

Lecture 9

Lecturer: Haim Permuter

Scribe: Imri Zviely and Itai Lishner

I. MAXIMUM ENTROPY, CHANNEL CODING AND SIDE INFORMATION

In the nature many phenomenons may be explained by the maximum entropy principle. For instance, the velocity of the particles in the air for a given temperature is distributed according to Gaussian distribution. The temperature is proportional to the kinetic energy of the particles. Given a specific temperature, namely, a mean square constraint on the velocity, we obtain that a Gaussian distribution maximizes the entropy. In this lecture we would learn how to calculate the maximum entropy, and more important, understand why this is the case in nature.

Consider the following problem: we would like to maximize the entropy $H(X)$ over $p(x)$ such that the following constraints will hold

- 1) $p(x) \geq 0, \quad \forall x \in \mathcal{X}$.
- 2) $\sum_{x \in \mathcal{X}} p(x) = 1$.
- 3) $\sum_{x \in \mathcal{X}} p(x) r_i(x) = \alpha_i, \quad i = (1, \dots, m)$,

where $r_i(x)$ are functions of x . Note that the constraint $\sum_{x \in \mathcal{X}} p(x) r_i(x) = \alpha_i$ is equivalent to $E[r_i(X)] = \alpha_i$.

Theorem 1 The optimal pmf $p(x)$ for the optimization problem given above is

$$p^*(x) = e^{\lambda_0 - 1 + \sum_{i=1}^m \lambda_i r_i(x)} \quad (1)$$

where $\lambda_i, i = 0, 1, 2, \dots, m$ is the one that satisfies

$$\begin{aligned} \sum_x p^*(x) &= 1 \\ \sum_x p^*(x) r_i(x) &= \alpha_i. \end{aligned} \quad (2)$$

Q: Is this problem a convex optimization problem ?

A: Yes.

We learned that a convex optimization problem is of the form

$$\min_x f_0(x) \quad (3)$$

$$s.t. \quad f_i(x) \leq 0, 1 \leq i \leq m, \quad (4)$$

$$h_j(x) = 0, 1 \leq j \leq l \quad (5)$$

where f_i , $i = 0, 1, 2, \dots, m$ are convex functions and h_j are affine functions.

In our case we have the problem

$$\max_{P_X} H(P_X) \quad (6)$$

$$s.t. \quad P_X(x) \geq 0, \quad \forall x \in \mathcal{X}, \quad (7)$$

$$\sum_{x \in \mathcal{X}} P_X(x) = 1, \quad (8)$$

$$\sum_{x \in \mathcal{X}} P_X(x) r_i(x) = \alpha_i, \quad i = (1, \dots, m). \quad (9)$$

Proof of Theorem 1

We solve the optimization problem using the Dual Lagrange principle. We first dismiss the constraint that $p(x) \geq 0$, however since we will obtain that the optimal solution $p^*(x)$ satisfies this condition, then $p^*(x)$ is optimal also with this constraint.

We form the Lagrange dual functional:

$$J(p(x), \bar{\lambda}) = -\sum p(x) \log(p(x)) + \lambda_0 (\sum p(x) - 1) + \sum_{i=1}^m \lambda_i (\sum_x p(x) r_i(x) - \alpha_i) \quad (10)$$

And differentiate with respect to $p(x)$:

$$\frac{\partial J(p(x), \bar{\lambda})}{\partial p(x)} = -\log(p(x)) - 1 + \lambda_0 + \sum_{i=1}^m \lambda_i r_i(x) = 0 \quad (11)$$

Hence we obtained that

$$p(x) = e^{\lambda_0 - 1 + \sum_{i=1}^m \lambda_i r_i(x)} \quad (12)$$

Now we need to find λ_0 and λ_i , $i = 1, 2, \dots, m$ such that

$$\begin{aligned} \sum_x p(x) &= 1 \\ \sum_x p(x) r_i(x) &= \alpha_i. \end{aligned} \quad (13)$$

