rides in NYC during January 2022 (NYC.gov, 2022). After preprocessing this table has $n = 18$ columns. For the LZ (ZSTD) compressor we used $|\mathcal{L}_r| = 9$ row latents and 15 column latents, for the ANS compressor we used $|\mathcal{L}_r| = 5$ row latents and 14 column latents.
- *Network*. Four social networks from SNAP Datasets, representing either friends as undirected edges for Facebook or directed following relationships on Google Plus (Leskovec & Krevl, 2014). We regard these as four distinct tables with 0 – 1 entries, with dimensions, respectively $m = n \in \{333, 747, 786, 1187\}$. For each table we used 5 row latents and 5 column latents.
- *Card transactions*. A table of simulated credit card transactions containing information like card ID, merchant city, zip, etc. This table has $m = 24,386,900$ rows and $n_0 = 15$ columns and was generated as described in (Altman, 2019) and downloaded from (IBM, 2022). After preprocessing the table has $n = 12$ columns. For this table we used 3 row latents and n column latents.
- *Business price index*. A table of the values of the consumer price index of various goods in New Zealand between 1996 and 2022. This is a table with $m = 72,750$ rows and $n_0 = 12$ columns from the Business price indexes: March 2022 quarter - CSV file from (stats.govt.nz, 2022). After preprocessing this table has $n = 10$ columns. Due to the highly correlated nature of consecutive rows, we first shuffle them before compressing. For the LZ method we used 6 row latents and 7 column latents, for the ANS method we used 2 row latents and 6 column latents.
- *Forest*. A table from the UCI data repository comprising $m = 581,011$ cartographic measurements with $n_0 = 55$ attributes, to predict forest cover type based on information gathered from US Geological Survey (Dua & Graff, 2017). It contains binary qualitative variables, and some continuous values like elevation and slope. After preprocessing this data has $n = 55$ columns. For the LZ method we used 6 row latents and 17 column latents, for the ANS method we used 6 row latents and 17 column latents.
- *US Census*. Another table from the UCI Machine Learning Repository (Dua & Graff, 2017) with $m = 2,458,285$ and $n_0 = 68$ categorical attributes related to demographic information, income, and occupation information. After preprocessing this data has $n = 68$ columns. For this data we used 9 row latents and n column latents.
- *Jokes*. A table containing ratings of a series of jokes by 24,983 users collected between April 1999 and May 2003 (Goldberg et al., 2001; Goldberg et al.). These ratings are real numbers on a scale from -10 to 10 , and a value of 99 is given to jokes that were not rated. There are $m = 23,983$ rows and $n_0 = 101$. The first column identifies how many jokes were rated by a user, and the rest of the columns contain the ratings. After preprocessing this data has $n = 101$ columns, all quantized. For the LZ method we used 2 row latents and n column latents, for the ANS method we used 8 row latents and 8 column latents.

A.2. Preprocessing

We preprocessed different columns as follows:

- If a column comprises $K \leq 256$ unique values, then we map the values to $\{0, \dots, K - 1\}$.
- If a column is numerical and comprises more than 256 unique values, we calculate the quartiles for the data and map each entry to its quartile membership (0 for the lowest quartile, 1 for the next largest, 2 for the next and 3 for the largest).
- If a column does not meet either of the above criteria, we discard it.

Finally, in some experiments we randomly permuted before compression. The rationale is that some of the above datasets have rows already ordered in a way that makes nearby rows highly correlated. In these cases, row reordering is –obviously– of limited use.

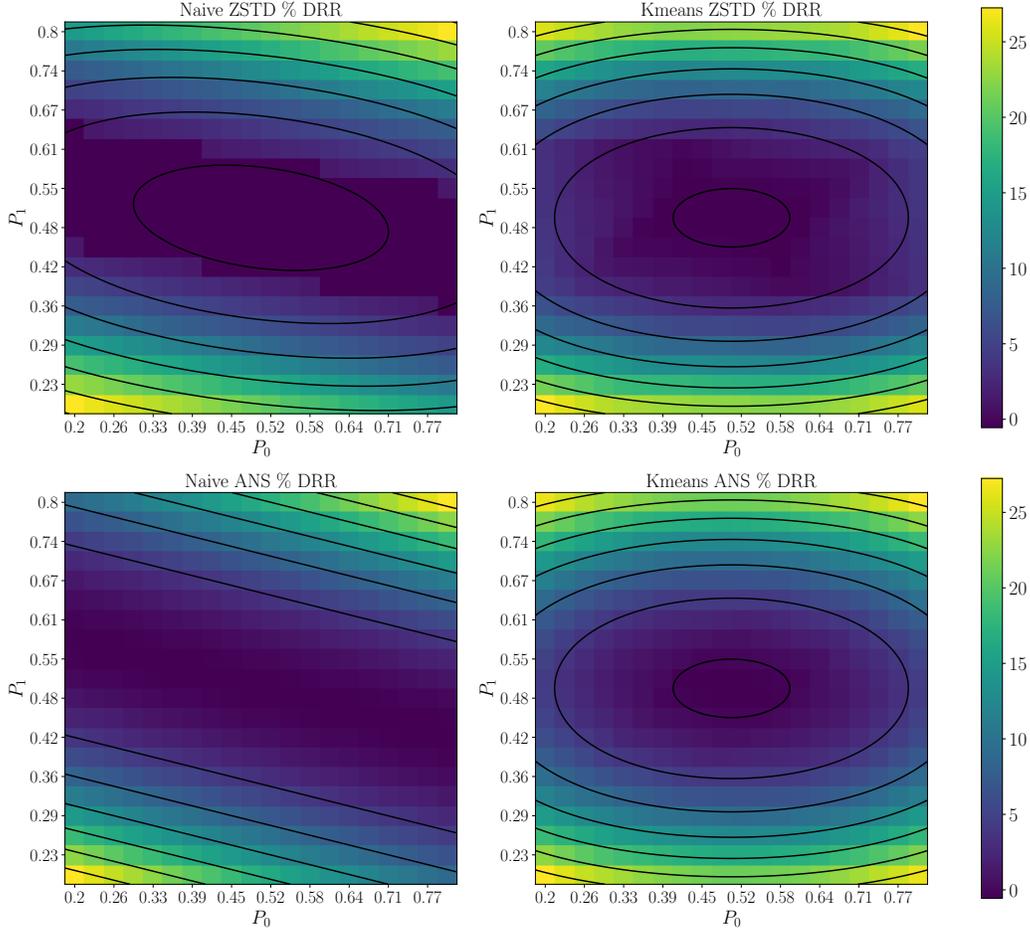


Figure 2: Comparing data reduction rate of naive coding and latent-based coding for data from SBM with $k = 5$ latents. Top row: ZSTD base compressor. Bottom row: ANS based compressor. Contour lines correspond to the theoretical predictions for various compression algorithms (cf. Example 5.4).

B. Further simulations

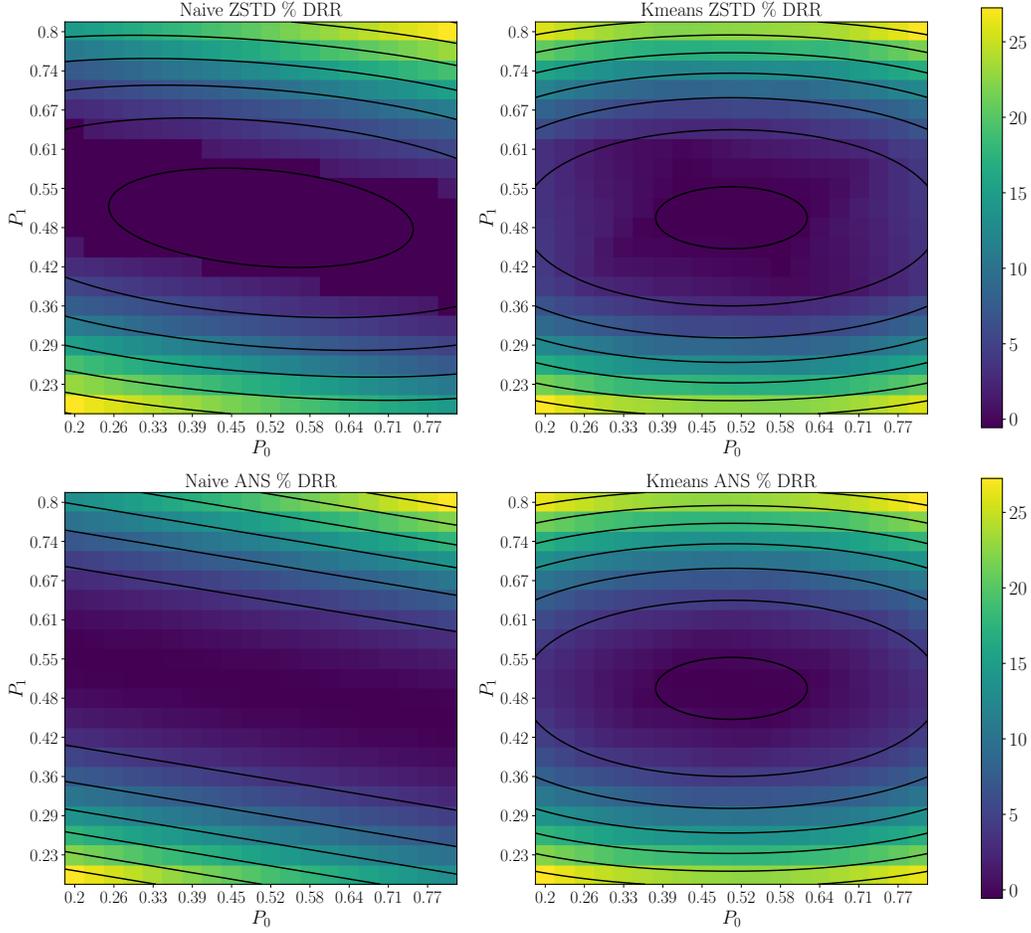
Figures 2 and 3 report empirical DRR values for ZSTD and ANS coding, for data generated according to the symmetric model $\mathcal{T}_{\text{SBM}}(p_0, p_1, k; m, n)$ of Section 4. We use $m = n = 1000$ as before, but now $k \in \{5, 7\}$. Results confirm the conclusions of Section 4.

C. A basic fact

Lemma C.1. *Let \mathcal{A} be a finite set and $F : \mathcal{A} \rightarrow \{0, 1\}^*$ be an injective map. Then, for any probability distribution p over \mathcal{A} ,*

$$\sum_{a \in \mathcal{A}} p(a) \text{len}(F(a)) \geq H(p) - \log_2 \log_2(|\mathcal{A}| + 2). \quad (\text{C.1})$$

Proof. Assume without loss of generality that $\mathcal{A} = \{1, \dots, M\}$, with $|\mathcal{A}| = K$, and that the elements of \mathcal{A} have all non-vanishing probability and are ordered by decreasing probability $p_1 \geq p_2 \geq \dots \geq p_K > 0$. Let $N_\ell := 2 + 4 + \dots + 2^\ell = 2\ell + 1 - 2$. Then the expected length is minimized any map F such that $\text{len}(F(a)) = \ell$ for $N_{\ell-1} \leq a \leq N_\ell$ with the


 Figure 3: Same as Figure 2 with $k = 7$.

maximum length ℓ_K being defined by $N_{\ell_K-1} < K \leq N_{\ell_K}$. For $A \sim p$, $L := \text{len}(F(A))$, we have

$$\begin{aligned}
 H(p) &:= H(A) \stackrel{(a)}{\leq} H(L) + H(A|L) \\
 &\leq \log_2 \ell_K + \sum_{\ell=1}^{\ell_K} \mathbb{P}(L = \ell) H(A|L = \ell) \\
 &\stackrel{(b)}{\leq} \log_2 \ell_K + \sum_{\ell=1}^{\ell_M} \mathbb{P}(L = \ell) \ell \\
 &\leq \log_2 \log_2(K+2) + \sum_{a \in \mathcal{A}} p(a) \text{len}(F(a)),
 \end{aligned}$$

where (a) is the chain rule of entropy and (b) follows because by injectivity, given $\text{len}(F(A)) = \ell$, A can take at most 2^ℓ values. \square

D. Proofs of results on ideal compression

D.1. Proof of Lemma 5.1

We begin by claiming that

$$\frac{1}{mn} H(\mathbf{X}^{m,n}) = H(X|U, V) + \frac{1}{mn} I(\mathbf{X}^{m,n}; \mathbf{U}^m, \mathbf{V}^n). \quad (\text{D.1})$$

Indeed, by the definition of mutual information, we have $H(\mathbf{X}_{m,n}) = H(\mathbf{X}_{m,n}|\mathbf{U}_m, \mathbf{V}_n) + I(\mathbf{X}_{m,n}; \mathbf{U}_m, \mathbf{V}_n)$. Equation (D.1) follows by noting that

$$\begin{aligned}
 H(\mathbf{X}^{m,n}|\mathbf{U}^m, \mathbf{V}^n) &= \sum_{\mathbf{u} \in \mathcal{L}^m} \sum_{\mathbf{v} \in \mathcal{L}^n} \mathbb{P}(\mathbf{U}^m = \mathbf{u}, \mathbf{V}^n = \mathbf{v}) H(\mathbf{X}^{m,n}|\mathbf{U}^m = \mathbf{u}, \mathbf{V}^n = \mathbf{v}) \\
 &\stackrel{(a)}{=} \sum_{i=1}^m \sum_{j=1}^n \sum_{\mathbf{u} \in \mathcal{L}^m} \sum_{\mathbf{v} \in \mathcal{L}^n} \mathbb{P}(\mathbf{U}^m = \mathbf{u}, \mathbf{V}^n = \mathbf{v}) H(X_{i,j}|U_i = u, V_j = v) \\
 &= \sum_{i=1}^m \sum_{j=1}^n \sum_{\mathbf{u} \in \mathcal{L}} \sum_{\mathbf{v} \in \mathcal{L}} \mathbb{P}(U_i = u_i, V_j = v_j) H(X_{i,j}|U_i = u_i, V_j = v_j) \\
 &= \sum_{i=1}^m \sum_{j=1}^n H(X_{i,j}|U_i, V_j) \stackrel{(b)}{=} mnH(X_{1,1}|U_1, V_1),
 \end{aligned}$$

where (a) follows from the fact that the $(X_{i,j})$ are conditionally independent given $\mathbf{U}^m, \mathbf{V}^n$, and since the conditional distribution of $X_{i,j}$ only depends on $\mathbf{U}^m, \mathbf{V}^n$ via U_i, V_j ; (b) holds because the triples $(X_{i,j}, U_i, V_j)$ are identically distributed.

The lower bound in Eq. (5.1) holds because mutual information is non-negative, and the upper bound because $I(\mathbf{X}^{m,n}; \mathbf{U}^m, \mathbf{V}^n) \leq H(\mathbf{U}^m, \mathbf{V}^n) = mH(U_1) + nH(V_1)$.

