

CS258 Information Theory Homework 4

Zhou Litao 518030910407 F1803016

March 24, 2020

Exercise 1 (Monotonicity of entropy per element) For a stationary stochastic process X_1, X_2, \dots, X_n , show that

$$\begin{aligned} \frac{H(X_1, X_2, \dots, X_n)}{n} &\leq \frac{H(X_1, X_2, \dots, X_{n-1})}{n-1} \\ \frac{H(X_1, X_2, \dots, X_n)}{n} &\geq H(X_n | X_{n-1}, \dots, X_1) \end{aligned}$$

Proof.

We first prove the second statement by properties of stationary process.

$$\begin{aligned} H(X_1, \dots, X_n) &= \sum_{i=1}^n H(X_i | X_{i-1}, \dots, X_2, X_1) \\ &\geq \sum_{i=1}^n H(X_i | X_{i-1}, X_{i-2}, \dots, X_{n-i+1}) \\ &= \sum_{i=1}^n H(X_n | X_{n-1}, \dots, X_2, X_1) \\ &= nH(X_n | X_1, X_2, \dots, X_{n-1}) \end{aligned} \tag{1}$$

By Equation 1 and entropy equalities we can prove the first statement.

$$\begin{aligned} (n-1)H(X_1, \dots, X_n) &\leq nH(X_1, \dots, X_n) - H(X_1, \dots, X_n) \\ &\leq n[H(X_1, \dots, X_n) - H(X_n | X_1, \dots, X_{n-1})] \\ &= nH(X_1, \dots, X_{n-1}) \end{aligned} \tag{2}$$

□

Exercise 2 (Initial conditions) Show, for a Markov chain, that

$$H(X_0 | X_n) \geq H(X_0 | X_{n-1})$$

Thus, initial conditions X_0 become more difficult to recover as the future X_n unfolds.

Proof. For Markov Chain, we have $p(x_0 | x_n, x_{n-1}) = p(x_0 | x_{n-1})$. Hence

$$\begin{aligned} H(X_0 | X_n, X_{n-1}) &= - \sum_{x_n, x_{n-1}} p(x_n, x_{n-1}) \sum_{x_0} p(x_0 | x_n, x_{n-1}) \log p(x_0 | x_n, x_{n-1}) \\ &= - \sum_{x_{n-1}} \left(\sum_{x_n} p(x_n, x_{n-1}) \right) \sum_{x_0} p(x_0 | x_{n-1}) \log p(x_0 | x_{n-1}) \\ &= - \sum_{x_{n-1}} p(x_{n-1}) \sum_{x_0} p(x_0 | x_{n-1}) \log p(x_0 | x_{n-1}) = H(X_0 | X_{n-1}) \end{aligned} \tag{3}$$

Since condition reduce entropy, we have $H(X_0 | X_{n-1}) = H(X_0 | X_n, X_{n-1}) \leq H(X_0 | X_n)$

□

Exercise 3 (The past has little to say about the future) For a stationary stochastic process $X_1, X_2, \dots, X_n, \dots$, show that

$$\lim_{n \rightarrow \infty} \frac{1}{2n} I(X_1, X_2, \dots, X_n; X_{n+1}, X_{n+2}, \dots, X_{2n}) = 0$$

Proof. First Note that

$$I(X_1, X_2, \dots, X_n; X_{n+1}, X_{n+2}, \dots, X_{2n}) = H(X_1, \dots, X_n) - H(X_{n+1}, \dots, X_{2n} | X_1, \dots, X_n) \quad (4)$$

By definition of entropy rate we know that

$$\frac{1}{2} \sum_{i=1}^n H(X_1, \dots, X_n) \rightarrow H(\mathcal{X})$$

□

Exercise 4 (Entropy rate) Let $\{X_i\}$ be a discrete stationary stochastic process with entropy rate $H(\mathcal{X})$. Show that

$$\frac{1}{n} H(X_n, \dots, X_1 | X_0, X_{-1}, \dots, X_{-k}) \rightarrow H(\mathcal{X})$$

for $k = 1, 2, \dots$

Proof.

$$H(X_n, \dots, X_1 | X_0, X_{-1}, \dots, X_{-k}) = \sum_{i=1}^n H(X_i | X_{i-1}, \dots, X_{-k}) \quad (5)$$

Note that

$$H'(X) = \lim_{n \rightarrow \infty} H(X_n | X_{n-1}, \dots, X_1) = \lim_{n \rightarrow \infty} H(X_n | X_{n-1}, \dots, X_{-k})$$

We already know that $H'(X) \rightarrow H(\mathcal{X})$.

