

PROBLEM 1.

- (a) It is easy to check that  $W$  is an i.i.d. process but  $Z$  is not. As  $W$  is i.i.d. it is also stationary. We want to show that  $Z$  is also stationary. To show this, it is sufficient to prove that the distribution of the process does not change by shift in the time domain.

$$\begin{aligned}
 p_Z(Z_m = a_m, Z_{m+1} = a_{m+1}, \dots, Z_{m+r} = a_{m+r}) & \\
 &= \frac{1}{2} p_X(X_m = a_m, X_{m+1} = a_{m+1}, \dots, X_{m+r} = a_{m+r}) \\
 &+ \frac{1}{2} p_Y(Y_m = a_m, Y_{m+1} = a_{m+1}, \dots, Y_{m+r} = a_{m+r}) \\
 &= \frac{1}{2} p_X(X_{m+s} = a_m, X_{m+s+1} = a_{m+1}, \dots, X_{m+s+r} = a_{m+r}) \\
 &+ \frac{1}{2} p_Y(Y_{m+s} = a_m, Y_{m+s+1} = a_{m+1}, \dots, Y_{m+s+r} = a_{m+r}) \\
 &= p_Z(Z_{m+s} = a_m, Z_{m+s+1} = a_{m+1}, \dots, Z_{m+s+r} = a_{m+r}),
 \end{aligned}$$

where we used the stationarity of the  $X$  and  $Y$  processes. This shows the invariance of the distribution with respect to the arbitrary shift  $s$  in time which implies stationarity.

- (b) For the  $Z$  process we have

$$\begin{aligned}
 H(Z) &= \lim_{n \rightarrow \infty} \frac{1}{n} H(Z_1, \dots, Z_n) \\
 &= \lim_{n \rightarrow \infty} \frac{1}{n} [H(Z_1, \dots, Z_n | \Theta) + (H(\Theta) - H(\Theta | Z_1, \dots, Z_n))] \\
 &\stackrel{(a)}{=} \lim_{n \rightarrow \infty} \frac{1}{n} H(Z_1, \dots, Z_n | \Theta) \\
 &= \frac{1}{2} H(X_0) + \frac{1}{2} H(Y_0) = 1.
 \end{aligned}$$

where (a) is because  $H(\Theta) - H(\Theta | Z_1, \dots, Z_n)$  is always less than or equal to 1 and it diminishes as it is multiplied by  $\frac{1}{n}$  when  $n$  goes to infinity.

$W$  process is an i.i.d process with the distribution  $p_W(a) = \frac{1}{2} p_X(a) + \frac{1}{2} p_Y(a)$ . From concavity of the entropy, it is easy to see that  $H(W) = H(W_0) \geq \frac{1}{2} H(X_0) + \frac{1}{2} H(Y_0) = 1$ . Hence, the entropy rate of  $W$  is greater than the entropy rate of  $Z$  and the equality holds if and only if  $X_0$  and  $Y_0$  have the same probability distribution function.

PROBLEM 2.

(a) Let  $p_i = \frac{n_i}{n}$ . Then

$$\begin{aligned}
1 &= (p_1 + p_2 + \dots + p_K)^n \stackrel{(a)}{=} \sum_{\substack{n_1, n_2, \dots, n_K \\ \text{s.t. } \sum n_i = n}} \binom{n}{n_1 n_2 \dots n_K} p_1^{n_1} p_2^{n_2} \dots p_K^{n_K} \\
&\geq \binom{n}{n_1 n_2 \dots n_K} p_1^{n_1} p_2^{n_2} \dots p_K^{n_K} \\
&\geq \binom{n}{n_1 n_2 \dots n_K} 2^{n_1 \log(p_1) + n_2 \log(p_2) + \dots + n_K \log(p_K)} \\
&\geq \binom{n}{n_1 n_2 \dots n_K} 2^{-nh(p_1, p_2, \dots, p_K)},
\end{aligned}$$

where (a) is the binomial expansion. This proves our claim.

(b) For a random sequence of length  $n$ ,  $U_1 U_2 \dots U_n$ , we encode the number of occurrences of the first  $(K - 1)$  letters, denoted  $N_1, N_2, \dots, N_{K-1}$ , since we get the last letter for free ( $N_K = n - N_1 - N_2 - \dots - N_{K-1}$ ). For each letter we need at most  $\lceil \log(n + 1) \rceil$  bits. Now that we know the number of times each letter appeared in the sequence we need to encode the index of this specific sequence among all sequences having the same numbers of letter occurrences  $(N_1, \dots, N_{K-1})$ . Since there are  $\binom{n}{N_1 N_2 \dots N_K}$  of those sequences then we need at most  $\lceil \log \left( \binom{n}{N_1 N_2 \dots N_K} \right) \rceil$  bits.