Finally, to prove Eq. (5.2), define

$$\pi_{\mathbf{U}, \mathbf{X}} := \arg \min_{\pi \in \mathfrak{S}_{\mathcal{L}}} \frac{1}{m} \sum_{i=1}^m \mathbf{1}_{\hat{u}_i \neq \pi(u_i)}, \quad \pi_{\mathbf{V}, \mathbf{X}} := \arg \min_{\pi \in \mathfrak{S}_{\mathcal{L}}} \frac{1}{n} \sum_{i=1}^n \mathbf{1}_{\hat{v}_i \neq \pi(v_i)}, \quad (\text{D.2})$$

If the minimizer is not unique, one can be chosen arbitrarily. We then have holds because

$$\begin{aligned}
 I(\mathbf{X}^{m,n}; \mathbf{U}^m, \mathbf{V}^n) &= H(\mathbf{U}^m, \mathbf{V}^n) - H(\mathbf{U}^m, \mathbf{V}^n|\mathbf{X}^{m,n}) \\
 &\geq mH(U_1) + nH(V_1) - H(\mathbf{U}^m|\mathbf{X}^{m,n}) - H(\mathbf{V}^n|\mathbf{X}^{m,n}) \\
 &\geq mH(U_1) + nH(V_1) - [H(\mathbf{U}^m|\mathbf{X}^{m,n}, \pi_{\mathbf{U}, \mathbf{X}}) + I(\mathbf{U}^m; \pi_{\mathbf{U}, \mathbf{X}}|\mathbf{X}^{m,n})] \\
 &\quad - [H(\mathbf{V}^n|\mathbf{X}^{m,n}, \pi_{\mathbf{V}, \mathbf{X}}) + I(\mathbf{V}^n; \pi_{\mathbf{V}, \mathbf{X}}|\mathbf{X}^{m,n})] \\
 &\geq mH(U_1) + nH(V_1) - H(\mathbf{U}^m|\mathbf{X}^{m,n}, \pi_{\mathbf{U}, \mathbf{X}}) - H(\mathbf{V}^n|\mathbf{X}^{m,n}, \pi_{\mathbf{V}, \mathbf{X}}) - 2|\mathcal{L}| \log_2(|\mathcal{L}|), \quad (\text{D.3})
 \end{aligned}$$

where in the last inequality we used the fact that $I(\mathbf{U}^m; \pi_{\mathbf{U}, \mathbf{X}}|\mathbf{X}^{m,n}) \leq H(\pi_{\mathbf{U}, \mathbf{X}}) \leq \log_2(|\mathcal{L}|)$.

Now, consider the term $H(\mathbf{U}^m|\mathbf{X}^{m,n}, \pi_{\mathbf{U}, \mathbf{X}})$. Letting $\mathbf{Y} := (\mathbf{X}^{m,n}, \pi_{\mathbf{U}, \mathbf{X}})$ the stated accuracy assumption implies that there exists an estimator $\hat{\mathbf{u}}^+ = \hat{\mathbf{u}}^+(\mathbf{Y})$ such that

$$\varepsilon_U = \frac{1}{m} \sum_{i=1}^m \varepsilon_{U,i}, \quad \varepsilon_{U,i} := \mathbb{P}(\hat{u}_i^+(\mathbf{Y}) \neq U_i).$$

By Fano's inequality

$$\begin{aligned}
 H(\mathbf{U}^m|\mathbf{X}^{m,n}, \pi_{\mathbf{U}, \mathbf{X}}) &\leq \sum_{i=1}^m H(U_i|\mathbf{X}^{m,n}, \pi_{\mathbf{U}, \mathbf{X}}) \\
 &\leq \sum_{i=1}^m [\mathfrak{h}(\varepsilon_{U,i}) + \varepsilon_{U,i} \log(|\mathcal{L}| - 1)] \\
 &\leq m[\mathfrak{h}(\varepsilon_U) + \varepsilon_U \log(|\mathcal{L}| - 1)],
 \end{aligned}$$

where the last step follows by Jensen's inequality. The claim (5.2) follows by substituting this bound in Eq. (D.4) and using a similar bound for $H(\mathbf{V}^n|\mathbf{X}^{m,n}, \pi_{\mathbf{V}, \mathbf{X}})$.

D.2. Proof of Lemma 5.3

We begin with a technical fact.

Lemma D.1. *Let $\xi = (\xi_{ij})_{i \leq m, j \leq n}$, $\sigma = (\sigma_i)_{i \leq m}$, $\tau = (\tau_j)_{j \leq n}$ be collections of mutually independent random variables taking values in a measurable space \mathcal{Z} . $x : \mathcal{Z}^3 \rightarrow \mathcal{Z}$, $F : \mathcal{Z}^{m \times n} \rightarrow \mathbb{R}$. Define $\mathbf{x}(\xi, \sigma, \tau) \in \mathbb{R}^{m \times n}$ via $\mathbf{x}(\xi, \sigma, \tau)_{ij} = x(\xi_{ij}, \sigma_i, \tau_j)$.*

Given a vector of independent random variables \mathbf{z} , we let $\text{Var}_{z_i}(f(\mathbf{z})) := \mathbb{E}_{z_i}[(f(\mathbf{z}) - \mathbb{E}_{z_i}f(\mathbf{z}))^2]$. Define the quantities

$$B_* := \max_{\substack{\mathbf{x}, \mathbf{x}' \in \mathcal{Z}^{m \times n} \\ d(\mathbf{x}, \mathbf{x}') \leq 1}} |F(\mathbf{x}) - F(\mathbf{x}')|, \quad (\text{D.5})$$

$$B_1 := \max_{\substack{\sigma, \sigma' \in \mathcal{Z}^m \\ d(\sigma, \sigma') \leq 1}} \max_{\tau \in \mathcal{Z}^n} |\mathbb{E}_\xi F(\mathbf{x}(\xi, \sigma, \tau)) - \mathbb{E}_\xi F(\mathbf{x}(\xi, \sigma', \tau))|, \quad (\text{D.6})$$

$$B_2 := \max_{\sigma \in \mathcal{Z}^m} \max_{\substack{\tau, \tau' \in \mathcal{Z}^n \\ d(\tau, \tau') \leq 1}} |\mathbb{E}_\xi F(\mathbf{x}(\xi, \sigma, \tau)) - \mathbb{E}_\xi F(\mathbf{x}(\xi, \sigma, \tau'))|, \quad (\text{D.7})$$

$$V_* := \sup_{\xi, \tau, \sigma} \sum_{i \leq m, j \leq n} \text{Var}_{\xi_{ij}} \{F(\mathbf{x}(\xi, \sigma, \tau))\}, \quad (\text{D.8})$$

$$V_1 := \sup_{\tau, \sigma} \sum_{i \leq m} \text{Var}_{\sigma_i} \{\mathbb{E}_\xi F(\mathbf{x}(\xi, \sigma, \tau))\}, \quad (\text{D.9})$$

$$V_2 := \sup_{\tau, \sigma} \sum_{j \leq n} \text{Var}_{\tau_j} \{\mathbb{E}_\xi F(\mathbf{x}(\xi, \sigma, \tau))\}. \quad (\text{D.10})$$

Then, for any $t \geq 0$, the following holds with probability at least $1 - 8e^{-t}$:

$$|F(\mathbf{x}(\xi, \sigma, \tau)) - \mathbb{E}F(\mathbf{x}(\xi, \sigma, \tau))| \leq 2 \max(\sqrt{2V_*t} + \sqrt{2V_1t} + \sqrt{2V_2t}; (B_* + B_1 + B_2)t). \quad (\text{D.11})$$

Proof. Let $\mathbf{z} \in \mathcal{Z}^N$ be a vector of independent random variables and $f : \mathcal{Z}^N \rightarrow \mathbb{R}$. Define the martingale $X_k := \mathbb{E}[f(\mathbf{z}) | \mathcal{F}_k]$ (where $\mathcal{F}_k := \sigma(z_1, \dots, z_k)$). Then we have

$$\text{ess sup} |X_k - X_{k-1}| \leq B_0 := \sup_{d(\mathbf{z}, \mathbf{z}') \leq 1} |f(\mathbf{z}) - f(\mathbf{z}')|, \quad (\text{D.12})$$

$$\sum_{k=1}^N \mathbb{E}[(X_k - X_{k-1})^2 | \mathcal{F}_{k-1}] = \sum_{k=1}^N \mathbb{E}[(\mathbb{E}[f | \mathbf{z}_{<k}, z_k] - \mathbb{E}_{z'_k} \mathbb{E}[f | \mathbf{z}_{<k}, z'_k])^2 | \mathbf{z}_{<k}] \quad (\text{D.13})$$

$$\leq V_0 := \sup_{\mathbf{z} \in \mathcal{Z}^N} \sum_{k=1}^N \text{Var}_{z_k}(f(\mathbf{z})). \quad (\text{D.14})$$

By Freedman's inequality, with probability at least $1 - 2e^{-t}$, we have

$$|f(\mathbf{z}) - \mathbb{E}f(\mathbf{z})| \leq \max\left(\sqrt{2V_0t}; \frac{2B_0t}{3}\right). \quad (\text{D.15})$$

Define $E(\sigma, \tau) := \mathbb{E}_\xi F(\mathbf{x}(\xi, \sigma, \tau))$, $L(\tau) := \mathbb{E}_{\sigma, bxi} F(\mathbf{x}(\xi, \sigma, \tau))$. Applying the above inequality, each of the following holds with probability at least $1 - 2e^{-t}$

$$|F(\mathbf{x}(\xi, \sigma, \tau)) - E(\sigma, \tau)| \leq \max\left(\sqrt{2V_*t}; \frac{2B_*t}{3}\right), \quad (\text{D.16})$$

$$|E(\sigma, \tau) - L(\tau)| \leq \max\left(\sqrt{2V_1t}; \frac{2B_1t}{3}\right), \quad (\text{D.17})$$

$$|L(\tau) - \mathbb{E}F(\mathbf{x})| \leq \max\left(\sqrt{2V_2t}; \frac{2B_2t}{3}\right), \quad (\text{D.18})$$

and the claim follows by union bound. \square

We next state and prove a more stronger version of Lemma 5.3.

Lemma D.2. For $\mathbf{X} \in \mathcal{X}^{m \times n}$, let $P(\mathbf{X}) = P_{Q, q_r, q_c; m, n}(\mathbf{X})$ the probability of table \mathbf{X} under the model $\mathcal{T}(Q, q_r, q_c; m, n)$, i.e.

$$P(\mathbf{X}) = \sum_{\mathbf{u} \in \mathcal{L}^m} \sum_{\mathbf{v} \in \mathcal{L}^n} \prod_{(i,j) \in [m] \times [n]} Q(X_{ij}|u_i, v_j) \prod_{i \in [m]} q_r(u_i) \prod_{j \in [n]} q_c(v_j). \quad (\text{D.19})$$

Define the following quantities:

$$M_* := \max_{x, x' \in \mathcal{X}} \max_{u, v \in \mathcal{L}} \left| \log \frac{Q(x|u, v)}{Q(x'|u, v)} \right|, \quad (\text{D.20})$$

$$M_1 := \max_{\tau, \sigma, \sigma'} \|Q(\cdot | \sigma, \tau) - Q(\cdot | \sigma', \tau)\|_{\text{TV}} \max_{u, v, x, x'} \left| \log \frac{Q(x|u, v)}{Q(x'|u, v)} \right|, \quad (\text{D.21})$$

$$M_2 := \max_{\tau, \tau', \sigma} \|Q(\cdot | \sigma, \tau) - Q(\cdot | \sigma, \tau')\|_{\text{TV}} \max_{u, v, x, x'} \left| \log \frac{Q(x|u, v)}{Q(x'|u, v)} \right|, \quad (\text{D.22})$$

$$s_* := \frac{1}{2} \max_{u_0, v_0 \in \mathcal{L}} \sum_{x, x' \in \mathcal{X}} Q(x|u_0, v_0) Q(x'|u_0, v_0) \max_{u, v \in \mathcal{L}} \left(\log \frac{Q(x|u, v)}{Q(x'|u, v)} \right)^2, \quad (\text{D.23})$$

$$s_1 := \frac{1}{2} \max_{u_0, u'_0, v_0 \in \mathcal{L}} \|Q(\cdot | u_0, v_0) - Q(\cdot | u'_0, v_0)\|_{\text{TV}} \max_{x, x' \in \mathcal{L}} \max_{u, v \in \mathcal{L}} \left(\log \frac{Q(x|u, v)}{Q(x'|u, v)} \right)^2, \quad (\text{D.24})$$

$$s_2 := \frac{1}{2} \max_{u_0, v_0, v'_0 \in \mathcal{L}} \|Q(\cdot | u_0, v_0) - Q(\cdot | u_0, v'_0)\|_{\text{TV}} \max_{x, x' \in \mathcal{L}} \max_{u, v \in \mathcal{L}} \left(\log \frac{Q(x|u, v)}{Q(x'|u, v)} \right)^2. \quad (\text{D.25})$$

Then, for $\mathbf{X} \sim \mathcal{T}(Q, q_r, q_c; m, n)$ and any $t \geq 0$ the following bound holds with probability at least $1 - 2e^{-t}$:

$$\left| -\log P(\mathbf{X}) - H(\mathbf{X}) \right| \leq 3 \max \left(\sqrt{s_* m n t} + \sqrt{s_1 m n^2 t} + \sqrt{s_2 m^2 n t}, M_* + M_1 n + M_2 m \right). \quad (\text{D.26})$$

Proof. Let $\sigma = (\sigma_i)_{i \leq m} \sim_{iid} r$, $\tau = (\tau_i)_{i \leq n} \sim_{iid} c$, $\xi = (\xi_{ij})_{i \leq m, j \leq n} \sim_{iid} \text{Unif}([0, 1])$, and $x : [0, 1] \times \mathcal{L} \times \mathcal{L} \rightarrow \mathcal{X}$ be such that $x(\xi_{ij}, \sigma_i, \tau_j) |_{\sigma_i, \tau_j} \sim Q(\cdot | \sigma_i, \tau_j)$. We define $F(\mathbf{x}) = -\log P(\mathbf{x})$, and will apply Lemma D.1 to this function. Using the notation from that lemma, we claim that $B_* \leq M_*$, $B_1 \leq M_1 n$, $B_2 \leq M_2 m$, and $V_* \leq m n s_*$, $V_1 \leq m n^2 s_1$, $V_2 \leq m^2 n s_2$.

Note that, if $(x_{ij}), (x'_{ij})$ differ only for entry i, j , then

$$F(\mathbf{x}) - F(\mathbf{x}') = -\log \mathbb{E}_{\mathbf{u}, \mathbf{v} | \mathbf{x}} \left\{ \frac{Q(x'_{ij} | u_i, v_j)}{Q(x_{ij} | u_i, v_j)} \right\}, \quad (\text{D.27})$$

where $\mathbb{E}_{\mathbf{u}, \mathbf{v} | \mathbf{x}}$ denotes expectation with respect to the posterior measure $P(\mathbf{u}, \mathbf{v} | \mathbf{X} = \mathbf{x})$. This immediately implies $B_* \leq M_*$.

Next consider the constant B_1 defined in Eq. (D.6). Using the exchangeability of the $(\xi_{i, \cdot}, \sigma_i)$, we get

$$\begin{aligned} B_1 &= \max_{\tau} \left| \mathbb{E}_{\xi} \log \mathbb{E}_{\mathbf{u}, \mathbf{v} | \mathbf{x}} \left\{ \prod_{j=1}^n \frac{Q(x(\xi_{1,j}, \sigma'_1, \tau_j) | u_1, v_j)}{Q(x(\xi_{1,j}, \sigma_1, \tau_j) | u_1, v_j)} \right\} \right| \\ &\leq \max_{\tau} \mathbb{E}_{\xi} \max_{\mathbf{u}, \mathbf{v}} \left| \log \left\{ \prod_{j=1}^n \frac{Q(x(\xi_{1,j}, \sigma'_1, \tau_j) | u_1, v_j)}{Q(x(\xi_{1,j}, \sigma_1, \tau_j) | u_1, v_j)} \right\} \right| \\ &\leq \max_{\tau} \sum_{j=1}^n \mathbb{E}_{\xi} \max_{\mathbf{u}, \mathbf{v}} \left| \log \left\{ \frac{Q(x(\xi_{1,j}, \sigma'_1, \tau_j) | u_1, v_j)}{Q(x(\xi_{1,j}, \sigma_1, \tau_j) | u_1, v_j)} \right\} \right| \\ &\leq n \max_{\tau, \sigma, \sigma'} \mathbb{E}_{\xi} \max_{\mathbf{u}, \mathbf{v}} \left| \log \frac{Q(x(\xi, \sigma'_1, \tau_j) | u, v)}{Q(x(\xi, \sigma_1, \tau_j) | u, v)} \right| \\ &\leq n \max_{\tau, \sigma, \sigma'} \|Q(\cdot | \sigma, \tau) - Q(\cdot | \sigma', \tau)\|_{\text{TV}} \max_{u, v, x, x'} \left| \log \frac{Q(x|u, v)}{Q(x'|u, v)} \right| = M_1. \end{aligned}$$

The bound $B_2 \leq M_2 m$ is proved analogously.