By Cesaro Mean,

$$\frac{1}{n} H(X_n, \dots, X_1 | X_0, X_{-1}, \dots, X_{-k}) \rightarrow H(\mathcal{X})$$

□

Exercise 5 (Markov's inequality and Chebyshev's inequality)

1. (Markov's inequality) For any nonnegative random variable X and any $t > 0$, show that

$$\Pr\{X \geq t\} \leq \frac{EX}{t} \quad (6)$$

Exhibit a random variable that achieves this inequality with equality.

2. (Chebyshev's inequality) Let Y be a random variable with mean μ and variance σ^2 . By letting $X = (Y - \mu)^2$, show that for any $\epsilon > 0$

$$\Pr\{|Y - \mu| > \epsilon\} \leq \frac{\sigma^2}{\epsilon^2} \quad (7)$$

3. (Weak law of large numbers) Let Z_1, Z_2, \dots, Z_n be a sequence of i.i.d. random variables with mean μ and variance σ^2 . Let $\bar{Z}_n = \frac{1}{n} \sum_{i=1}^n Z_i$ be the sample mean. Show that

$$\Pr\{|\bar{Z}_n - \mu| > \epsilon\} \leq \frac{\sigma^2}{n\epsilon^2} \quad (8)$$

Thus, $\Pr\{|\bar{Z}_n - \mu| > \epsilon\} \rightarrow 0$ as $n \rightarrow \infty$. This is known as the weak law of large numbers.

Proof.

1. Note that for any t ,

$$t \cdot 1_{\{X \geq t\}} \leq X \quad (9)$$

By taking expectation at both sides we have

$$E(1_{\{X \geq t\}}) = \Pr(X \geq t) \leq \frac{EX}{t} \quad (10)$$

2. Note $EX = DY = \sigma^2$. By letting $X = (Y - \mu)^2$ and $t = \epsilon^2$ in the Markov inequality, we can derive that

$$\Pr\{|\bar{Z}_n - \mu| > \epsilon\} \leq \frac{\sigma^2}{n\epsilon^2} \quad (11)$$

3. Note $E\bar{Z}_n = \mu$ and $D\bar{Z}_n = \frac{1}{n^2} \sum_{i=1}^n DZ_i = \frac{\sigma^2}{n}$. By applying Chebyshev's inequality on \bar{Z}_n we have

$$\Pr\{|\bar{Z}_n - \mu| > \epsilon\} \leq \frac{\sigma^2}{n\epsilon^2} \quad (12)$$

□

Exercise 6 (Piece of cake) A cake is sliced roughly in half, the largest piece being chosen each time, the other pieces discarded. We will assume that a random cut creates pieces of proportions

$$P = \begin{cases} \left(\frac{2}{3}, \frac{1}{3}\right) & \text{with probability } \frac{3}{4} \\ \left(\frac{3}{5}, \frac{2}{5}\right) & \text{with probability } \frac{1}{4} \end{cases}$$

Thus, for example, the first cut (and choice of largest piece) may result in a piece of size $\frac{3}{5}$. Cutting and choosing from this piece might reduce it to size $\left(\frac{3}{5}\right)\left(\frac{2}{3}\right)$ at time 2, and so on. How large, to first order in the exponent, is the piece of cake after n cuts?

Solution. Let C_1, C_2, \dots, C_n denote the choice of each cut. Then after n cuts, the size of the cake $W_n = \prod_{i=1}^n C_i$. By taking the logarithm at the equation, we have $\log W_n = \sum_{i=1}^n \log C_i$. Since C_i s are i.i.d., we can apply the law of large numbers as follows.

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log W_n = E(\log C) = \frac{3}{4} \log \frac{2}{3} + \frac{1}{4} \log \frac{3}{5} \approx -0.494 \quad (13)$$

Note that the equation above indicates that

$$\frac{1}{n} \log W_n = -0.494 + o(1) \quad W_n = 2^{-0.494n + o(n)} \quad (14)$$

So the first order in the exponent is $\frac{3}{4} \log \frac{2}{3} + \frac{1}{4} \log \frac{3}{5} \approx -0.494$. □

Exercise 7 (AEP) Let X_i be iid $\sim p(x), x \in \{1, 2, \dots, m\}$. Let $\mu = EX$ and $H = -\sum p(x) \log p(x)$. Let $A^n = \{x^n \in \mathcal{X}^n : |-\frac{1}{n} \log p(x^n) - H| \leq \epsilon\}$. Let $B^n = \{x^n \in \mathcal{X}^n : |\frac{1}{n} \sum_{i=1}^n X_i - \mu| \leq \epsilon\}$