For the continuous alphabet where we want to maximize the differential entropy over differential distributions $f(x)$ that satisfies the constraints $E[r_i(X)] = \alpha_i$, $i = 1, 2, \dots, m$ we obtain

$$f(x) = e^{\lambda_0 - 1 + \sum_{i=1}^m \lambda_i r_i(x)} \quad (14)$$

such that

$$\int_x f(x) dx = 1 \quad (15)$$

$$\int_x f(x) r_i(x) dx = \alpha_i. \quad (16)$$

$$(17)$$

Example 1 (dice) Consider a dice with six faces $[1, 2, \dots, 6]$. What would be the probability of the dice that maximize the entropy?

$$\left. \begin{array}{l} p(x) = e^{\lambda_0 - 1} \\ \sum_x p(x) = 1 \end{array} \right\} \Rightarrow p^*(x) = \frac{1}{6} \quad (18)$$

Note that we obtain that $p(x)$ does not depend on x , namely is constant. Therefore, a uniform distribution maximizes the entropy.

Example 2 : $X \in [-\infty, \infty]$ Let the constraints be $E[X] = 0$ and $E[X^2] = \sigma^2$. Then the form of the maximizing distribution is

$$f(x) = e^{\lambda_0 + \lambda_1 x + \lambda_2 x^2} \quad (19)$$

To find the appropriate constants, we first recognize that this distribution has the same form as normal distribution. Hence, the density that satisfies the constraints and also maximizes the entropy is the $\mathcal{N}(0, \sigma^2)$ distribution:

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{x^2}{2\sigma^2}} \quad (20)$$

Alternative proof to Theorem 1 Let $P_X(x) = e^{\lambda_0 - 1 + \sum_{i=1}^m \lambda_i r_i(x)}$, and let $Q_X(x)$ be any distribution that satisfies the conditions $E[r_i(x)] = \alpha_i$.

We need to show that $H(P_X) \geq H(Q_X)$. Consider,

$$D(Q||P) = \sum_x Q(x) \log \frac{Q(x)}{P(x)} \quad (21)$$

$$= \sum_x Q(x) \log(Q(x)) - \sum_x Q(x) \log(P(x)) \quad (22)$$

$$= -H(Q) - \sum_x Q(x) \left[\lambda_0 - 1 + \sum_{i=1}^m \lambda_i r_i(x) \right] \quad (23)$$

$$\stackrel{(a)}{=} -H(Q) - \sum_x P(x) \left[\lambda_0 - 1 + \sum_{i=1}^m \lambda_i r_i(x) \right] \quad (24)$$

$$\stackrel{(b)}{=} -H(Q) + H(P), \quad (25)$$

$$(26)$$

and since $D(P||Q) \geq 0$ we obtained that $H(P) \geq H(Q)$.

Example 3 : Dice x , $\mathcal{X} = [1, 2, \dots, 6]$

Suppose that n dice are thrown on the table and we are told that the sum of the results is $n\alpha$. What is the probability $p(x)$?

It is easy to show that $E[x] = \alpha$.

We define n_1, n_2, \dots, n_6 , when n_i represents the number of throws where $x = i$.

Number of sequences with (n_1, n_2, \dots, n_6) is equal to:

$$\binom{n}{n_1 \ n_2 \ \dots \ n_6} = \frac{n!}{n_1! \ n_2! \ \dots \ n_6!} \quad (27)$$

In order to find the most probable state we wish to maximize $\frac{n!}{n_1! \ n_2! \ \dots \ n_6!}$ under the constraint:

$$\sum_{i=1}^6 \frac{n_i}{n} \cdot i = \alpha \quad ; \quad E[x] = \alpha \quad (28)$$