Consider now the quantity V_* of Eq. (D.8). Denote by $\xi_{(ij)}(t)$ the array obtained by replacing entry (i, j) in ξ by t , and by $\mathbf{x}(t) = \mathbf{x}(\xi_{(ij)}(t), \sigma, \tau)$. Then we have

$$\begin{aligned} \text{Var}_{\xi_{ij}}(F(\mathbf{x})) &= \frac{1}{2} \mathbb{E}_{\xi', \xi''} \left\{ \left(F(\mathbf{x}(\xi_{(ij)}(\xi')), \sigma, \tau) - F(\mathbf{x}(\xi_{(ij)}(\xi''), \sigma, \tau)) \right)^2 \right\} \\ &= \frac{1}{2} \mathbb{E}_{\xi', \xi''} \left\{ \left(\log \mathbb{E}_{\mathbf{u}, \mathbf{v} | \mathbf{x}(\xi')} \left\{ \frac{Q(\mathbf{x}(\xi'', \sigma_i, \tau_j) | u_i, v_j)}{Q(\mathbf{x}(\xi', \sigma_i, \tau_j) | u_i, v_j)} \right\} \right)^2 \right\} \\ &\leq \frac{1}{2} \mathbb{E}_{\xi', \xi''} \max_{u, v} \left(\log \left\{ \frac{Q(\mathbf{x}(\xi'', \sigma_i, \tau_j) | u, v)}{Q(\mathbf{x}(\xi', \sigma_i, \tau_j) | u, v)} \right\} \right)^2 \\ &= \frac{1}{2} \sum_{x, x'} Q(x | \sigma, \tau) Q(x' | \sigma, \tau) \max_{u, v} \left(\log \left\{ \frac{Q(x | u, v)}{Q(x' | u, v)} \right\} \right)^2. \end{aligned}$$

We then have, as claimed

$$\begin{aligned} V_* &\leq \max_{\xi, \sigma, \tau} \sum_{i \leq m, j \leq n} \text{Var}_{\xi_{ij}} \{ F(\mathbf{x}) \} \\ &\leq mn \max_{\xi, \sigma, \tau} \text{Var}_{\xi_{ij}} \{ F(\mathbf{x}) \} \leq mns_*. \end{aligned}$$

Finally consider the quantity V_1 of Eq. (D.9) (the argument is similar for V_2). Denote by $\sigma_{(i)}(t)$ the vector obtained by replacing entry i in σ by t . Proceeding as above, we have

$$\begin{aligned} \text{Var}_{\sigma_i}(\mathbb{E}_{\xi} F(\mathbf{x})) &= \frac{1}{2} \mathbb{E}_{\sigma', \sigma''} \left\{ \left(\mathbb{E}_{\xi} F(\mathbf{x}(\xi, \sigma_{(i)}(\sigma')), \tau) - \mathbb{E}_{\xi} F(\mathbf{x}(\xi, \sigma_{(i)}(\sigma''), \tau)) \right)^2 \right\} \\ &= \frac{1}{2} \mathbb{E}_{\sigma', \sigma''} \left\{ \left(\mathbb{E}_{\xi} \log \mathbb{E}_{\mathbf{u}, \mathbf{v} | \mathbf{x}(\sigma')} \left\{ \prod_{j=1}^n \frac{Q(\mathbf{x}(\xi_{ij}, \sigma'', \tau_j) | u_i, v_j)}{Q(\mathbf{x}(\xi_{ij}, \sigma', \tau_j) | u_i, v_j)} \right\} \right)^2 \right\} \\ &\leq \frac{1}{2} \mathbb{E}_{\sigma', \sigma''} \left\{ \left(\mathbb{E}_{\xi} \log \left\{ \prod_{j=1}^n \max_{u, v} \frac{Q(\mathbf{x}(\xi_{ij}, \sigma'', \tau_j) | u, v)}{Q(\mathbf{x}(\xi_{ij}, \sigma', \tau_j) | u, v)} \right\} \right)^2 \right\} \\ &\leq \frac{1}{2} \mathbb{E}_{\sigma', \sigma''} \left\{ \left(\sum_{j=1}^n E_{\xi} \log \left\{ \max_{u, v} \frac{Q(\mathbf{x}(\xi, \sigma'', \tau_j) | u, v)}{Q(\mathbf{x}(\xi, \sigma', \tau_j) | u, v)} \right\} \right)^2 \right\} \\ &\leq \frac{n^2}{2} \max_{\tau} \mathbb{E}_{\sigma', \sigma''} \left\{ \left(\mathbb{E}_{\xi} \log \left\{ \max_{u, v} \frac{Q(\mathbf{x}(\xi, \sigma'', \tau) | u, v)}{Q(\mathbf{x}(\xi, \sigma', \tau) | u, v)} \right\} \right)^2 \right\} \\ &\leq \frac{n^2}{2} \max_{\tau, \sigma, \sigma'} \|Q(\cdot | \sigma, \tau) - Q(\cdot | \sigma', \tau)\|_{\text{TV}} \max_{x, x'} \left(\mathbb{E}_{\xi} \log \left\{ \max_{u, v} \frac{Q(x' | u, v)}{Q(x | u, v)} \right\} \right)^2 = n^2 s_1. \end{aligned}$$

Therefore

$$V_1 = \max_{\sigma, \tau} \sum_{i=1}^m \text{Var}_{\sigma_i}(\mathbb{E}_{\xi} F(\mathbf{x})) \leq mn^2 s_1.$$

This finishes the proof. \square

E. Proofs for finite state encoders

Recall from Section 5.2 that a finite-state encoder is defined by a triple (Σ, f, g) . Formally we can define the action of f, g on $\mathbf{X}^n \in \mathcal{X}^n$ recursively via (recall that \oplus denotes concatenation)

$$f_{m+1}(\mathbf{X}^{m+1}, s_0) = f_m(\mathbf{X}^m, s_0) \oplus f(X_{m+1}, g(\mathbf{X}^m, s_0)), \quad (\text{E.1})$$

$$g_{m+1}(\mathbf{X}^{m+1}, s_0) = g(X_{m+1}, g(\mathbf{X}^m, s_0)), \quad (\text{E.2})$$

and the encoder is thus given by $E(\mathbf{X}^n) = f_n(\mathbf{X}^n, s_{\text{init}})$.

We say that the state space Σ is non-degenerate if, for each $s_1 \in \Sigma$ there exists $m, \mathbf{X}^m \in \mathcal{X}^m$ such that $g_m(\mathbf{X}^m, s_{\text{init}}) = s_1$. Notice that if state space is degenerate, we could always remove one or more symbols from Σ without changing the encoder, and making the state-space non-degenerate. For this reason, we will hereafter assume non-degeneracy without mentioning it.

We say that the FS encoder is information lossless (IL) if for any $n \in \mathbb{N}$, $\mathbf{X}^n \mapsto f_n(\mathbf{X}^n, s_{\text{init}})$ is injective.

Remark E.1. An information-lossless encoder satisfies a stronger condition: for any $m \in \mathbb{N}$ and any $s_* \in \Sigma$, the map $\mathbf{X}^m \mapsto f_m(\mathbf{X}^m, s_*)$ is injective.

Indeed, assume this were not the case. Then there would exist two distinct inputs $\mathbf{X}^m, \tilde{\mathbf{X}}_1^m \in \mathcal{X}^m$ and a state $s_* \in \Sigma$ such that $f_m(\mathbf{X}^m, s_*) = f_m(\tilde{\mathbf{X}}_1^m, s_*)$. By non-degeneracy, there exists $a_1^\ell \in \mathcal{X}^\ell$ such that $s_* = g_\ell(a_1^\ell, s_{\text{init}})$. Defining $n = \ell + m$, $\mathbf{Y}^n = a_1^\ell \oplus \mathbf{X}^m$, $\tilde{\mathbf{Y}}^n = a_1^\ell \oplus \tilde{\mathbf{X}}_1^m$, it is not hard to check that these inputs are distinct but $f_n(\mathbf{Y}^n, s_{\text{init}}) = f_n(\tilde{\mathbf{Y}}^n, s_{\text{init}})$.

Proposition E.1. Define the compression rate on input x_1^n as $R(\mathbf{X}^n) = \text{len}(f_n(\mathbf{X}^n, s_{\text{init}})) / (n \log_2 |\mathcal{X}|)$. Then for any $\ell \geq 1$, the following holds (where $n' := n - 2\ell$ and we recall that $M := |\Sigma|$):

$$R(\mathbf{X}^n) \geq \frac{n - 2\ell}{n\ell \log_2 |\mathcal{X}|} H(\hat{p}_{\mathbf{X}_1^{n'}}^\ell) - \frac{1}{\ell \log_2 |\mathcal{X}|} (\log_2(|\Sigma|\ell) + \log_2 \log_2 |\mathcal{X}|). \quad (\text{E.3})$$

Proof. We will denote by $L(\mathbf{X}^m; s_*)$ the length of the encoding of \mathbf{X}^m when starting in state s_* :

$$L(\mathbf{X}^m; s_*) := \text{len}(f_m(\mathbf{X}^m, s_*)). \quad (\text{E.4})$$

We then have, for any $b \in \{0, \dots, \ell - 1\}$, and setting by convention $s_0 = s_{\text{init}}$, we get

$$R(\mathbf{X}^n) \geq \frac{1}{n \log_2 |\mathcal{X}|} \sum_{k=0}^{\lfloor n/\ell \rfloor - 2} L(\mathbf{X}_{k\ell+b+1}^{(k+1)\ell+b}; s_{k\ell+b}). \quad (\text{E.5})$$

By averaging over b , and introducing the shorthand $n' := n - 2\ell$, we get

$$R(\mathbf{X}^n) \geq \frac{1}{n\ell \log_2 |\mathcal{X}|} \sum_{m=1}^{(\lfloor n/\ell \rfloor - 1)\ell} L(\mathbf{X}_m^{m+\ell-1}; s_{m-1}) \quad (\text{E.6})$$

$$\geq \frac{n - 2\ell}{n\ell \log_2 |\mathcal{X}|} \sum_{s \in \Sigma} \sum_{u_1^\ell \in \mathcal{X}^\ell} \hat{p}_{\mathbf{X}_1^{n'}}^\ell(u_1^\ell, s) L(u_1^\ell; s) \quad (\text{E.7})$$

$$\stackrel{(a)}{\geq} \frac{n - 2\ell}{n\ell \log_2 |\mathcal{X}|} \sum_{s \in \Sigma} \left\{ \hat{p}_{\mathbf{X}_1^{n'}}^\ell(s) H(\hat{p}_{\mathbf{X}_1^{n'}}^\ell(\cdot | s)) - \log_2 \log_2(|\mathcal{X}|^\ell) \right\}, \quad (\text{E.8})$$

where (a) holds by Lemma C.1. By the chain rule of entropy (recalling that $M := |\Sigma|$), we have:

$$\begin{aligned} \sum_{s \in \Sigma} \hat{p}_{\mathbf{X}_1^{n'}}^\ell(s) H(\hat{p}_{\mathbf{X}_1^{n'}}^\ell(\cdot | s)) &= H(\mathbf{X}_1^\ell | S) = H(\mathbf{X}_1^\ell) + H(S | \mathbf{X}_1^\ell) - H(S) \\ &\geq H(\mathbf{X}_1^\ell) - \log_2 M = H(\hat{p}_{\mathbf{X}_1^{n'}}^\ell) - \log_2 M. \end{aligned}$$

The claim (E.3) follows by using the last inequality in Eq. (E.8). \square

Theorem E.2. Let $\mathbf{X}^{m,n} \sim \mathcal{T}(Q, q_r, q_c; m, n)$ and (Σ, f, g) be an information lossless finite state encoder. With an abuse of notation, denote $f_{mn}(\mathbf{X}^{m \times n}, s_{\text{init}}) \in \{0, 1\}^*$ the binary sequence obtained by applying the finite state encoder to the vector $\text{vec}(\mathbf{X}^{m \times n}) \in \mathcal{X}^{mn}$ obtained by scanning $\mathbf{X}^{m \times n}$ in row-first order. Define the compression rate by

$$R(\mathbf{X}^{m,n}) := \frac{\text{len}(f_{mn}(\mathbf{X}^{m \times n}, s_{\text{init}}))}{mn \log_2 |\mathcal{X}|}. \quad (\text{E.9})$$

Assuming $m > 10$, $|\Sigma| \geq |\mathcal{X}|$, and $\log_2 |\Sigma| \leq n \log_2 |\mathcal{X}| / 9$, the expected compression rate is lower bounded as follows

$$\mathbb{E} R(\mathbf{X}^{m,n}) \geq \frac{H(X|U)}{\log_2 |\mathcal{X}|} - 10 \sqrt{\frac{\log |\Sigma|}{n \log |\mathcal{X}|}} \cdot \log(n \log |\Sigma|). \quad (\text{E.10})$$

Proof. We let $N := mn$, $N' = mn - 2\ell$ where we $\ell \leq n/3$ will be selected later. We write $\mathbf{X}^N := \text{vec}(\mathbf{X}^{m,n})$ for the vectorization $\mathbf{X}^{m,n}$, $\mathbf{X}^{N'}$ for the vector comprising its first N' entries. Recall the definition of empirical distribution. For any fixed $\mathbf{w} \in \mathcal{X}^\ell$

$$\hat{p}_{\mathbf{X}^{N'}}^\ell(\mathbf{w}) := \frac{1}{N' - \ell + 1} \sum_{i=1}^{N' - \ell + 1} \mathbf{1}_{\mathbf{X}_i^{i+\ell-1} = \mathbf{w}}.$$

Let $S := \{i \in [N' - \ell + 1] : [i, i + \ell - 2] \cap n\mathbb{N} = \emptyset\}$. In words, these are the subset of blocks of length ℓ that do not cross the end of a line in the table. Since for each line break there are at most $\ell - 1$ such blocks, we have $|S| \geq N' - \ell + 1 - (m - 1)(\ell - 1)$. We will consider the following modified empirical distribution

$$\bar{p}_{\mathbf{X}^{N'}}^\ell(\mathbf{w}) := \frac{1}{|S|} \sum_{i \in S} \mathbf{1}_{\mathbf{X}_i^{i+\ell-1} = \mathbf{w}}.$$

Then by construction

$$\begin{aligned} \hat{p}_{\mathbf{X}^{N'}}^\ell(\mathbf{w}) &= (1 - \eta_\ell) \bar{p}_{\mathbf{X}^{N'}}^\ell(\mathbf{w}) + \eta_\ell q_{\mathbf{X}^{N'}}^\ell(\mathbf{w}), \\ \eta_\ell &:= 1 - \frac{|S|}{N' - \ell + 1} = \frac{(m - 1)(\ell - 1)}{N' - \ell + 1}, \end{aligned}$$

where $q_{\mathbf{X}^{N'}}^\ell$ is the empirical distribution of blocks that do cross the line. By concavity of the entropy, we have

$$H(\hat{p}_{\mathbf{X}^{N'}}^\ell) \geq (1 - \eta_\ell) H(\bar{p}_{\mathbf{X}^{N'}}^\ell) + \eta_\ell H(q_{\mathbf{X}^{N'}}^\ell) \geq (1 - \eta_\ell) H(\bar{p}_{\mathbf{X}^{N'}}^\ell). \quad (\text{E.11})$$

Further, since $\ell \leq n/3$,

$$\begin{aligned} \eta_\ell &= \frac{(m - 1)(\ell - 1)}{mn - 3\ell + 1} \\ &\leq \frac{(m - 1)\ell}{mn - 3\ell} \leq \frac{(m - 1)\ell}{(m - 1)n} \leq \frac{\ell}{n}. \end{aligned} \quad (\text{E.12})$$

Now let the row latents $\mathbf{u} := (u_i)_{i \leq m}$ be fixed, and denote by $\hat{r}_{\mathbf{u}}^S$ their weighted empirical distribution, defined as follows:

$$\hat{r}_{\mathbf{u}}^S(s) := \sum_{i=1}^m \frac{|S \cap [(i - 1)n + 1, in]|}{|S|} \mathbf{1}_{u_i = s}.$$

In words, $\hat{r}_{\mathbf{u}}^S$ is the empirical distribution of the latents $(u_i)_{i \leq m}$ where row i is weighted by its contribution to S . Note that all the weights are equal to $(n - 2(\ell - 1))/|S|$ except, potentially, for the last one.