1. Does $\Pr\{X^n \in A^n\} \rightarrow 1$?
2. Does $\Pr\{X^n \in A^n \cap B^n\} \rightarrow 1$?
3. Show that $|A^n \cap B^n| \leq 2^{n(H+\epsilon)}$ for all n
4. Show that $|A^n \cap B^n| \geq \left(\frac{1}{2}\right) 2^{n(H-\epsilon)}$ for n sufficiently large.

Solution.

1. Yes. By the Large Number Law,

$$-\frac{1}{n} \log p(x^n) = -\frac{1}{n} \sum_{i=1}^n \log p(X_i) \rightarrow H(X) \quad \text{in probability} \quad (15)$$

, which implies that

$$\Pr \{X^n \in A^n\} = \Pr \left\{ \left| -\frac{1}{n} \log p(x^n) - H \right| \leq \epsilon \right\} \rightarrow 1 \quad (16)$$

2. Yes. Part (1) implies that $\lim_{n \rightarrow \infty} \Pr \{X^n \in A^n\} = 1$.

By strong law of large number we have $\lim_{n \rightarrow \infty} \Pr \{X^n \in B^n\} = 1$

For arbitrary $\delta > 0$, there exists N_1 , such that $\Pr \{X^n \in A^n\} > 1 - \frac{\delta}{2}$ for all $n > N_1$, and there exists N_2 , such that $\Pr \{X^n \in B^n\} > 1 - \frac{\delta}{2}$ for all $n > N_2$. We take $N = \max \{N_1, N_2\}$, for any $n > N$, we have

$$\begin{aligned} \Pr \{X^n \in A^n \cap B^n\} &= \Pr \{X^n \in A^n\} + \Pr \{X^n \in B^n\} - \Pr \{X^n \in A^n \cup B^n\} \\ &> 1 - \frac{\delta}{2} + 1 - \frac{\delta}{2} - 1 = 1 - \delta \end{aligned} \quad (17)$$

, which indicates that $\Pr \{X^n \in A^n \cap B^n\} \rightarrow 1$

3. From $x^n \in A^n$ we know that

$$2^{-n(H+\epsilon)} \leq p(x^n) \leq 2^{-n(H-\epsilon)} \quad (18)$$

$$\begin{aligned} 1 &= \sum_{x^n \in \mathcal{S}^n} p(x) \\ &\geq \sum_{x^n \in A^n \cap B^n} p(x) \\ &\geq 2^{-n(H(X)+\epsilon)} |A^n \cap B^n| \end{aligned} \quad (19)$$

It follows that $|A^n \cap B^n| \leq 2^{-n(H+\epsilon)}$.

4. From $\Pr \{X^n \in A^n \cap B^n\} \rightarrow 1$, we take $\delta = \frac{1}{2}$, then there exists sufficiently large n , such that

$$\begin{aligned} \frac{1}{2} &\leq \Pr(X^n \in A^n \cap B^n) \\ &\leq \sum_{x^n \in A^n \cap B^n} p(x^n) \\ &\leq |A^n \cap B^n| 2^{-n(H(X)-\epsilon)} \end{aligned} \quad (20)$$

It follows that $|A^n \cap B^n| \geq \left(\frac{1}{2}\right) 2^{-n(H-\epsilon)}$.

□

Exercise 8 (Doubly stochastic matrices) An $n \times n$ matrix $P = [P_{ij}]$ is said to be doubly stochastic if $P_{ij} \geq 0$ and $\sum_j P_{ij} = 1$ for all i and $\sum_i P_{ij} = 1$ for all j . An $n \times n$ matrix P is said to be a permutation matrix if it is doubly stochastic and there is precisely one $P_{ij} = 1$ in each row and each column. It can be shown that every doubly stochastic matrix can be written as the convex combination of permutation matrices.

1. Let $\mathbf{a}^t = (a_1, a_2, \dots, a_n)$, $a_i \geq 0$, $\sum a_i = 1$, be a probability vector. Let $\mathbf{b} = \mathbf{a}P$, where P is doubly stochastic. Show that \mathbf{b} is a probability vector and that $H(b_1, b_2, \dots, b_n) \geq H(a_1, a_2, \dots, a_n)$. Thus, stochastic mixing increases entropy.
2. Show that a stationary distribution μ for a doubly stochastic matrix P is the uniform distribution.
3. Conversely, prove that if the uniform distribution is a stationary distribution for a Markov transition matrix P , then P is doubly stochastic.