For large values of n using Stirling's approximation, $n! \approx (\frac{n}{e})^n$, we define that:

$$N = \frac{n!}{n_1! \ n_2! \ \dots \ n_6!} \approx \frac{(\frac{n}{e})^n}{(\frac{n_1}{e})^{n_1} \ \dots \ (\frac{n_6}{e})^{n_6}} \quad (29)$$

$$\log(N) = n \log(n) - n_1 \log(n_1) - \dots - n_6 \log(n_6) \quad (30)$$

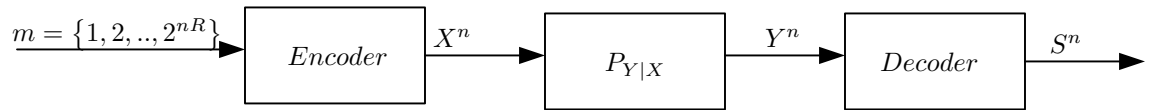
$$n = n_1 + n_2 + \dots + n_6 \quad (31)$$

$$\log(N) = - \sum_i n_i \log \frac{n_i}{n} \quad (32)$$

$$N = 2^{- \sum_i n_i \log \frac{n_i}{n}} = 2^{-n \sum_i \frac{n_i}{n} \log \frac{n_i}{n}} = 2^{nH(\frac{n_1}{n}, \frac{n_2}{n}, \dots, \frac{n_6}{n})} \quad (33)$$

II. CHANNEL CODING WITH SIDE INFORMATION

Channel coding:



$$C = \max_{p(x)} I(X; Y) \quad (34)$$

Fig. 1. Communication system

Channel coding with Side Information:

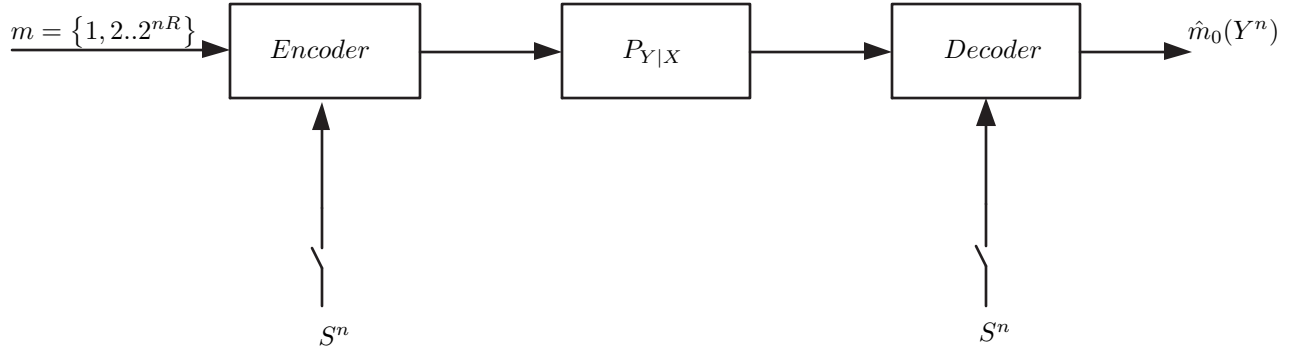


Fig. 2. Communication system with Side Information

$$P(y_i|x^i, s^i, y^{i-1}) = P(y_i|x_i, s_i) \quad (35)$$

$$P(y^n|x^n, s^n) = \prod_{i=1}^n P(y_i|x_i, s_i) \quad (\text{without feedback}) \quad (36)$$

Case I : state information known only to the Decoder:

$$C = \max_{p(x)} I(X; Y|S) \stackrel{(a)}{=} I(X; Y, S) \quad (37)$$

$$a) p(s)p(x)p(y|x, s) = p(x, y, s)$$

We define the problem:

$$\begin{aligned} \text{Encoder} & : f : \{1, 2, \dots, 2^{nR}\} \rightarrow X^n \\ \text{Decoder} & : g : Y^n, S^n \rightarrow \{1, 2, \dots, 2^{nR}\} \\ S^n & \sim p(s) \text{ i.i.d} \end{aligned}$$

For memoryless channel:

$$\begin{aligned} P(y_i, s_i|x^i, y^{i-1}, s^{i-1}) &= P(s_i)P(y_i|x^i, y^{i-1}, s^i) = \\ &= P(s_i)P(y_i|x_i, s_i) = P(y_i, s_i|x_i) \end{aligned}$$