We have

$$p_*^\ell(\mathbf{w}) := \mathbb{E}[\bar{p}_{\mathbf{X}^{N'}}^\ell(\mathbf{w})] = \sum_{u \in \mathcal{L}} \hat{r}_{\mathbf{u}}^S(u) \prod_{i=1}^{\ell} Q_{x|u}(w_i|u), \quad Q_{x|u}(w|u) := \sum_{v \in \mathcal{L}} Q(w|u, v) q_{\mathbf{c}}(v).$$

Using Eq. (E.11), (E.12) and concavity of the entropy, we get

$$\mathbb{E}[H(\hat{p}_{\mathbf{X}^{N'}}^\ell) | \mathbf{u}] \geq \left(1 - \frac{\ell}{n}\right) H(p_*^\ell). \quad (\text{E.13})$$

By Proposition E.1, we get

$$\begin{aligned} \mathbb{E}[R(\mathbf{X}^{m,n}) | \mathbf{u}] &\geq \frac{mn - 2\ell}{mn\ell \log_2 |\mathcal{X}|} \left(1 - \frac{\ell}{n}\right) H(p_*^\ell) - \frac{1}{\ell \log_2 |\mathcal{X}|} (\log_2(|\Sigma|\ell) + \log_2 \log_2 |\mathcal{X}|) \\ &\geq \frac{1}{\ell \log_2 |\mathcal{X}|} H(p_*^\ell) - \frac{2\ell}{n} - \frac{1}{\ell \log_2 |\mathcal{X}|} (\log_2(|\Sigma|\ell) + \log_2 \log_2 |\mathcal{X}|), \end{aligned}$$

Algorithm 3 Lempel-Ziv

Input: Data $\mathbf{X}^N \in \mathcal{X}^N = \{0, \dots, |\mathcal{X}| - 1\}^N$
Output: Binary string $\mathbf{Z} \in \{0, 1\}^*$
for $k = 1$ **to** N **do**
 if $\exists j < k : X_j = X_k$ **then**
 $L_k \leftarrow \max\{\ell \geq 1 : \exists j \in \{1, \dots, k-1\} \text{ s.t. } \mathbf{X}_j^{j+\ell-1} = \mathbf{X}_k^{k+\ell-1}\}$
 $T_k \leftarrow \max\{j \in \{1, \dots, k-1\} \text{ s.t. } \mathbf{X}_j^{j+L_k-1} = \mathbf{X}_k^{k+L_k-1}\}$
 else
 $L_k \leftarrow 1$
 $T_k \leftarrow (-X_k)$
 end if
 $\mathbf{Z} \leftarrow \mathbf{Z} \oplus \text{plain}(T_k) \oplus \text{elias}(L_k)$
 $k \leftarrow k + L_k$
 if $\text{len}(\mathbf{Z}) \leq \text{len}(\text{plain}(\mathbf{X}^N))$ **then**
 return \mathbf{Z}
 else
 return $\text{plain}(\mathbf{X}^N)$
 end if
end for

where in the last inequality we used the fact that $H(p_*^\ell) \leq \ell \log_2 |\mathcal{X}|$. We choose

$$\ell = \sqrt{\frac{n \log_2 |\Sigma|}{\log_2 |\mathcal{X}|}} \leq \frac{n}{3}, \quad (\text{E.14})$$

Substituting and simplifying, we get

$$\mathbb{E}[\mathbf{R}(\mathbf{X}^{m,n})|\mathbf{u}] \geq \frac{H(p_*^\ell)}{\ell \log_2 |\mathcal{X}|} - \frac{10}{\sqrt{n}} \cdot \sqrt{\frac{\log |\Sigma|}{\log |\mathcal{X}|}} \cdot \log(n \log |\Sigma|). \quad (\text{E.15})$$

Finally, letting $(W_1, \dots, W_\ell, U) \in \mathcal{X}^\ell \times \mathcal{L}$ be random variables with joint distribution $\hat{r}_\mathbf{u}^S(u) \prod_{i=1}^\ell Q_{x|u}(w_i|u)$. Then

$$H(p_*^\ell) \geq \sum_{u \in \mathcal{L}} \hat{r}_\mathbf{u}^S(u) H(Q_{x|u}^{\otimes \ell}(\cdot|u)) \quad (\text{E.16})$$

$$\geq \ell \sum_{u \in \mathcal{L}} \hat{r}_\mathbf{u}^S(u) H(X|U = u), \quad (\text{E.17})$$

and therefore $\mathbb{E}H(p_*^\ell) \geq H(X|U)$, finishing the proof. \square

F. Proofs for Lempel-Ziv coding

The pseudocode of the Lempel-Ziv algorithm that we will analyze is given here. For ease of presentation, we identify \mathcal{X} with a set of integers.

Note that if a simbol X_k never appeared in the past, we point to $T_k = -X_k$ and set $L_k = 1$. This is essentially equivalent to prepending a sequence of distinct $|\mathcal{X}|$ symbols to \mathbf{X}^N .

It is useful to define for each $k \leq N$,

$$L_k(\mathbf{X}^N) := \max\{\ell \geq 1 : \exists j \in \{1, \dots, k-1\} \text{ s.t. } \mathbf{X}_j^{j+\ell-1} = \mathbf{X}_k^{k+\ell-1}\}, \quad (\text{F.1})$$

$$T_k(\mathbf{X}^N) := \min\{j \in \{1, \dots, k-1\} \text{ s.t. } \mathbf{X}_j^{j+L_k-1} = \mathbf{X}_k^{k+L_k-1}\}. \quad (\text{F.2})$$

F.1. Proof of Theorem 5.6

Lemma F.1. *Under Assumption 5.5, there exists a constant C such that the following holds with probability at least $1 - N^{-10}$:*

$$\max_{k \leq N} L_k(\mathbf{X}^N) \leq C \log N. \quad (\text{F.3})$$

Proof. We begin by considering a slightly different setting, and will then show that our question reduces to this setting. Let $(Z_i)_{i \geq 1}$ be independent random variables with $Z_i \sim q_i$ a probability distribution over \mathcal{X} . Further assume $\max_{x \in \mathcal{X}} q_i(x) \leq 1 - c$ for all $i \geq 1$. Then we claim that, for any $t, \ell \geq 1$, we have

$$\mathbb{P}(Z_1^\ell = Z_{t+1}^{t+\ell}) \leq (1 - c)^\ell. \quad (\text{F.4})$$

Indeed, condition on the event $Z_1^t = x_1^t$ for some $x_1, \dots, x_t \in \mathcal{X}$. Then the event $Z_1^\ell = Z_{t+1}^{t+\ell}$ implies that, for $i \in \{t+1, \dots, t+\ell\}$, $Z_i = x_{\pi(i)}$ where $\pi(i) = i \bmod t$, $\pi(i) \in [1, t]$. Then

$$\begin{aligned} \mathbb{P}(Z_1^\ell = Z_{t+1}^{t+\ell}) &\leq \max_{x_1^t \in \mathcal{X}^t} \mathbb{P}(Z_1^\ell = Z_{t+1}^{t+\ell} | Z_1^t = x_1^t) \\ &\leq \max_{x_1^t \in \mathcal{X}^t} \mathbb{P}(Z_i = x_{\pi(i)} \forall i \in \{t+1, \dots, t+\ell\} | Z_1^t = x_1^t) \\ &\leq \max_{x_1^t \in \mathcal{X}^t} \prod_{i=t+1}^{t+\ell} \mathbb{P}(Z_i = x_{\pi(i)}) \leq (1 - c)^\ell. \end{aligned}$$

This proves claim (F.4).

Let us now reconsider our original setting:

$$\begin{aligned} \mathbb{P}\left(\max_{k \leq N} L_k(\mathbf{X}^N) \geq \ell\right) &= \mathbb{P}(\exists i < j \leq N : X_i^{i+\ell-1} = X_j^{j+\ell-1}) \\ &\leq N^2 \max_{i < j \leq N} \mathbb{P}(X_i^{i+\ell-1} = X_j^{j+\ell-1}) \\ &\leq N^2 \max_{\mathbf{u}^m \in \mathcal{L}^m, \mathbf{v}^n \in \mathcal{L}^n} \max_{i < j \leq N} \mathbb{P}(X_i^{i+\ell-1} = X_j^{j+\ell-1} | \mathbf{u}^m, \mathbf{v}^n) \\ &\leq N^2 (1 - c)^\ell, \end{aligned}$$

where the last inequality follows from claim (F.4), since the $(X_i)_{i \leq N}$ are conditionally independent given the latents $\mathbf{u}^m, \mathbf{v}^n$, with probability mass function upper bounded by $1 - c$. The thesis follows by taking $\ell = 12 \log N / \log(1/(1 - c))$. \square

For $i \in [m], j \in [n]$, we define $\langle ij \rangle := (i - 1)n + j$. In words, $k = \langle ij \rangle$ is the of entry at row i column j when the table $\mathbf{X}^{m,n}$ is scanned in row first order. For $\ell \geq 1$, define the events

$$\mathcal{E}_{i,j}(\ell) := \left\{ \exists i' \in [m], j' \in [n] : \langle i'j' \rangle < \langle ij \rangle, |j' - j| \geq \ell, \mathbf{X}_{\langle i'j' \rangle}^{\langle i'j' \rangle + \ell - 1} = \mathbf{X}_{\langle ij \rangle}^{\langle ij \rangle + \ell - 1} \right\}, \quad (\text{F.5})$$

$$\mathcal{F}_{i,j}(\ell) := \left\{ \exists i' \in [m], j' \in [n] : \langle i'j' \rangle < \langle ij \rangle, |j' - j| < \ell, \mathbf{X}_{\langle i'j' \rangle}^{\langle i'j' \rangle + \ell - 1} = \mathbf{X}_{\langle ij \rangle}^{\langle ij \rangle + \ell - 1} \right\}. \quad (\text{F.6})$$

Then we have

$$\mathbb{P}(L_{\langle ij \rangle}(\mathbf{X}^N) \geq \ell) \leq \mathbb{P}(\mathcal{E}_{i,j}(\ell)) + \mathbb{P}(\mathcal{F}_{i,j}(\ell)). \quad (\text{F.7})$$

The next two lemmas control the probabilities of these events.

Lemma F.2. *Let $\ell(\delta, u) := \lceil (1 + \delta)(\log N) / H(X|U = u) \rceil$, $n' = n - \max_{u \in \mathcal{L}} \ell(\delta, u)$, and $m_0 = m^{1 - o_n(1)}$. Under Assumption 5.5, for any $\delta > 0$, there exist constants $C, \varepsilon > 0$ independent of $\mathbf{u} \in \mathcal{L}^m$, such that the following hold*

$$\max_{i \leq m, j \leq n'} \mathbb{P}(\mathcal{E}_{i,j}(\ell(\delta, u_i))) \leq C N^{-\varepsilon}, \quad (\text{F.8})$$

$$\min_{m_0 \leq i \leq m, j \leq n'} \mathbb{P}(\mathcal{E}_{i,j}(\ell(-\delta, u_i))) \geq 1 - C N^{-\varepsilon}. \quad (\text{F.9})$$

Lemma F.3. Let $\ell_c(\delta, u) := \lceil (1 + \delta)(\log m)/H(X|U = u, V) \rceil$, $n'_c = n - \max_{u \in \mathcal{L}} \ell_c(\delta, u)$, and $m_0 = m^{1-o_n(1)}$. Under Assumption 5.5, for any $\delta > 0$, there exist constants $C, \varepsilon > 0$ independent of $\mathbf{u} \in \mathcal{L}^m$, such that the following hold

$$\max_{i \leq m, j \leq n'_c} \mathbb{P}(\mathcal{F}_{i,j}(\ell_c(\delta, u_i))) \leq C m^{-\varepsilon}, \quad (\text{F.10})$$

$$\min_{m_0 \leq i \leq m, j \leq n'_c} \mathbb{P}(\mathcal{F}_{i,j}(\ell_c(-\delta, u_i))) \geq 1 - C m^{-\varepsilon}. \quad (\text{F.11})$$

We are now in position to prove Theorem 5.6.

Proof of Theorem 5.6. We denote by $(k(1), \dots, k(M))$ the values taken by k in the while loop of the Lempel-Ziv pseudocode. In particular

$$k(1) = 1, \quad (\text{F.12})$$

$$k(\ell + 1) = k(\ell) + L_{k(\ell)}(\mathbf{X}^N), \quad (\text{F.13})$$

$$k(M) = N. \quad (\text{F.14})$$

Therefore the total length of the code is

$$\text{len}(\text{LZ}(\mathbf{X}^{m,n})) = M \lceil \log_2(N + |\mathcal{X}|) \rceil + \sum_{\ell=1}^M \text{len}(\text{elias}(L_{k(\ell)})) \quad (\text{F.15})$$

By Lemma F.1 (and recalling that $\text{len}(\text{elias}(L)) \leq 2 \log_2 L + 1$) we have, with high probability, $\max_{\ell \leq M} \text{len}(\text{elias}(L_{k(\ell)})) \leq 2 \log_2(C \log N)$. Letting \mathcal{G} denote the ‘good’ event that this bound holds, we have, on \mathcal{G}

$$M \log_2 N \leq \text{len}(\text{LZ}(\mathbf{X}^{m,n})) \leq M \lceil \log_2(N + |\mathcal{X}|) \rceil + 2M \log_2(C \log N) \quad (\text{F.16})$$

Since $|\mathcal{X}|$ is a constant, this means that for any $\eta > 0$, there exists $N_0(\eta)$ such that, for all $N \geq N_0(\eta)$, with probability at least $1 - \eta$:

$$M \cdot \mathbf{1}_{\mathcal{G}} \log_2 N \leq \text{len}(\text{LZ}(\mathbf{X}^{m,n})) \leq (1 + \eta)M \cdot \mathbf{1}_{\mathcal{G}} \log_2 N + N \cdot \mathbf{1}_{\mathcal{G}^c} \log_2 |\mathcal{X}|, \quad (\text{F.17})$$

where on the right $\text{len}(\text{LZ}\mathbf{X}^{m,n}) \leq N \log_2 |\mathcal{X}|$ by construction. We thus have

$$\mathbb{E}\{M \cdot \mathbf{1}_{\mathcal{G}}\} \frac{\log_2 N}{N \log_2 |\mathcal{X}|} \leq \mathbb{E} R_{\text{LZ}}(\mathbf{X}^{m,n}) \leq (1 + \eta) \mathbb{E}\{M \cdot \mathbf{1}_{\mathcal{G}}\} \frac{\log_2 N}{N \log_2 |\mathcal{X}|} + \eta, \quad (\text{F.18})$$

that is

$$\liminf_{m,n \rightarrow \infty} \mathbb{E} R_{\text{LZ}}(\mathbf{X}^{m,n}) \geq \liminf_{m,n \rightarrow \infty} \mathbb{E}\{M \cdot \mathbf{1}_{\mathcal{G}}\} \cdot \frac{\log_2 N}{N \log_2 |\mathcal{X}|}, \quad (\text{F.19})$$

$$\limsup_{m,n \rightarrow \infty} \mathbb{E} R_{\text{LZ}}(\mathbf{X}^{m,n}) \leq \limsup_{m,n \rightarrow \infty} \mathbb{E}\{M \cdot \mathbf{1}_{\mathcal{G}}\} \cdot \frac{\log_2 N}{N \log_2 |\mathcal{X}|}. \quad (\text{F.20})$$

We are therefore left with the task of bounding $\mathbb{E}\{M \cdot \mathbf{1}_{\mathcal{G}}\}$

We begin by the lower bound. Define the set of ‘bad indices’ $B(\mathbf{X}^{m,n}, \delta) \subseteq [m] \times [n]$,

$$B(\mathbf{X}^{m,n}, \delta) := \left\{ (i, j) \in [m] \times [n] : \mathcal{E}_{i,j}(\ell(\delta, u_i)) \text{ or } \mathcal{F}_{i,j}(\ell_c(\delta, u_i)) \right\} \quad (\text{F.21})$$

We will drop the arguments $\mathbf{X}^{m,n}, \delta$ for economy of notation, and write $B := B(\mathbf{X}^{m,n}, \delta)$. We further define

$$S(u) = S(u; \mathbf{X}^{m,n}) := \{(i, j) \in [m] \times [n] : u_i = u \text{ and } \exists \ell \leq M : \langle ij \rangle = k(\ell)\}. \quad (\text{F.22})$$

In words, $S(u)$ is the set of positions (i, j) of the table $\mathbf{X}^{m,n}$ where words in the LZ parsing begin.