Hence the total length  $L$  of the codeword is

$$L = (K - 1) \lceil \log(n + 1) \rceil + \left\lceil \log \left( \binom{n}{N_1 N_2 \dots N_K} \right) \right\rceil.$$

The expected length is

$$\begin{aligned}
\mathbb{E}(L) &= (K - 1) \lceil \log(n + 1) \rceil + \mathbb{E} \left( \left\lceil \log \left( \binom{n}{N_1 N_2 \dots N_K} \right) \right\rceil \right) \\
&\leq (K - 1) (\log(n + 1) + 1) + \mathbb{E} \left( \log \left( \binom{n}{N_1 N_2 \dots N_K} \right) \right) + 1 \\
&\stackrel{(a)}{\leq} (K - 1) \log(n + 1) + K + \mathbb{E} \left( nh \left( \frac{N_1}{n}, \frac{N_2}{n}, \dots, \frac{N_K}{n} \right) \right) \\
&\stackrel{(b)}{\leq} (K - 1) \log(n + 1) + K + nh \left( \mathbb{E} \left( \frac{N_1}{n}, \frac{N_2}{n}, \dots, \frac{N_K}{n} \right) \right)
\end{aligned}$$

where (a) is due to the first part of the exercise and (b) is due to Jensen's inequality and concavity of the entropy function.

For a random sequence  $U_1 U_2 \dots U_n$ , the number of occurrences of a particular letter  $u_i$  is a random variables that can be written as the sum of indicator functions which take the value 1 with probability  $q_i$

$$N_i = \sum_{j=1}^n 1_{\{U_j = u_i\}}.$$

Hence

$$\mathbb{E} \left( \frac{N_i}{n} \right) = \frac{\sum_{j=1}^n \mathbb{E} (1_{\{U_j = u_i\}})}{n} = q_i.$$

Therefore, the expected codeword length per letter is

$$\frac{1}{n}\mathbb{E}(L) \leq (K-1)\frac{\log(n+1)}{n} + \frac{K}{n} + h(q_1, q_2, \dots, q_K).$$

This shows that  $\lim_{n \rightarrow \infty} \frac{1}{n}\mathbb{E}(L) \leq h(q_1, q_2, \dots, q_K) = H(U)$ . Since the source is i.i.d then  $H(U) = \lim_{n \rightarrow \infty} \frac{1}{n}H(U_1 U_2 \dots U_n)$ . This proves the optimality of this compression code for i.i.d sources.

- (c) If the source is not i.i.d then  $H(U) \neq \lim_{n \rightarrow \infty} \frac{1}{n}H(U_1 U_2 \dots U_n)$ . Hence, the code is not necessarily optimal.

PROBLEM 3.

- (a) We have for all  $k, n$  and  $\epsilon$ ,  $P((U_1, \dots, U_n) \in T(n, p_k, \epsilon)) \leq P((U_1, \dots, U_n) \in T(n, \epsilon))$  as  $T(n, \epsilon) \supseteq T(n, p_k, \epsilon)$ . This implies that for any  $\epsilon > 0$ , with  $k$  and  $p$  such that  $p_k = p$ , we have

$$\begin{aligned} \lim_{n \rightarrow \infty} \Pr((U_1, \dots, U_n) \in T(n, p_k, \epsilon)) &\leq \lim_{n \rightarrow \infty} \Pr((U_1, \dots, U_n) \in T(n, \epsilon)) \\ &1 \leq \lim_{n \rightarrow \infty} \Pr((U_1, \dots, U_n) \in T(n, \epsilon)). \end{aligned}$$

where the second line is due to the property of typical sets.

As we also have  $\lim_{n \rightarrow \infty} \Pr((U_1, \dots, U_n) \in T(n, \epsilon)) \leq 1$ , with these inequalities we prove the statement.

- (b) For typical sets, we know that  $|T(n, p_k, \epsilon)| \leq 2^{(1+\epsilon)Hkn} \leq 2^{(1+\epsilon)Hn}$ . Hence, we obtain the following upper bound.

$$|T(n, \epsilon)| = \left| \bigcup_k T(n, p_k, \epsilon) \right| \leq \sum_k |T(n, p_k, \epsilon)| \leq K 2^{(1+\epsilon)Hn}.$$

By taking logarithm and dividing by  $n$  the above expression, we have

$$\frac{1}{n} \log |T(n, \epsilon)| \leq (1 + \epsilon)H + \frac{\log K}{n}.$$

This implies that for any  $n \geq \log K/\epsilon$  we have

$$\frac{1}{n} \log |T(n, \epsilon)| \leq (1 + \epsilon)H + \epsilon.$$

- (c) Let us use the construction of prefix-free code for typical set given in the lectures. First, take an injective function  $f_{\epsilon, n} : T(n, \epsilon) \rightarrow \{0, 1\}^{\lceil n(1+\epsilon)H + n\epsilon \rceil}$ , this function exists for large enough  $n$  due to our result in (b). Now take another injective function  $g_n : \mathcal{U}^n \rightarrow \{0, 1\}^{\lceil n \log |\mathcal{U}| \rceil}$ . We define  $c_{\epsilon, n}(x)$  as  $0 \parallel f_{\epsilon, n}(x)$  if  $x \in T(n, \epsilon)$  and  $1 \parallel g_n$  otherwise, where  $\parallel$  is the concatenation operator.