Case II : state information known to both Encoder and Decoder:

$$C = \max_{p(x|s)} I(X; Y|S) \quad (38)$$

We define the problem:

$$\text{Encoder} : f : \{1, 2, \dots, 2^{nR}\} \times S^n \rightarrow X^n$$

$$\text{Decoder} : g : Y^n, S^n \rightarrow \{1, 2, \dots, 2^{nR}\}$$

proof :

$$C = \max_{p(x|s)} \sum_s p(s) I(X; Y | S = s) \quad (39)$$

We will split the message $\{1, 2, \dots, 2^{nR}\} : \{1, 2, \dots, 2^{nR_0}\} \times \{1, 2, \dots, 2^{nR_1}\}$

$$R_0 = I(X; Y | S = 0) \quad (40)$$

$$R_1 = I(X; Y | S = 1) \quad (41)$$

Doing that we split the channel into two separate channels:

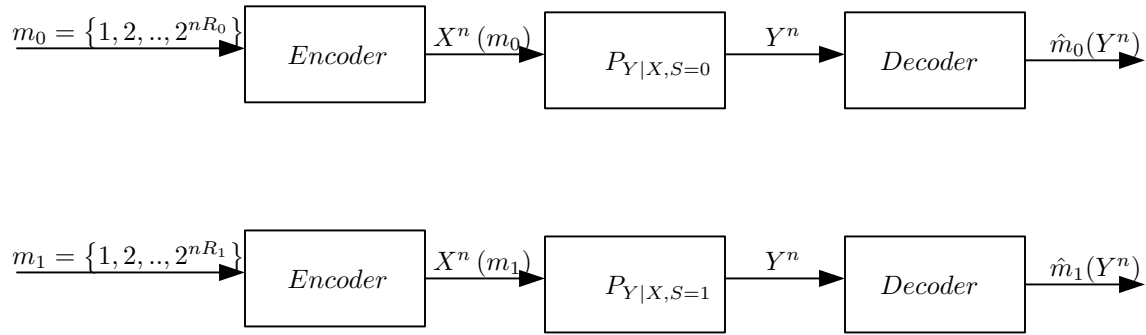


Fig. 3. Splitted Channel

such that:

$$R_0 = n \cdot p(S = 0) I(X; Y | S = 0)$$

$$\text{BlockSize}(S = 0) = n \cdot p(S = 0)$$

$$R_1 = n \cdot p(S = 1) I(X; Y | S = 1)$$

$$\text{BlockSize}(S = 1) = n \cdot p(S = 1)$$

Calculating the total rate of the channel:

$$R = \frac{1}{n} [np(S = 0)I(X; Y | S = 0) + np(S = 1)I(X; Y | S = 1)] \stackrel{(a)}{=} I(X; Y | S) \quad (42)$$

(a) Law of large numbers

Converse :

$$\begin{aligned}
 nR &= H(M) && (43) \\
 &\stackrel{(a)}{=} H(M|S^n) \\
 &= H(M|S^n) - H(M|S^n, Y^n) + H(M|S^n, Y^n) \\
 &\stackrel{(b)}{\leq} I(M; Y^n|S^n) + n\epsilon_n \\
 &= I(M, X^n(M, S^n); Y^n|S^n) + n\epsilon_n \\
 &= H(Y^n|S^n) - H(Y^n|S^n, X^n, M) + n\epsilon_n \\
 &\stackrel{(c)}{=} \sum_{i=1}^n H(Y_i|Y^{i-1}, S^n) - H(Y_i|Y^{i-1}, S^n, X^n, M) + n\epsilon_n \\
 &\leq \sum_{i=1}^n H(Y_i|S_i) - H(Y_i|S_i, X_i) + n\epsilon_n \\
 &= \sum_{i=1}^n I(Y_i; X_i|S_i) + n\epsilon_n \\
 &\leq \left[\max_{p(x|s)} I(Y; X|S) + \epsilon_n \right] n
 \end{aligned}$$

(a) $S^n \perp M$

(b) Fano inequality

(c) Chain rule