Proof.

1. From $\mathbf{b} = \mathbf{a}P$ we know that

$$b_j = \sum_{i=1}^m a_i p_{ij} \quad (21)$$

Then we have.

$$\begin{aligned} H(\mathbf{b}) - H(\mathbf{a}) &= - \sum_{j=1}^m \left(\sum_{i=1}^m a_i p_{ij} \right) \log b_j + \sum_{i=1}^m a_i \log a_i \\ &= - \sum_{i=1}^m a_i \left(\sum_{j=1}^m p_{ij} \log b_j \right) + \sum_{i=1}^m a_i \left(\sum_{j=1}^m p_{ij} \log a_i \right) \\ &= \sum_{i=1}^m \sum_{j=1}^m a_i p_{ij} \log \frac{a_i}{b_j} \\ &= \sum_{i=1}^m \sum_{j=1}^m b_j p_{ij} \left(\frac{a_i}{b_j} \log \frac{a_i}{b_j} \right) \\ &\geq \left(\sum_{i=1}^m \sum_{j=1}^m b_j p_{ij} \frac{a_i}{b_j} \right) \log \left(\sum_{i=1}^m \sum_{j=1}^m b_j p_{ij} \frac{a_i}{b_j} \right) \\ &= \left(\sum_{i=1}^m \sum_{j=1}^m p_{ij} a_i \right) \log \left(\sum_{i=1}^m \sum_{j=1}^m p_{ij} a_i \right) = 1 \cdot \log 1 = 0 \end{aligned} \quad (22)$$

2. By condition we have $\mu_i = \frac{1}{m}$ for any i . Since for any j ,

$$\sum_{i=1}^m \mu_i p_{ij} = \frac{1}{m} \sum_{i=1}^m p_{ij} = \frac{1}{m} = \mu_i \quad (23)$$

We have that $\mu P = \mu$. The uniform distribution is a stationary distribution for a doubly stochastic matrix.

3. From $\mu P = \mu$ and $\mu = \frac{1}{m}$, we know that

$$\sum_{i=1}^m \frac{1}{m} p_{ij} = \frac{1}{m} \quad (24)$$

holds for any j , which implies

$$\sum_{i=1}^m p_{ij} = 1 \quad \text{for any } j \quad (25)$$

Then P is doubly stochastic.

□

Exercise 9 (Shuffles increase entropy) Argue that for any distribution on shuffles T and any distribution on card positions X that

$$\begin{aligned} H(TX) &\geq H(TX|T) \\ &= H(T^{-1}TX|T) \\ &= H(X|T) \\ &= H(X) \end{aligned}$$

if X and T are independent.

Proof. The first line holds because condition reduces entropy. The second line holds since T^{-1} can be given by the condition T . The last line holds if X and T are independent, which finishes the proof. \square

Exercise 10 (Entropy rates of Markov chains)

1. Find the entropy rate of the two-state Markov chain with transition matrix

$$P = \begin{bmatrix} 1 - p_{01} & p_{01} \\ p_{10} & 1 - p_{10} \end{bmatrix}$$

2. What values of p_{01}, p_{10} maximize the entropy rate?
3. Find the entropy rate of the two-state Markov chain with transition matrix

$$P = \begin{bmatrix} 1 - p & p \\ 1 & 0 \end{bmatrix}$$

4. Find the maximum value of the entropy rate of the Markov chain of part (c). We expect that the maximizing value of p should be less than $\frac{1}{2}$, since the 0 state permits more information to be generated than the 1 state.
5. Let $N(t)$ be the number of allowable state sequences of length t for the Markov chain of part (c). Find $N(t)$ and calculate

$$H_0 = \lim_{t \rightarrow \infty} \frac{1}{t} \log N(t)$$

[Hint: Find a linear recurrence that expresses $N(t)$ in terms of $N(t-1)$ and $N(t-2)$. Why is H_0 an upper bound on the entropy rate of the Markov chain? Compare H_0 with the maximum entropy found in part (d).]

Solution.