We have that

$$\begin{aligned} \Pr(U^n \in T(n, \epsilon)) &= \Pr(\text{length}(c_{\epsilon, n}(U^n)) = \lceil n(1 + \epsilon)H + n\epsilon \rceil + 1) \\ &\leq \Pr(\text{length}(c_{\epsilon, n}(U^n)) \leq n(1 + \epsilon)H + n\epsilon + 2). \end{aligned}$$

From (a) we know that there exists an  $n_a(\epsilon, \delta)$  such that  $1 - \delta < \Pr(U^n \in T(n, \epsilon))$  for all  $n \geq n_a(\epsilon, \delta)$ . From (b) we require  $n \geq \log K/\epsilon = n_b(K, \epsilon)$ . To get the form required in the problem statement, we need that :

$$n((1 + \epsilon)H + \epsilon + 2/n) < nR$$

Since  $1/n \leq \epsilon$  for  $n \geq n_b(K, \epsilon)$ , the following inequality will also work

$$n((1 + \epsilon)H + 3\epsilon) < nR.$$

The above inequality satisfied by choosing an appropriate  $\epsilon$  (i.e.,  $0 \leq \epsilon < \frac{R-H}{H+3}$ ).

Therefore, for a code  $c_{\epsilon,n}$  constructed as above and  $\epsilon$  chosen small enough, we have

$$\begin{aligned} \Pr(\text{length}(c_{\epsilon^*,n}(U^n)) < nR) &\geq \Pr(\text{length}(c_{\epsilon,n}(U^n)) \leq n(1 + \epsilon)H + n\epsilon + 1) \\ &\geq \Pr(U^n \in T(n, \epsilon)) \\ &> 1 - \delta \end{aligned}$$

for all  $n \geq \max\{n_a(\epsilon, \delta), n_b(K, \epsilon)\}$ .

**PROBLEM 4.**

- (a) We have  $\rho(X_1^\infty) = 0$ . We show this by showing that  $\rho(X_1^\infty) \leq \delta$  for any  $\delta > 0$ . To see the last statement, build an invertible FSM which “recognizes” a string of type “ab...ab” for a particular even length, call it  $L$ , and outputs lets say “0” at the end of this string and returns to the starting state. Hence this machine will output an infinite string of “0” when the input is  $X_1^\infty$ . From each state (including the starting state) of the chain which recognizes the special string make an edge back to the starting state in the case the next input is not the correct one. The output for each such edge is  $1 + \lceil \log L \rceil$  bits long, the first bit is 1 to indicate that it is not the special path and on the next  $\lceil \log L \rceil$  bits we give the index of the state (in binary representation) from which the return edge is drawn. This machine is clearly lossless and has a compressibility of  $1/L$  for the desired sequence.
- (b) A machine as described above will have  $\rho_M(X_1^\infty) = 1/4$ . In fact, one cannot do better than this. Consider a cycle, when from a given state we get back to the same state. During such a cycle we have to output at least one symbol, because the machine has to be information lossless. In an  $L$  state machine we eventually create such a cycle within at most  $L$  steps. This means that we output at least one symbol for every  $L$  input symbols, so  $\rho_M(X_1^\infty) \geq 1/L$ .
- (c) We have  $\rho_{LZ} = 0$  since compressibility is non-negative and we know that the compressibility of LZ is at least as good as that of any FSM, i.e., we know that  $\rho_{LZ}(X_1^\infty) \leq \rho(X_1^\infty)$ .
- (d) The dictionary increases by 1 every time and has size 2 in the beginning. Hence, if we look at lets say  $c$  steps of the algorithm then we need in total

$$\sum_{i=1}^c \lceil \log(1 + i) \rceil \leq c \log(2(c + 1))$$

bits to describe the output.

What are the words which we are using. Note that the parsing is  $a, b, ab, aba, ba, bab, \dots$ . Note that in average at most every second step the length of the used dictionary word increases by 1, i.e., we have a linear increase in the used dictionary words. Therefore, if we compute the total length which we have parsed after  $c$  steps, this length increases like the square of  $c$ .

It follows that the ratio of the total number of bits used divided by the total length described behaves like  $1/c$ , i.e., it tends to 0.