We also write $N(u) = n \cdot |\{i \in [m] : u_i = u\}|$ for the total number of rows in $\mathbf{X}^{m,n}$ with row latent equal to u_i and L_i^- for the length of the first segment in row i initiated in row $i - 1$:

$$\begin{aligned}
 N(u) &\leq \sum_{(i,j) \in S(u)} L_{\langle ij \rangle} + \sum_{i \leq m: u_i = u} L_i^- \\
 &\leq \sum_{(i,j) \in S(u) \cap B^c} L_{\langle ij \rangle} + \sum_{(i,j) \in B} L_{\langle ij \rangle} + \sum_{i \leq m: u_i = u} L_i^- \\
 &\leq \sum_{(i,j) \in S(u)} \ell(u; \delta) \vee \ell_c(u; \delta) + (|B| + m) \cdot C \log N \\
 &\leq |S(u)| \ell(u; \delta) \vee \ell_c(u; \delta) + (|B| + m) \cdot C \log N,
 \end{aligned}$$

where the last inequality holds on event \mathcal{G} . By taking expectation on this event, we get

$$\mathbb{E}\{N(u) \cdot \mathbf{1}_{\mathcal{G}}\} \leq \mathbb{E}\{|S(u)| \cdot \mathbf{1}_{\mathcal{G}}\} \cdot \ell(u; \delta) \vee \ell_c(u; \delta) + (\mathbb{E}|B| + m) \cdot C \log N.$$

By Lemmas F.2 and F.2,

$$\begin{aligned}
 \mathbb{E}(|B|) &\leq m_0 n + \sum_{m_0 \leq i \leq m, j \leq n'} \mathbb{P}\left(\mathcal{E}_{i,j}(\ell(\delta, u_i)) \cup \mathcal{F}_{i,j}(\ell_c(\delta, u_i))\right) + Cm \log N \\
 &\leq m_0 n + Cm^{1-\varepsilon} n + Cm \log N \\
 &\leq \frac{CN}{(\log N)^2}.
 \end{aligned}$$

$\mathbb{E}(|B|) \leq Cm^{1-\varepsilon} n + Cm \log n$. Further $\mathbb{E}\{N(u)\} = Nq_{\mathbf{r}}(u)$ and $\mathbb{E}\{N(u) \cdot \mathbf{1}_{\mathcal{G}}\} \geq \mathbb{E}\{N(u)\} - N\mathbb{P}(\mathcal{G}^c)$, whence

$$\liminf_{m,n \rightarrow \infty} \frac{1}{N} \mathbb{E}\{|S(u)| \cdot \mathbf{1}_{\mathcal{G}}\} \cdot \ell(u; \delta) \vee \ell_c(u; \delta) \geq q_{\mathbf{r}}(u). \quad (\text{F.23})$$

Recalling the definition of $\ell(u; \delta)$, $\ell_c(u; \delta)$ and the fact that δ is arbitrary, the last inequality yields

$$\liminf_{m,n \rightarrow \infty} \mathbb{E}\{|S(u)| \cdot \mathbf{1}_{\mathcal{G}}\} \frac{\log_2 N}{N} \geq q_{\mathbf{r}}(u) \left[H(X|U = u) \wedge \left(\frac{1 + \alpha}{\alpha} \right) H(X|U = u, V) \right]. \quad (\text{F.24})$$

Summing over u , noting that $\sum_{u \in \mathcal{L}} |S(u)| = M$, and substituting in Eq. (F.19) yields the lower bound on the rate in Eq. (5.7).

Finally, the upper bound is proved by a similar strategy as for the lower bound. Define the set of ‘bad indices’ $B_- = B_-(\mathbf{X}^{m,n}, \delta) \subseteq [m] \times [n]$,

$$B_-(\mathbf{X}^{m,n}, \delta) := \left\{ (i, j) \in [m] \times [n] : \mathcal{E}_{i,j}^c(\ell(-\delta, u_i)) \text{ or } \mathcal{F}_{i,j}^c(\ell_c(-\delta, u_i)) \right\} \quad (\text{F.25})$$

We also denote by L_i^+ the length of the last segment in row i . We then have

$$\begin{aligned}
 N(u) &\geq \sum_{(i,j) \in S(u)} L_{\langle ij \rangle} - \sum_{i \leq m: u_i = u} L_i^+ \\
 &\geq \sum_{(i,j) \in S(u) \cap B_-^c} L_{\langle ij \rangle} - \sum_{i \leq m: u_i = u} L_i^+ \\
 &\geq \sum_{(i,j) \in S(u) \cap B_-^c} \ell(u; -\delta) \vee \ell_c(u; -\delta) - \sum_{i \leq m: u_i = u} L_i^+ \\
 &\geq |S(u)| \ell(u; -\delta) \vee \ell_c(u; -\delta) - (|B_-| + m) \cdot C \log N,
 \end{aligned}$$

where the last inequality holds on event \mathcal{G} . By taking expectation on this event, we get

$$\mathbb{E}\{N(u) \cdot \mathbf{1}_{\mathcal{G}}\} \geq \mathbb{E}\{|S(u)| \cdot \mathbf{1}_{\mathcal{G}}\} \cdot \ell(-\delta, u) \vee \ell_c(-\delta, u) - (\mathbb{E}|B_-| + m) \cdot C \log N.$$

By Lemmas F.2 and F.2,

$$\begin{aligned} \mathbb{E}(|B_-|) &\leq m_0 n + \sum_{m_0 \leq i \leq m, j \leq n'} \mathbb{P}\left(\mathcal{E}_{i,j}^c(\ell(-\delta, u_i)) \cup \mathcal{F}_{i,j}^c(\ell_c(-\delta, u_i))\right) + Cm \log N \\ &\leq m_0 n + Cm^{1-\varepsilon} n + Cm \log N \\ &\leq \frac{CN}{(\log N)^2}. \end{aligned}$$

The proof is completed exactly as for the lower bound. \square

F.2. Proof of Lemma F.2

We will use the following standard lemmas.

Lemma F.4. *Let X be a centered random variable with $\mathbb{P}(X \leq x_0) = 1$, $x_0 > 0$. Then, letting $c(x_0) = (e^{x_0} - 1 - x_0)/x_0^2$, we have*

$$\mathbb{E}(e^X) \leq 1 + c(x_0)\mathbb{E}(X^2). \quad (\text{F.26})$$

Proof. This simply follows from $\exp(x) \leq 1 + x + c(x_0)x^2$ for $x \leq x_0$. \square

Lemma F.5. *Let $(p_i)_{i \geq 1}$, $(q_i)_{i \geq 1}$, be probability distributions on \mathcal{X} , with $\sup_{i \geq 1} \max_{x \in \mathcal{X}} p_i(x) \leq 1 - c$, and $\sup_{i \geq 1} \sum_{x \in \mathcal{X}} p_i(x)(\log p_i(x))^2 \leq C$ for constants c, C .*

Let $(X_i)_{i \leq \ell}$ be independent random variables with $X_i \sim p_i$, and set $\mathbf{X} = (X_1, \dots, X_\ell)$. Let $\mathbf{Y}(j) \in \mathcal{X}^\ell$, $j \geq 1$ be a sequence of i.i.d. random vectors, with $(Y_i(j))_{i \leq \ell}$ independent and $Y_i(j) \sim q_i$. Finally, let $T := \min\{t \geq 1 : \mathbf{Y}(t) = \mathbf{X}\}$.

Then, for any $\varepsilon > 0$, there exists $\delta = \delta(\varepsilon, c, C) > 0$ such that (letting $\bar{H}(p) := \ell^{-1} \sum_{i=1}^{\ell} H(p_i)$)

$$\mathbb{P}(T \leq e^{\ell[\bar{H}(p) - \varepsilon]}) \leq e^{-\delta \ell}. \quad (\text{F.27})$$

Further, the same bound holds (with a different $\delta(\varepsilon, c, C)$) $(\mathbf{Y}(j))_{j \geq 1}$ are independent not identically distributed, if there exist a finite set $(q_i^a)_{i \geq 1, a \in [K]}$, $K \leq \ell^{C_0}$, and a map $b : \mathbb{N} \rightarrow [K]$ such that $\mathbf{Y}(j) \sim q_1^{b(j)} \otimes \dots \otimes q_\ell^{b(j)}$.

Proof. We denote by \mathbf{Y} a vector distributed as $\mathbf{Y}(i)$. Conditional on $\mathbf{X} = \mathbf{x}$, T is a geometric random variables with mean $1/(1 - \mathbb{P}(\mathbf{Y} = \mathbf{x}))$. Hence, for $t_\ell(\varepsilon) := e^{\ell[\bar{H}(p) - \varepsilon]}$,

$$\begin{aligned} \mathbb{P}(T \leq t_\ell(\varepsilon) | \mathbf{X} = \mathbf{x}) &= 1 - (1 - \mathbb{P}(\mathbf{Y} = \mathbf{x}))^{t_\ell(\varepsilon)} \\ &\leq t_\ell(\varepsilon) \mathbb{P}(\mathbf{Y} = \mathbf{x}). \end{aligned}$$

Hence

$$\mathbb{P}(T \leq t_\ell(\varepsilon)) \leq e^{-\ell\varepsilon/2} + \mathbb{P}(\mathbb{P}(\mathbf{Y} = \mathbf{X} | \mathbf{X}) \geq t_\ell(\varepsilon/2)^{-1}) \quad (\text{F.28})$$

$$= e^{-\ell\varepsilon/2} + P_\ell(\varepsilon/2), \quad (\text{F.29})$$

$$P_\ell(u) := \mathbb{P}\left(\sum_{i=1}^{\ell} \log \frac{1}{q_i(X_i)} < \sum_{i=1}^{\ell} H(p_i) - \ell u\right). \quad (\text{F.30})$$

By Chernoff bound, for any $\lambda \geq 0$, $P_\ell(u) \leq \exp\{-\ell\phi(\lambda, u)\}$, where

$$\phi(\lambda, u) := \lambda u - \frac{1}{\ell} \sum_{i=1}^{\ell} [\lambda H(p_i) + \log \mathbb{E}[q_i(X_i)^\lambda]]. \quad (\text{F.31})$$

By Hölder inequality, for $\lambda \in [0, 1]$ we have $\mathbb{E}[q_i(X_i)^\lambda] \leq (\sum_x p(x)^\beta)^{1/\beta}$ where $\beta = 1/(1 - \lambda)$. Therefore

$$\begin{aligned} \psi(\lambda; p) &:= \lambda H(p) + (1 - \lambda) \log \left(\sum_{x \in \mathcal{X}} p(x)^{1/(1-\lambda)} \right) \\ &= (1 - \lambda) \log \mathbb{E}_{X \sim p} \exp \left(\frac{\lambda}{1 - \lambda} (\log p(X) + H(p)) \right). \end{aligned}$$

Consider the random variable $Z_i := \frac{\lambda}{1-\lambda} (\log p_i(X_i) + H(p))$ where $X_i \sim p_i$. Under the assumptions of the lemma, for $\lambda \in [0, 1/2]$ we have $\mathbb{E}(Z_i) = 0$ and

$$Z_i \leq \log(1-c) + H(p) \leq \log[|\mathcal{X}|(1-c)], \quad (\text{F.32})$$

$$\mathbb{E}[Z_i^2] \leq \left(\frac{\lambda}{1-\lambda}\right)^2 \sum_{x \in \mathcal{X}} p_i(x) (\log p_i(x))^2 \leq 4C\lambda^2, \quad (\text{F.33})$$

Using Lemma F.4, we get

$$\psi(\lambda; p_i) = (1-\lambda) \log \mathbb{E} e^{Z_i} \quad (\text{F.34})$$

$$\leq (1-\lambda) \log(1 + c_0 \mathbb{E}(Z_i^2)) \quad (\text{F.35})$$

$$\leq \log(1 + c_* \lambda^2), \quad (\text{F.36})$$

whence

$$\phi(\lambda, u) \geq \lambda u - \log(1 + c_* \lambda^2).$$

By maximizing this expression over λ , we find that $P_\ell(\varepsilon/2) \leq \exp(-\delta_0(\varepsilon)\ell)$ which completes the proof for the case of i.i.d. vectors $\mathbf{Y}(j)$.

The case of non-identically distributed vectors follows by union bound over $a \in [K]$. \square

Lemma F.6. *Let $(p_i)_{i \geq 1}$, be probability distributions on \mathcal{X} , with $\sup_{i \geq 1} \max_{x \in \mathcal{X}} p_i(x) \leq 1-c$, and $\sup_{i \geq 1} \sum_{x \in \mathcal{X}} p_i(x) (\log p_i(x))^2 \leq C$ for constants c, C .*

Let $(X_i)_{i \leq \ell}$ be independent random variables with $X_i \sim p_i$, $\mathbf{X} = (X_1, \dots, X_\ell)$. Let $\mathbf{Y}(j) \in \mathcal{X}^\ell$, $j \geq 1$ be a sequence of i.i.d. copies of \mathbf{X} . Finally, let $T := \min\{t \geq 1 : \mathbf{Y}(t) = \mathbf{X}\}$.