1. We first calculate the stationary distribution μ .

$$\begin{cases} \mu P = \mu \\ \mu \mathbf{1}^T = 1 \end{cases} \Rightarrow \mu = \left[\frac{p_{10}}{p_{01} + p_{10}}, \frac{p_{01}}{p_{01} + p_{10}} \right] \quad (26)$$

By the entropy rate of Markov Chain,

$$\begin{aligned} H(\mathcal{X}) &= \sum_{i=1,2} \mu_i \left(\sum_{j=1,2} -p_{ij} \log p_{ij} \right) \\ &= \frac{p_{10}H(p_{01}) + p_{01}H(p_{10})}{p_{01} + p_{10}} \end{aligned} \quad (27)$$

2. Note that the entropy is upper bounded by its alphabet, since

$$\frac{1}{n} H(X_1, X_2, \dots, X_n) \leq \frac{1}{n} \log |\mathcal{X}|^n = \log |\mathcal{X}| \quad (28)$$

Hence the maximal entropy rate for this problem is $\log 2 = 1$. This can be obtained when $p_{01} = p_{10} = \frac{1}{2}$, where

$$\begin{aligned} H(\mathcal{X}) &= \frac{p_{10}H(p_{01}) + p_{01}H(p_{10})}{p_{01} + p_{10}} \\ &= \frac{1}{2} + \frac{1}{2} = 1 \end{aligned} \quad (29)$$

Therefore, $p_{01} = p_{10} = \frac{1}{2}$ maximize the entropy rate.

3. We first calculate the stationary distribution μ .

$$\begin{cases} \mu P = \mu \\ \mu \mathbf{1}^T = 1 \end{cases} \Rightarrow \mu = \left[\frac{1}{p+1}, \frac{p}{p+1} \right] \quad (30)$$

By the entropy rate of Markov Chain,

$$\begin{aligned} H(\mathcal{X}) &= \sum_{i=1,2} \mu_i \left(\sum_{j=1,2} -p_{ij} \log p_{ij} \right) \\ &= \frac{-p \log p - (1-p) \log(1-p)}{p+1} \end{aligned} \quad (31)$$

4. We take the derivative of $H(\mathcal{X})$.

$$\frac{dH(\mathcal{X})}{dp} = \frac{\log(1-p) - \log p + \log(1-p)}{(p+1)^2} := 0 \Rightarrow p = \frac{3-\sqrt{5}}{2} \quad (32)$$

The maximal entropy rate is $\log \frac{1+\sqrt{5}}{2} \approx 0.6942$.

5. The transition matrix implies that there is no possibility that the last state is 1 with the last but one state to be 1. We can calculate the $N(t)$ recursively. If the last state is 0, the previous state possibilities add up to $N(t-1)$. If the last state is 1, however, the last but one state can only be 0. Then all the previous state possibilities will add up to $N(t-2)$. Further more, we can manually check the initial length that $N(1) = 2$, $N(2) = 3$. Now we have

$$N(t) = N(t-1) + N(t-2) \quad (33)$$

By solving the characteristic equation we know that $N(t)$ must be in the form like

$$N(t) = C_1 \left(\frac{1+\sqrt{5}}{2} \right)^t + C_2 \left(\frac{1-\sqrt{5}}{2} \right)^t \quad (34)$$

Then

$$\begin{aligned} H_0 &= \lim_{t \rightarrow \infty} \frac{1}{t} \log N(t) \\ &= \lim_{t \rightarrow \infty} \frac{1}{t} \log \left(C_1 \left(\frac{1+\sqrt{5}}{2} \right)^t + C_2 \left(\frac{1-\sqrt{5}}{2} \right)^t \right) \\ &= \lim_{t \rightarrow \infty} \frac{1}{t} \log C_1 \left(\frac{1+\sqrt{5}}{2} \right)^t \\ &= \lim_{t \rightarrow \infty} \frac{1}{t} \log \frac{1+\sqrt{5}}{2} = \log \frac{1+\sqrt{5}}{2} \end{aligned} \quad (35)$$

We find that H_0 is the upper bound of the entropy rate, and can be obtained with the maximum entropy found in part (d). This is because by the property of entropy,

$$H(X_1, X_2, \dots, X_n) \leq \log |X_1, X_2, \dots, X_n| = N(t) \quad (36)$$

□

Exercise 11 (Maximal entropy graphs) Consider a random walk on a connected graph with four edges.

1. Which graph has the highest entropy rate?
2. Which graph has the lowest?

Solution. In a random walk, the next vertex will be arbitrarily chosen from the adjacent vertices of the current vertex. That is to say, all the edges are given the same weight.

We can formulate this problem as follows. By the inclusion-exclusion principle, there are at most 5 vertices in the graph given four edges. Their adjacent relation can be represented in a 5×5 0-1 matrix W , where $W_{ij} = 0$ or 1 and $W_{ij} = W_{ji}$.

Then the transition probability matrix P and the stationary distribution μ can be calculated.

$$P_{ij} = \frac{W_{ij}}{\sum_k W_{ik}} \quad (37)$$

$$\mu_i = \frac{W_i}{2W} \quad (38)$$

where $W_i = \sum_j w_{ij}$ and $W = \sum_i \frac{W_i}{2}$

Then the entropy rate can be calculated as

$$\begin{aligned} H(\mathcal{X}) &= H(X_2|X_1) \\ &= H(X_2, X_1) - H(X_1) \\ &= H\left(\frac{1}{8}, \frac{1}{8}, \dots, \frac{1}{8}\right) - H(\mu) \\ &= 3 - H(\mu) \end{aligned} \quad (39)$$

For 4 edges, there are five possible graphs, as has been shown in Figure 1.

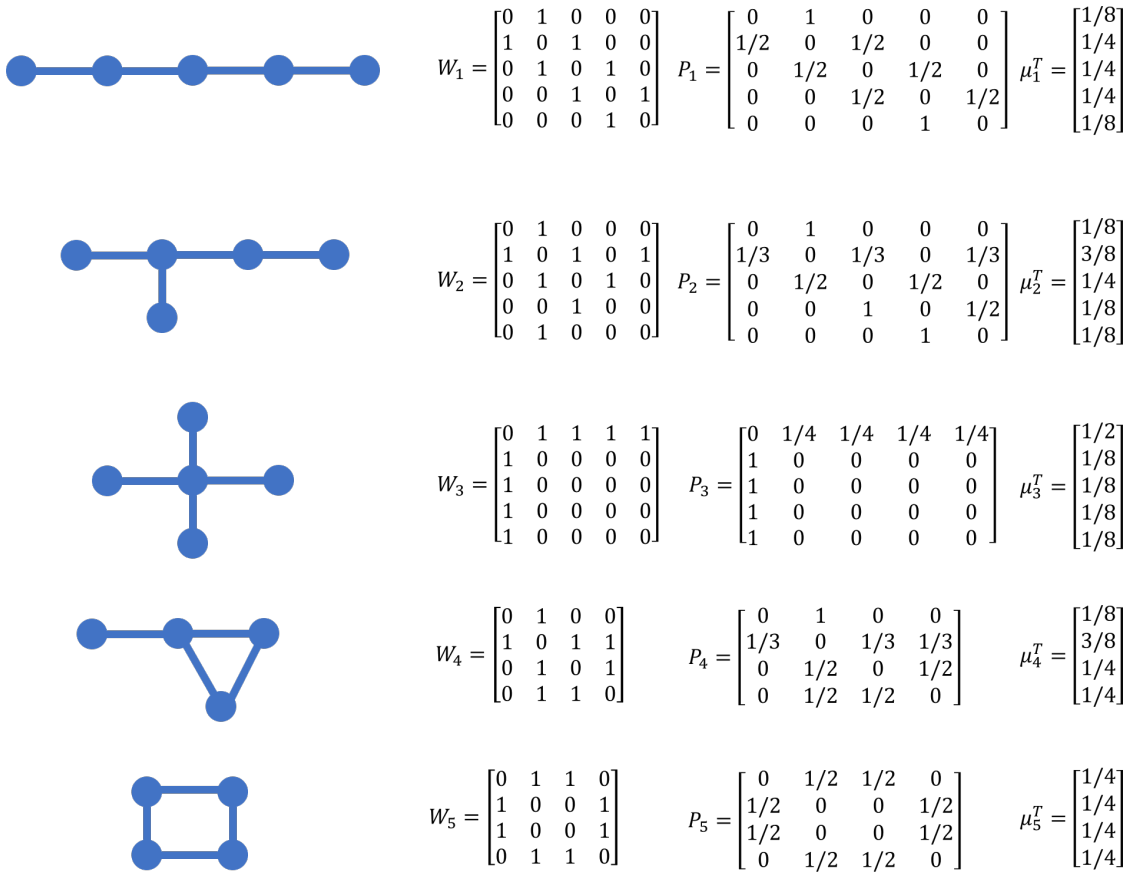


Figure 1: Five Possible Graphs for Four Edges

Their corresponding entropy rates are 0.75, 0.84436, 1, 1.09436 and 1.

1. The fourth graph has the largest entropy rate.
2. The first graph has the lowest entropy rate.

□