Then, for any $\varepsilon > 0$, there exists $\delta = \delta(\varepsilon, c, C) > 0$, such that (letting $\bar{H}(p) := \ell^{-1} \sum_{i=1}^\ell H(p_i)$)

$$\mathbb{P}(T \geq e^{\ell[\bar{H}(p)+\varepsilon]}) \leq e^{-\delta\ell}. \quad (\text{F.37})$$

Proof. The proof follows the same argument as for Lemma F.5. Denote by \mathbf{Y} a vector distributed as $\mathbf{Y}(i)$. and define $t_\ell(\varepsilon) := e^{\ell[\bar{H}(p)+\varepsilon]}$,

$$\begin{aligned} \mathbb{P}(T \geq t_\ell(\varepsilon) | \mathbf{X} = \mathbf{x}) &= (1 - \mathbb{P}(\mathbf{Y} = \mathbf{x}))^{t_\ell(\varepsilon)} \\ &\leq \exp\left(-t_\ell(\varepsilon) \mathbb{P}(\mathbf{Y} = \mathbf{x})\right). \end{aligned}$$

Hence

$$\mathbb{P}(T \geq t_\ell(\varepsilon)) \leq \exp\left\{-e^{\ell\varepsilon/2}\right\} + \mathbb{P}\left(\mathbb{P}(\mathbf{Y} = \mathbf{X} | \mathbf{X}) \leq t_\ell(\varepsilon/2)^{-1}\right) \quad (\text{F.38})$$

$$\leq e^{-\ell\varepsilon/2} + \tilde{P}_\ell(\varepsilon/2), \quad (\text{F.39})$$

$$\tilde{P}_\ell(u) := \mathbb{P}\left(\sum_{i=1}^\ell \log \frac{1}{p_i(X_i)} \geq \sum_{i=1}^\ell H(p_i) + \ell u\right). \quad (\text{F.40})$$

We claim that, for each $u > 0$, $\tilde{P}_\ell(u) \leq e^{-\delta_0(u)\ell}$ for some $\delta_0(u) > 0$. Indeed, using again Chernoff's bound, we get, for any $\lambda \geq 0$, $\tilde{P}_\ell(u) \leq e^{-\ell\tilde{\phi}(\lambda, u)}$, where

$$\tilde{\phi}(\lambda, u) := \lambda u - \frac{1}{\ell} \sum_{i=1}^\ell \tilde{\psi}(\lambda; p_i), \quad (\text{F.41})$$

$$\tilde{\psi}(\lambda; p_i) := \log \mathbb{E} \exp(\lambda W_i), \quad W_i := \log \frac{1}{p_i(X_i)} - H(p_i). \quad (\text{F.42})$$

where in the last line $X_i \sim p_i$. Under the assumptions of the lemma, $W_i \leq C$ almost surely and applying again Lemma F.4, we get $\tilde{\psi}(\lambda; p_i) \leq \log(1 + c_* \lambda^2)$ for $\lambda \leq 1$. The proof is completed by selecting for each $u > 0$, $\lambda > 0$ so that $\lambda u - \log(1 + c_* \lambda^2) > 0$. \square

We are now in position to prove Lemma F.2.

Proof of Lemma F.2. We begin by proving the bound (F.8).

Fix $i \leq m, j \leq n', \mathbf{u} \in \mathcal{L}^m, \delta > 0$, and write $\ell = \ell(\delta, u_i)$. Define $R_{ij} := \{(i, j') : \max(1, j - \ell) \leq j' \leq j - 1\}$ and $S_{ij} := \{(i', j') : i' < i \text{ or } i' = i, j' < j - \ell\}$. Finally, for $t \in \{0, \dots, \ell - 1\}$, let $S_{ij}(t) := S_{ij} \cap \{(i', j') : \langle i' j' \rangle = t \pmod{\ell}\}$.

By union bound

$$\begin{aligned} \mathbb{P}(\mathcal{E}_{i,j}(\ell) | \mathbf{u}) &\leq A + \sum_{t=0}^{\ell-1} B(t), \\ A &:= \sum_{(rs) \in R_{ij}} \mathbb{P}(\mathbf{X}_{(rs)}^{\langle rs \rangle + \ell - 1} = \mathbf{X}_{(ij)}^{\langle ij \rangle + \ell - 1} | \mathbf{u}), \\ B(t) &:= \mathbb{P}(\exists (r, s) \in S_{ij}(t) : \mathbf{X}_{(rs)}^{\langle rs \rangle + \ell - 1} = \mathbf{X}_{(ij)}^{\langle ij \rangle + \ell - 1} | \mathbf{u}). \end{aligned}$$

Now, by the bound of Eq. (F.4),

$$A \leq \ell \cdot (1 - c)^\ell \leq C N^{-\varepsilon}, \quad (\text{F.43})$$

for suitable constants C, ε .

Next, for any $t \in \{0, \dots, \ell - 1\}$, the vectors $\{\mathbf{X}_{(rs)}^{\langle rs \rangle + \ell - 1}\}$ are mutually independent and independent of $\{\mathbf{X}_{(ij)}^{\langle ij \rangle + \ell - 1}\}$.

Conditional on \mathbf{u} , the coordinates of $\mathbf{X}_{(rs)}^{\langle rs \rangle + \ell - 1} = (X_{(rs)}, \dots, X_{(rs) + \ell - 1})$ are independent with marginal distributions $X_{\langle r' s' \rangle} \sim Q_{x|u}(\cdot | u_{r'})$ (note that independence of the coordinates holds because $\ell < m/2$ and therefore $\mathbf{X}_{(rs)}^{\langle rs \rangle + \ell - 1}$ does not include two entries in the same column). Note that the collection of marginal distributions $Q_{x|u}(\cdot | u)$, $u \in \mathcal{L}$ satisfies the conditions of Lemma F.5 by assumption. Further, the vector $\mathbf{X}_{(rs)}^{\langle rs \rangle + \ell - 1}$ can have at most one of $K = |\mathcal{L}|^2(\ell + 1)$ distributions (depending on the latents value and the occurrence of a line break in the block.)

Applying Lemma F.5, we obtain:

$$B(t) \leq e^{-\varepsilon_0 \ell} \leq C N^{-\varepsilon} \quad (\text{F.44})$$

Summing over $t \in \{0, \dots, \ell - 1\}$ and adjusting the constants yields the claim (F.8).

Next consider the bound (F.9). Fix $\mathbf{u} \in \mathcal{L}^m, i \leq m, j \leq n'$, and write $\ell = \ell(-\delta, u_i)$ for brevity below

$$\mathbb{P}(\mathcal{E}_{i,j}^c(\ell) | \mathbf{u}) \leq \mathbb{P}(\forall (i', j') \in S_{ij}(t) \text{ s.t. } u_{i'} = u_i, j' < n' : \mathbf{X}_{(i' j')}^{\langle i' j' \rangle + \ell - 1} \neq \mathbf{X}_{(ij)}^{\langle ij \rangle + \ell - 1} | \mathbf{u}). \quad (\text{F.45})$$

Here $t \in \{0, \dots, \ell - 1\}$ can be chosen arbitrarily. Let $S_{ij}(t; \mathbf{u}) := \{(i', j') \in S_{ij}(t) \text{ s.t. } u_{i'} = u_i, j' < n'\}$. Conditional on \mathbf{u} , the vectors $(\mathbf{X}_{(i' j')}^{\langle i' j' \rangle + \ell - 1})_{(i', j') \in S_{ij}(t; \mathbf{u})}$ are i.i.d. and independent of $\mathbf{X}_{(ij)}^{\langle ij \rangle + \ell - 1}$. Further, they are distributed as $\mathbf{X}_{(ij)}^{\langle ij \rangle + \ell - 1}$. Finally,

$$\begin{aligned} N_{ij}(\mathbf{u}) &:= |S_{ij}(t; \mathbf{u})| \geq \frac{n(m_i(u) - C \log N)}{\ell} \\ &\geq \frac{m_i(u)n}{C \log N} - C'n. \end{aligned}$$

where $m_i(u)$ is the number rows $i' < i$ such that $u_{i'} = u$. Since $i \geq i$ and $\mathbb{P}(u_{i'} = u) \geq \min_{u'} q_{\mathbf{r}}(u') > 0$, by Chernoff bound there exist constants C, c_0 such that, for all m, n large enough (since $i \geq m_0$)

$$\mathbb{P}(N_{ij}(\mathbf{u}) \geq \frac{c_0 m_0 n}{\log N}) \geq 1 - C e^{-m_0/C}. \quad (\text{F.46})$$

Further, for any $\delta > 0$ we can choose positive constants $\varepsilon_0, \varepsilon_1 > 0$ such that the following holds for all m, n , large enough

$$\frac{c_0 m_0 n}{\log N} \geq N^{1-\varepsilon_1} \geq e^{\ell[H(X|U=u_i) + \varepsilon_0]} \quad (\text{F.47})$$

Let T_{ij} be the rank of the first (i', j') (moving backward) in the set defined above such that $\mathbf{X}_{\langle i'j' \rangle}^{\langle i'j' \rangle + \ell - 1} = \mathbf{X}_{\langle ij \rangle}^{\langle ij \rangle + \ell - 1}$, and $T_{ij} = \infty$ if no such vector exists. We can continue from Eq. (F.45) to get

$$\begin{aligned} \mathbb{P}(\mathcal{E}_{i,j}^c(\ell)) &\leq \mathbb{P}(T_{ij} \geq N_{ij}(\mathbf{u})) \\ &\leq \mathbb{P}(T_{ij} \geq N_{ij}(\mathbf{u}); N_{ij}(\mathbf{u}) \geq e^{\ell[H(X|U=u_i)+\varepsilon_0]}) + Ce^{-m_0/C} \\ &\stackrel{(a)}{\leq} \exp\left\{-\delta_0 \min_{u \in \mathcal{L}}(\ell(-\delta; u))\right\} + Ce^{-m_0/C} \leq CN^{-\varepsilon}, \end{aligned}$$

where in (a) we used Lemma F.6. This completes the proof of Eq. (F.9). \square

F.3. Proof of Lemma F.3

We begin by considering the bound (F.10).

Fix $i \leq m$, $j \leq n'_c$, $\mathbf{u} \in \mathcal{L}^m$, $\mathbf{v} \in \mathcal{L}^n$, $\delta > 0$, and write $\ell = \ell_c(\delta, u_i)$, $n' = n'_c$. By union bound:

$$\begin{aligned} \mathbb{P}(\mathcal{F}_{i,j}(\ell)|\mathbf{u}, \mathbf{v}) &= \mathbb{P}\left(\bigcup_{s \in [n], |j-s| < \ell} \mathcal{B}(s) \mid \mathbf{u}, \mathbf{v}\right), \\ \mathcal{B}(s) &:= \left\{ \exists r < i : \mathbf{X}_{\langle rs \rangle}^{\langle rs \rangle + \ell - 1} = \mathbf{X}_{\langle ij \rangle}^{\langle ij \rangle + \ell - 1} \right\}. \end{aligned}$$

Note that for a fixed s , and conditional on \mathbf{u}, \mathbf{v} , the vectors $(\mathbf{X}_{\langle rs \rangle}^{\langle rs \rangle + \ell - 1})_{1 \leq s \leq i-1}$ are mutually independent and independent of $\mathbf{X}_{\langle ij \rangle}^{\langle ij \rangle + \ell - 1}$. Further, $\mathbf{X}_{\langle rs \rangle}^{\langle rs \rangle + \ell - 1}$ has independent coordinates with marginals $X_{\langle r's' \rangle} \sim Q_{x|u}(\cdot | u_{r'}, v_{s'})$ (recall that we are conditioning both on \mathbf{u} and \mathbf{v}). In particular, the marginal distributions satisfy the assumption of Lemma F.5 and the law of $\mathbf{X}_{\langle rs \rangle}^{\langle rs \rangle + \ell - 1}$ can take one of $K = |\mathcal{L}|^2(\ell + 1)$ possible values. Letting $i - T(s)$ the last row at which $\mathbf{X}_{\langle rs \rangle}^{\langle rs \rangle + \ell - 1} = \mathbf{X}_{\langle ij \rangle}^{\langle ij \rangle + \ell - 1}$ (with $T(s) \geq i$ if no such row exists), we have, for some constants $C, c_0 > 0$,

$$\begin{aligned} \mathbb{P}(\mathcal{F}_{i,j}(\ell)|\mathbf{u}, \mathbf{v}) &\leq \mathbb{P}\left(\bigcup_{s \in [n], |j-s| < \ell} \{T(s) \leq i-1\} \cap \{i-1 \leq e^{\ell[\bar{H}-\varepsilon]}\} \mid \mathbf{u}, \mathbf{v}\right) + \mathbf{1}(i-1 > e^{\ell[\bar{H}-\varepsilon]}) \\ &\stackrel{(a)}{\leq} 2\ell e^{-\ell\varepsilon} + \mathbf{1}(m > e^{\ell[\bar{H}-\varepsilon]}) \\ &\leq Cm^{-c_0\varepsilon} + \mathbf{1}(m > e^{\ell[\bar{H}-\varepsilon]}), \end{aligned}$$

where in (a) we used Lemma F.5, and we defined $\bar{H} := \ell^{-1} \sum_{k=j}^{j+\ell-1} H(X|U = u_i, V = v_k)$.

Taking expectation with respect to \mathbf{v} , we get

$$\begin{aligned} \mathbb{P}(\mathcal{F}_{i,j}(\ell)|\mathbf{u}) &\leq Cm^{-c_0\varepsilon} + \mathbb{P}\left(\frac{1}{\ell} \sum_{k=j}^{j+\ell-1} H(X|U = u_i, V = v_k) < \frac{1}{1+\delta}(H(X|U = u_i, V) + \varepsilon)\right) \\ &\stackrel{(a)}{\leq} Cm^{-c_0\varepsilon} + e^{-\ell\varepsilon} \leq C'm^{-c_0\varepsilon}, \end{aligned}$$

where in (a) we used Chernoff bound. This completes the proof of Eq. (F.10).

Finally, the proof Eq. (F.11) is similar to the one of Eq. (F.9). We fix $\mathbf{u} \in \mathcal{L}^m$, $i \leq m$, $j \leq n'_c$, and write $\ell = \ell_c(-\delta, u_i)$.

$$\mathbb{P}(\mathcal{F}_{i,j}^c(\ell)|\mathbf{u}) \leq \mathbb{P}\left(\forall i' < i u_{i'} = u_i : \mathbf{X}_{\langle i'j \rangle}^{\langle i'j \rangle + \ell - 1} \neq \mathbf{X}_{\langle ij \rangle}^{\langle ij \rangle + \ell - 1} \mid \mathbf{u}\right). \quad (\text{F.48})$$

Let $S_{ij}^c(\mathbf{u}) := \{(i', j) \text{ s.t. } u_{i'} = u_i, i' < i\}$. Conditional on \mathbf{u}, \mathbf{v} , the vectors $(\mathbf{X}_{\langle i'j \rangle}^{\langle i'j \rangle + \ell - 1})_{(i', j) \in S_{ij}^c(\mathbf{u})}$ are i.i.d. and independent copies of $\mathbf{X}_{\langle ij \rangle}^{\langle ij \rangle + \ell - 1}$. Finally, $N_i^c(\mathbf{u}) := |S_{ij}^c(\mathbf{u})|$ is the number rows $i' < i$ such that $u_{i'} = u_i$. By Chernoff bound there exist constants C, c_0 such that, for all m, n large enough (recalling that we need only to consider $i \geq m_0$)

$$\mathbb{P}\left(N_i^c(\mathbf{u}) \geq c_0 m_0\right) \geq 1 - Ce^{-m_0/C}. \quad (\text{F.49})$$

Since $m_0 \geq m^{1-o_n(1)}$, for any $\delta > 0$ we can choose constants $\varepsilon_0, \varepsilon_1 > 0$ so that

$$c_0 m_0 \geq m^{1-\varepsilon_1} \geq e^{\ell[H(X|U=u_i, V)+2\varepsilon_0]}. \quad (\text{F.50})$$

Recall the definition $\bar{H} := \ell^{-1} \sum_{k=j}^{j+\ell-1} H(X|U = u_i, V = v_k)$. By an application of Chernoff bound

$$\mathbb{P}\left(N_i^c(\mathbf{u}) \geq e^{\ell[H(X|U=u_i, V)+\varepsilon_0]}\right) \geq 1 - Cm^{-\varepsilon} - Ce^{-m_0/C}. \quad (\text{F.51})$$

Let T_i be the rank of the first i' (moving backward) in the set defined above such that $\mathbf{X}_{\langle i'j \rangle}^{(i'j)+\ell-1} = \mathbf{X}_{\langle ij \rangle}^{(ij)+\ell-1}$, and $T_i = \infty$ if no such vector exists. From Eq. (F.48) we get

$$\begin{aligned} \mathbb{P}(\mathcal{F}_{i,j}^c(\ell)) &\leq \mathbb{P}(T_i \geq N_i(\mathbf{u})) \\ &\leq \mathbb{P}(T_i \geq N_i(\mathbf{u}); N_i(\mathbf{u}) \geq e^{\ell[H(X|U=u_i, V)+\varepsilon_0]}) + Cm^{-\varepsilon} \\ &\stackrel{(a)}{\leq} \exp\left\{-\delta_0 \min_{u \in \mathcal{L}}(\ell_c(-\delta; u))\right\} + Cm^{-\varepsilon} \leq 2Cm^{-\varepsilon}, \end{aligned}$$

where in (a) we used Lemma F.6.

G. Proofs for latent-based encoders

G.1. Proof of Lemma 5.7

G.1.1. GENERAL BOUND (5.8)

From Eq. (1.3), we get

$$R_{\text{lat}}(\mathbf{X}) = \frac{1}{mn \log_2 |\mathcal{X}|} \left\{ \text{len}(\text{header}) + \text{len}(Z_{\mathcal{L}}(\hat{\mathbf{u}})) + \text{len}(Z_{\mathcal{L}}(\hat{\mathbf{v}})) + \sum_{u,v \in \mathcal{L}} \text{len}(Z_{\mathcal{X}}(\hat{\mathbf{X}}(u, v))) \right\} \wedge 1, \quad (\text{G.1})$$

where $\hat{\mathbf{X}}(u, v) := \text{vec}(X_{ij} : \hat{u}_i(\mathbf{X}) = u, \hat{v}_j(\mathbf{X}) = v)$ are the estimated blocks of \mathbf{X} . Note that this rate depends on the base compressors $Z_{\mathcal{L}}, Z_{\mathcal{X}}$ but we will omit these from our notations.

Define the ‘ideal’ expected compression rate (i.e. the rate achieved by a compressor that is given the latents):

$$R_{\#} := \frac{1}{mn \log_2 |\mathcal{X}|} \left\{ \mathbb{E}[\text{len}(\text{header})] + \mathbb{E}[\text{len}(Z(\mathbf{u}))] + \mathbb{E}[\text{len}(Z(\mathbf{v}))] + \sum_{u,v \in \mathcal{L}} \mathbb{E}[\text{len}(Z(\mathbf{X}(u, v)))] \right\}.$$

Since $R_{\text{lat}}(\mathbf{X}) \leq 1$ by construction, we have

$$\begin{aligned} \mathbb{E} R_{\text{lat}}(\mathbf{X}) &\leq \mathbb{E} \left\{ R_{\text{lat}}(\mathbf{X}) \mathbf{1}_{\text{Err}_U(\mathbf{X}; \hat{\mathbf{u}})=1} \mathbf{1}_{\text{Err}_V(\mathbf{X}; \hat{\mathbf{v}})=1} \right\} + \mathbb{P}(\text{Err}_U(\mathbf{X}; \hat{\mathbf{u}}) < 1) + \mathbb{P}(\text{Err}_V(\mathbf{X}; \hat{\mathbf{v}}) < 1) \\ &\stackrel{(*)}{\leq} R_{\#} + \mathbb{P}(\text{Err}_U(\mathbf{X}; \hat{\mathbf{u}}) < 1) + \mathbb{P}(\text{Err}_V(\mathbf{X}; \hat{\mathbf{v}}) < 1), \end{aligned}$$

where in step (*) we bounded $\mathbb{E}[\text{len}(Z(\hat{\mathbf{u}})) \mathbf{1}_{\text{Err}_U(\mathbf{X}; \hat{\mathbf{u}})=1}] = \mathbb{E}[\text{len}(Z(\mathbf{u})) \mathbf{1}_{\text{Err}_U(\mathbf{X}; \hat{\mathbf{u}})=1}] \leq \mathbb{E}[\text{len}(Z(\mathbf{u}))]$, because, on the event $\{\text{Err}_U(\mathbf{X}; \hat{\mathbf{u}}) = 1\}$, $\hat{\mathbf{u}}$ coincides with \mathbf{u} up to relabelings, and the compressed length is invariant under relabelings. Similar arguments were applied to $\text{len}(Z(\mathbf{v}))$ and $\text{len}(Z(\mathbf{X}(u, v)))$.

We now have, by the definition of $\Delta_Z(N; k)$ in Eq. (5.12),

$$\frac{\mathbb{E}[\text{len}(Z(\mathbf{u}))]}{mn \log_2 |\mathcal{X}|} \leq \frac{H(U)}{n \log_2 |\mathcal{X}|} + \frac{1}{n} \Delta_Z(m \wedge n; \{r, c\}), \quad (\text{G.2})$$

$$\frac{\mathbb{E}[\text{len}(Z(\mathbf{v}))]}{mn \log_2 |\mathcal{X}|} \leq \frac{H(V)}{m \log_2 |\mathcal{X}|} + \frac{1}{m} \Delta_Z(m \wedge n; \{r, c\}), \quad (\text{G.3})$$

$$\frac{\mathbb{E}[\text{len}(Z(\mathbf{X}(u, v))) | \mathbf{u}, \mathbf{v}]}{mn \log_2 |\mathcal{X}|} \leq \hat{q}_r(u) \hat{q}_c(v) \frac{H(X|U=u, V=v)}{\log_2 |\mathcal{X}|} + \Delta_Z(c \cdot mn; \{Q(\cdot | u, v)\}_{i,v \in \mathcal{L}}), \quad (\text{G.4})$$

where in the last line \hat{r} is the empirical distribution of the row latents and \hat{c} is the empirical distribution of the column latents. By taking expectation in the last expression, we get

$$\sum_{u,v \in \mathcal{L}} \frac{\mathbb{E}[\text{len}(Z(\mathbf{X}(u, v))) | \mathbf{u}, \mathbf{v}]}{mn \log_2 |\mathcal{X}|} \leq \frac{H(X|U, V)}{\log_2 |\mathcal{X}|} + |\mathcal{L}|^2 \Delta_Z(c \cdot mn; \{Q(\cdot | u, v)\}_{u,v \in \mathcal{L}}). \quad (\text{G.5})$$

Finally, the header contains $|\mathcal{L}|^2 + 2$ integers of maximum size mn , whence $\text{len}(\text{header}) \leq 4 \log_2(mn)$. We conclude that

$$\begin{aligned} R_{\#} \leq & \frac{1}{\log_2 |\mathcal{X}|} \left\{ H(X|U, V) + \frac{1}{n} H(U) + \frac{1}{n} H(V) \right\} + \frac{2 \log_2(mn)}{mn} \\ & + |\mathcal{L}|^2 \Delta_Z(c \cdot mn; \{Q(\cdot | u, v)\}_{u, v \in \mathcal{L}}) + 2 \Delta_Z(m \wedge n; \{r, c\}). \end{aligned}$$

The claim (5.8) follows from the first bound in Eq. (5.2) noticing that, under the stated assumptions on m, n ,

$$\frac{1}{n} [\mathfrak{h}(\varepsilon_U) + \varepsilon_u \log(|\mathcal{L}| - 1)] \leq \varepsilon_U \leq \mathbb{P}(\text{Err}_U(\mathbf{X}^{m,n}; \hat{\mathbf{u}}) < 1). \quad (\text{G.6})$$

G.1.2. REDUNDANCY BOUNDS FOR SPECIFIC ENCODERS: EQS. (5.9)–(5.11)

LZ coding. Let $\mathbf{X}^N = (X_1, \dots, X_N)$ be a vector with i.i.d. symbols $X_i \sim q$ with q a probability distribution over \mathcal{X} . The analysis is similar to the one in Appendix F, and we will adopt the same notations here. There are two important differences: data are i.i.d. (not matrix-structured) and we want to derive a sharper estimate (not just the entropy term, but bounding the overhead as well).

We define $L_k(\mathbf{X}^N), T_k(\mathbf{X}^N)$ as per Eqs. (F.1), (F.2). We let $(k(1), \dots, k(M_N))$ be the values taken by k in the while loop of the Lempel-Ziv pseudocode of Section 5.2.2. In particular

$$k(1) = 1, \quad (\text{G.7})$$

$$k(\ell + 1) = k(\ell) + L_{k(\ell)}(\mathbf{X}^N), \quad (\text{G.8})$$

$$k(M_N) = N. \quad (\text{G.9})$$

(We set $k(0) = 0$ by convention.) Therefore the total length of the code is

$$\begin{aligned} \text{len}(\text{LZ}(\mathbf{X}^N)) &= M_N \lceil \log_2(N + |\mathcal{X}|) \rceil + \sum_{\ell=1}^{M_N} \text{len}(\text{elias}(L_{k(\ell)})) \\ &\leq M_N \lceil \log_2(N + |\mathcal{X}|) \rceil + 2 \sum_{\ell=1}^{M_N} \log_2(L_{k(\ell)}) \\ &\leq M_N \lceil \log_2(N + |\mathcal{X}|) \rceil + 2M_N \log_2(N/M_N), \end{aligned}$$

where the last step follows by Jensen's inequality. By one more application of Jensen, we obtain

$$\mathbb{E} R_{\text{LZ}}(\mathbf{X}^N) \leq \frac{1}{\log_2 |\mathcal{X}|} \cdot \frac{\mathbb{E} M_N}{N} \cdot \{ \lceil \log_2(N + |\mathcal{X}|) \rceil + 2 \log_2(N/\mathbb{E} M_N) \}. \quad (\text{G.10})$$

Define the set of break points and bad positions as

$$S_N := \{k(1), k(2), \dots, k(M_N)\}, \quad (\text{G.11})$$

$$B_N(\ell) := \{k \in [N/2, N] : L_k(\mathbf{X}^N) < \ell\}. \quad (\text{G.12})$$

Note that $S_N = S_N^{\leq} \cup S_N^{\gt}$ where:

$$S_N^{\leq} := \{k(j) : j \leq M_N, k(j-1) \leq \lfloor N/2 \rfloor\}, \quad S_N^{\gt} := \{k(j) : j \leq M_N, k(j-1) > \lfloor N/2 \rfloor\}. \quad (\text{G.13})$$

Further $|S_N^{\leq}| \stackrel{\text{d}}{=} M_{\lfloor N/2 \rfloor}$ and, for any $\ell \in \mathbb{N}$,

$$\frac{N}{2} \geq \sum_{k \in S_N^{\gt}} L_k \geq (|S_N^{\gt}| - |B_N(\ell)|) \ell. \quad (\text{G.14})$$

Therefore,

$$\begin{aligned} \mathbb{E} |S_N^>| &\leq \frac{N}{2\ell} + \mathbb{E} |B_N(\ell)| \\ &\leq \frac{N}{2\ell} + \sum_{k=\lceil N/2 \rceil}^N \mathbb{P}(L_k(\mathbf{X}^N) < \ell). \end{aligned} \quad (\text{G.15})$$

We claim that this implies, for $C_0 = 20c_* \log |\mathcal{X}|$ and $\log N \geq (2 \log(2/H(q)))^2$,

$$\frac{1}{N} \mathbb{E} |S_N^>| \leq \frac{H(q)}{2 \log_2 N} + C_0 \frac{(\log \log_2 N)^{1/2}}{(\log_2 N)^{3/2}} =: \psi(\log_2 N). \quad (\text{G.16})$$

Before proving this claim, let us show that it implies the thesis. Recall that $M_N = |S_N|$ and $S_N = S_N^{\leq} \cup S_N^>$ where $|S_N^{\leq}| \stackrel{d}{=} M_{\lfloor N/2 \rfloor}$. Therefore, we have proven

$$\begin{aligned} \mathbb{E} M_N &\leq N \psi(\log_2 N) + \mathbb{E} M_{\lfloor N/2 \rfloor} \\ &\leq \sum_{\ell=0}^{K-1} N_\ell \psi(\log_2 N_\ell) + \mathbb{E} M_{N_K}, \end{aligned} \quad (\text{G.17})$$

where we defined recursively $N_0 = N$, $N_{\ell+1} = \lfloor N/2 \rfloor$, and $K := \min\{\ell : \log_2 N_\ell < (2 \log(2/H(q)))^2\}$. Of course, $M_{N_K} \leq N_K \leq \exp((2 \log(2/H(q)))^2)$. Further $\underline{N}_\ell \leq N_\ell \leq \bar{N}_\ell$, where $\underline{N}_0 = \bar{N}_0 = N$ and $\bar{N}_{\ell+1} = \bar{N}_\ell/2$, $\underline{N}_{\ell+1} = (\underline{N}_\ell - 1)/2$ for $\ell \geq 0$. We thus get $\underline{N}_\ell = (N+1)2^{-\ell} - 1$, $\bar{N}_\ell = N 2^{-\ell}$ and therefore

$$\begin{aligned} \frac{1}{N} \sum_{\ell=0}^{K-1} N_\ell \psi(\log_2 N_\ell) &\leq \frac{1}{N} \sum_{\ell=0}^{\infty} \bar{N}_\ell \psi(\log_2 \bar{N}_\ell) \\ &\leq \frac{H(q)}{2} \sum_{\ell=0}^{\infty} 2^{-\ell} \frac{1}{\log_2(N 2^{-\ell} - 1)} + C_0 \sum_{\ell=0}^{\infty} 2^{-\ell} \frac{(\log \log_2 N)^{1/2}}{(\log_2(N 2^{-\ell} - 1))^{3/2}} \\ &\leq \frac{H(q)}{\log_2 N} + 2C_0 \frac{(\log \log_2 N)^{1/2}}{(\log_2 N)^{3/2}}. \end{aligned}$$

Substituting in Eq. (G.17), we get

$$\begin{aligned} \frac{1}{N} \mathbb{E} M_N &\leq \frac{H(q)}{\log_2 N} + 2C_0 \frac{(\log \log_2 N)^{1/2}}{(\log_2 N)^{3/2}} + \frac{1}{N} \exp\{(2 \log(2/H(q)))^2\} \\ &\leq \frac{H(q)}{\log_2 N} + 3C_0 \frac{(\log \log_2 N)^{1/2}}{(\log_2 N)^{3/2}}, \end{aligned}$$

where the last inequality follows for $N \geq \exp\{(4 \log(2/H(q)))^2\}$ (noting that $C_0 > 1$). Finally, the desired bound (5.9) follows by substituting the last estimate in Eq. (G.10).

We are left with the task of proving claim (G.16). Fix any k , $\lceil N/2 \rceil \leq k \leq N$ and write q^ℓ for the product distribution $q \times \dots \times q$ (ℓ times). Setting $H = H_{\text{nats}}(q)$ (measuring here entropy in nats), for any $\delta > 0$,

$$\begin{aligned} \mathbb{P}(L_k(\mathbf{X}^N) < \ell) &= \mathbb{P}(\mathbf{X}_i^{i+\ell-1} \neq \mathbf{X}_k^{k+\ell-1} \forall i < k) \\ &\leq \sum_{\mathbf{x}^\ell \in \mathcal{X}^\ell} \mathbb{P}(Z_k(\mathbf{x}^\ell) = 0) \cdot \mathbb{P}(\mathbf{X}_k^{k+\ell-1} = \mathbf{x}^\ell) \\ &\leq \sum_{\mathbf{x}^\ell \in \mathcal{X}^\ell} q^\ell(\mathbf{x}^\ell) (1 - q^\ell(\mathbf{x}^\ell))^{N/2\ell} \\ &\leq \sum_{\mathbf{x}^\ell \in \mathcal{X}^\ell} q^\ell(\mathbf{x}^\ell) \mathbf{1}(q^\ell(\mathbf{x}^\ell) \leq e^{-\ell[H+\delta]}) + \exp\left\{-\frac{N}{2\ell} \cdot e^{-\ell[H+\delta]}\right\} \\ &=: P_{\leq}(\ell; \delta) + P_{>}(\ell, N; \delta). \end{aligned} \quad (\text{G.18})$$

By Chernoff bound

$$P_{\leq}(\ell; \delta) \leq e^{-\ell \max_{\lambda > 0} \psi_{\delta}(\lambda)},$$

$$\psi_{\delta}(\lambda) := \lambda[H + \delta] - \log \left\{ \sum_{x \in \mathcal{X}} q(x)^{1-\lambda} \right\}.$$

Note that $\lambda \mapsto \psi_{\delta}(\lambda)$ is continuous, concave, with $\psi'_{\delta}(0) = \delta$, $\psi_{\delta}(0) = 0$, $\psi_{\delta}(1) = \delta + H - \log |\mathcal{X}|$. Hence (assuming $H < \log |\mathcal{X}|$ because otherwise there is nothing to prove) for all δ small enough ψ is maximized for $\lambda \in (0, 1)$. Further, defining the random variable $Q = q(x)$ for $x \sim \text{Unif}(\mathcal{X})$,

$$\psi''_{\delta}(\lambda) = -\frac{\mathbb{E}[Q^{1-\lambda}(\log Q)]}{\mathbb{E}[Q^{1-\lambda}]} + \left[\frac{\mathbb{E}[Q^{1-\lambda}(\log Q)^2]}{\mathbb{E}[Q^{1-\lambda}]} \right]^2 \quad (\text{G.20})$$

$$\geq -\frac{\mathbb{E}[Q^{1-\lambda}(\log Q)^2]}{\mathbb{E}[Q^{1-\lambda}]} \quad (\text{G.21})$$

$$\geq -\mathbb{E}[(\log Q)^2] =: -c_*(q). \quad (\text{G.22})$$

Here the last inequality holds because $Q \mapsto Q^{1-\lambda}$ is monotone increasing (for $\lambda \in [0, 1]$) and $Q \mapsto (\log Q)^2$ is monotone decreasing over $Q \in [0, 1]$, and therefore $\mathbb{E}[Q^{1-\lambda}(\log Q)^2] \leq \mathbb{E}[Q^{1-\lambda}] \mathbb{E}[(\log Q)^2]$. In what follows, we set $c_* := c_*(q) \wedge 1$. Hence $\psi_{\delta}(\lambda) \geq \delta\lambda - c_*\lambda^2/2$ for $\lambda \in [0, 1]$ and therefore using Eq. (G.20),

$$P_{\leq}(\ell; \delta) \leq \exp \left\{ -\ell \min \left(\frac{\delta^2}{2c_*}, \delta - \frac{c_*}{2} \right) \right\}.$$

Substituting in Eq. (G.19), and using this in Eq. (G.15), we get, for $\delta \in [0, c_*]$:

$$\frac{1}{N} \mathbb{E} |S_N^>| \leq \frac{1}{2\ell} + \exp \left\{ -\frac{\ell\delta^2}{2c_*} \right\} + \exp \left\{ -\frac{N}{2\ell} \cdot e^{-\ell[H+\delta]} \right\}$$

We set

$$\ell = \frac{\log N}{H} (1 - \varepsilon), \quad \delta = \frac{1}{2} H \varepsilon,$$

for $\varepsilon \leq (2c_*/H) \wedge (1/2)$. Substituting in the previous bound, we get

$$\frac{1}{N} \mathbb{E} |S_N^>| \leq \frac{H}{2 \log N} (1 + 2\varepsilon) + \exp \left\{ -\frac{H\varepsilon^2 \log N}{16c_*} \right\} + \exp \left\{ -\frac{HN^{(\varepsilon+\varepsilon^2)/2}}{2 \log N} \right\}.$$

We finally select $\varepsilon = c_0(c_* \log |\mathcal{X}|/H)(\log \log N / \log N)^{1/2}$, with c_0 a sufficiently small absolute constant. Substituting above,

$$\begin{aligned} \frac{1}{N} \mathbb{E} |S_N^>| - \frac{H}{2 \log N} &\leq c_0 c_* \log |\mathcal{X}| \frac{(\log \log N)^{1/2}}{(\log N)^{3/2}} + \exp \left\{ -\frac{c_0^2 (\log |\mathcal{X}|)^2 c_*}{16H} \log \log N \right\} \\ &\quad + \exp \left\{ -\frac{H}{2 \log N} e^{(c_0 c_* \log |\mathcal{X}|/2H)\sqrt{\log N}} \right\} \\ &\leq c_0 c_* \log |\mathcal{X}| \frac{(\log \log N)^{1/2}}{(\log N)^{3/2}} + (\log N)^{-c_0^2/16} + \exp \left\{ -\frac{H}{2 \log N} e^{(c_0/2)\sqrt{\log N}} \right\}. \end{aligned}$$

Setting $c_0 = 5$, we get

$$\frac{1}{N} \mathbb{E} |S_N^>| - \frac{H}{2 \log N} \leq 6 c_* \log |\mathcal{X}| \frac{(\log \log N)^{1/2}}{(\log N)^{3/2}} + \exp \left\{ -\frac{H}{2 \log N} e^{2\sqrt{\log N}} \right\},$$

whence, the claim (G.16) follows for $\log N \geq (\log(2/H))^2$.

Arithmetic Coding. In Arithmetic Coding (AC) we encode the empirical distribution of $\mathbf{X}^N \hat{q}_N(x) := N^{-1} \sum_{i \leq N} \mathbf{1}_{X_i=x}$, and then encode \mathbf{X}^N in at most $-\log_2 \hat{q}_N(\mathbf{X}^N) + 1$ bits. The encoding of \hat{q}_N amounts to encoding the $|\mathcal{X}| - 1$ integers $N\hat{q}_N(x)$, $x \in \mathcal{X} \setminus \{0\}$ (assuming that $0 \in \mathcal{X}$, one of the counts can be obtained by difference). We thus have

$$\begin{aligned} \text{len}(\text{AC}(\mathbf{X}^N)) &\leq -\log_2 \hat{q}_N^{\otimes N}(\mathbf{X}^N) + 1 + \sum_{x \in \mathcal{X}} \text{len}(\text{elias}(N\hat{q}_N(x))) \\ &\leq -\log_2 \hat{q}_N^{\otimes N}(\mathbf{X}^N) + 2|\mathcal{X}| \log_2 N \\ &= \sum_{i=1}^N -\log_2 \hat{q}_N(X_i) + 2|\mathcal{X}| \log_2 N \\ &= N H(\hat{q}_N) + 2|\mathcal{X}| \log_2 N. \end{aligned}$$

Taking expectations

$$\begin{aligned} \mathbb{E}R_{\text{AC}}(\mathbf{X}^N) &\leq \frac{\mathbb{E} H(\hat{q}_N)}{\log_2 |\mathcal{X}|} + \frac{2|\mathcal{X}| \log_2 N}{N \log_2 |\mathcal{X}|} \\ &\leq \frac{H(q)}{\log_2 |\mathcal{X}|} + \frac{2|\mathcal{X}| \log_2 N}{N \log_2 |\mathcal{X}|}. \end{aligned}$$

ANS Coding. The bound (5.11) follows for range ANS coding from the analysis of (Duda, 2009; 2013; Kosolobov, 2022), where encoding of empirical distributions are analyzed as for AC coding.

G.2. Proof of Theorem 5.8

The proof consists in applying Lemma 5.7 and showing that $\mathbb{P}(\text{Err}_U(\mathbf{X}^{m,n}; \hat{\mathbf{u}}) > 0) \leq \log(mn)/mn$, $\mathbb{P}(\text{Err}_V(\mathbf{X}^{m,n}; \hat{\mathbf{v}}) > 0) \leq \log(mn)/mn$.

In what follows we will assume without loss of generality $m \leq n$, and recall that $|\mathcal{L}| = k$, identifying $\mathcal{L} = \{1, \dots, k\}$. We will assume k fixed. We will use C, c, c', \dots for constants that might depend on k , as well as the constant c_0 in the statement in ways that we do not track.

We will show that these bounds hold conditional on \mathbf{u}, \mathbf{v} , on the events $\min_u \hat{q}_r(u), \min_v \hat{q}_r(v) \geq c/2$ which holds with probability at least $1 - \exp(-c'm) \geq 1 - \log(mn)/mn$. Hence, hereafter we will treat \mathbf{u}, \mathbf{v} as deterministic. Recall that $\mathbf{M} \in \mathbb{R}^{m \times n}$ is the matrix entries with

$$M_{ij} = \psi(X_{i,j}),$$

and let $\mathbf{M}_* = \mathbb{E}\{\mathbf{M}\}$. We collect a few facts about \mathbf{M} and its expectation.

Singular values. Note that \mathbf{M}_* takes the form

$$\mathbf{M}_* = \mathbf{L}\Psi\mathbf{R}^\top, \tag{G.23}$$

where $\Psi \in \mathbb{R}^{r \times r}$ is a matrix with entries $\Psi_{u,v} = \bar{\psi}(u, v)$, $\mathbf{L} \in \{0, 1\}^{m \times r}$, with $L_{ij} = 1 \Leftrightarrow u_i = j$, and $\mathbf{R} \in \{0, 1\}^{n \times r}$, with $R_{ij} = 1 \Leftrightarrow v_i = j$. Define $\mathbf{L} = \mathbf{L}_0 \mathbf{D}_L^{1/2}$ where \mathbf{D}_L is a diagonal matrix with $(\mathbf{D}_L)_{ii} = m\hat{q}_r(i)$, and analogously $\mathbf{R} = \mathbf{R}_0 \mathbf{D}_R^{1/2}$, and introduce the singular value decomposition $\mathbf{D}_L^{1/2} \Psi \mathbf{D}_R^{1/2} = \bar{\mathbf{A}} \Sigma \bar{\mathbf{B}}^\top$. We then have the singular value decomposition

$$\mathbf{M}_* = \mathbf{A}_* \Sigma \mathbf{B}_*^\top, \quad \mathbf{A}_* = \mathbf{L}_0 \bar{\mathbf{A}}, \quad \mathbf{B}_* = \mathbf{R}_0 \bar{\mathbf{B}}. \tag{G.24}$$

Therefore $\sigma_k(\mathbf{M}_*) \geq \sigma_{\min}(\mathbf{D}_L)^{1/2} \sigma_{\min}(\Psi) \sigma_{\min}(\mathbf{D}_R)^{1/2}$ (here and below σ_k denotes the k -th largest singular value) and using the assumptions on \hat{q}_r, \hat{q}_c ,

$$\sigma_k(\mathbf{M}_*) \geq c\mu\sqrt{mn}. \tag{G.25}$$

Concentration. $M - M_*$ is a centered matrix with independent entries with variance bounded by σ^2 and entries bounded by 1 (by the assumption $|\psi(x)| \leq 1$). By matrix Bernstein inequality there exists a universal constant C such that the following holds with probability at least $1 - (100n)^{-2}$:

$$\|M - M_*\|_{\text{op}} \leq C \max\left(\sigma\sqrt{n \log n}; \log n\right). \quad (\text{G.26})$$

Incoherence. Since all the entries of M_* are bounded by 1, we get

$$\|M_*\|_{2 \rightarrow \infty} \vee \|M_*^\top\|_{2 \rightarrow \infty} \leq \nu\sqrt{n}. \quad (\text{G.27})$$

Row concentration. For any $i \leq m$, and any $\mathbf{W} \in \mathbb{R}^{n \times k}$ fixed, with probability at least $1 - (100n)^{-5}$:

$$\|(M - bM_*)_{i, \cdot} \mathbf{W}\|_2 \leq C \max\left(\sigma\|\mathbf{W}\|_F \sqrt{\log n}; \|\mathbf{W}\|_{2 \rightarrow \infty} \log n\right).$$

Defining $\Delta_* := \sigma_k(M_*) \geq c\mu\sqrt{mn}$, this implies

$$\begin{aligned} \|(M - M_*)_{i, \cdot} \mathbf{W}\|_2 &\leq \Delta_* \|\mathbf{W}\|_{2 \rightarrow \infty} \varphi\left(\frac{\|\mathbf{W}\|_F}{\sqrt{n}\|\mathbf{W}\|_{2 \rightarrow \infty}}\right), \\ \varphi(x) &:= \frac{C}{\mu\sqrt{mn}} \max\left(x\sigma\sqrt{n \log n}; \log n\right). \end{aligned}$$

Given these, we apply (Abbe et al., 2020)[Corollary 2.1], with the following estimates of various parameters (see (Abbe et al., 2020) for definitions):

$$\begin{aligned} \Delta_* &\asymp \mu\sqrt{mn}, \\ \gamma &\lesssim \max\left(\frac{\sigma}{\mu} \sqrt{\frac{\log n}{m}}; \frac{\log n}{\mu\sqrt{mn}}\right) \lesssim \frac{\sigma}{\mu} \sqrt{\frac{\log n}{m}}, \\ \|M_*\|_{2 \rightarrow \infty} \vee \|M_*^\top\|_{2 \rightarrow \infty} &\leq \nu\sqrt{n} \lesssim \gamma\Delta_*, \\ \varphi(1) &\lesssim \frac{\sigma}{\mu} \sqrt{\frac{\log n}{m}}, \\ \varphi(\gamma) &\lesssim \max\left(\frac{\sigma^2 \log n}{\mu^2 m}; \frac{\log n}{\mu\sqrt{mn}}\right), \\ \kappa &\asymp 1. \end{aligned}$$

Then (Abbe et al., 2020)[Corollary 2.1] implies that there exists a $k \times k$ orthogonal matrix $\tilde{\mathbf{Q}}$ such that

$$\frac{\|\mathbf{A} - \mathbf{A}_* \tilde{\mathbf{Q}}\|_{2 \rightarrow \infty}}{\|\mathbf{A}_*\|_{2 \rightarrow \infty}} \lesssim (1 + \varphi(1))(\gamma + \varphi(\gamma)) + \varphi(1) \quad (\text{G.28})$$

$$\lesssim \frac{\sigma}{\mu} \sqrt{\frac{\log m}{n}}. \quad (\text{G.29})$$

Recall that $\mathbf{A}_* = \mathbf{L}_0 \bar{\mathbf{A}}$ and $\bar{\mathbf{A}}$ is an orthogonal matrix. Therefore, there exists an orthogonal matrix \mathbf{Q} such that (with the desired probability):

$$\frac{\|\mathbf{A} - \mathbf{L}_0 \mathbf{Q}\|_{2 \rightarrow \infty}}{\|\mathbf{L}_0\|_{2 \rightarrow \infty}} \lesssim \frac{\sigma}{\mu} \sqrt{\frac{\log m}{n}}. \quad (\text{G.30})$$

Further, the i -th row of \mathbf{L}_0 is

$$(\mathbf{L}_0)_{i, \cdot} = \frac{1}{\sqrt{m\hat{q}_T(u_i)}} \mathbf{e}_{u_i}^\top =: z(u_i) \mathbf{e}_{u_i}^\top. \quad (\text{G.31})$$

Hence, for any i , $\sqrt{c_0} \|\mathbf{L}_0\|_{2 \rightarrow \infty} \leq \|(\mathbf{L}_0)_{i,\cdot}\| \leq \|\mathbf{L}_0\|_{2 \rightarrow \infty}$. Denoting by \mathbf{q}_j the j -th row of \mathbf{Q} , we thus get for all i ,

$$\frac{\|\mathbf{a}_i - z(u_i)\mathbf{q}_{u_i}\|_2}{\|\mathbf{a}_i\|_2} \lesssim \frac{\sigma}{\mu} \sqrt{\frac{\log m}{n}}. \quad (\text{G.32})$$

The claim follows immediately using the fact that $\sqrt{c_0} \max_u z(u) \leq \min_u z(u) \leq \max_u z(u)$.