

1 Motivation: Indexing

Alice has a vector $x \in \{0, 1\}^n$. Bob has $i \in \{1, n\}$. $\text{INDEX}(x, i) = x_i$. Alice sends a message M to Bob based only on the vector x , and Bob computes $f(M, i) = \hat{x}_i$. If $x_i = \hat{x}_i$ with probability at least $\frac{2}{3}$, then we want to show that $|M| = \Omega(n)$. The technique used to prove such lower bounds is **information theory**.

2 Preliminaries

Consider distributions p over a finite support of size n . $p = (p_1, \dots, p_n)$, $p_i \in [0, 1]$, $\sum_i p_i = 1$. X is a random variable with distribution p i.e., $\Pr[X = i] = p_i$.

Definition (Entropy). Let X be an RV with dist p on n items. The entropy $H(X) = \sum_i p_i \log_2(\frac{1}{p_i})$. If $p_i = 0$, then $p_i \log_2(\frac{1}{p_i}) = 0$ (as in the limit computation).

Entropy is a way to capture the uncertainty of a distribution. For example, a unif. dist. would have high entropy, but a constant one would have low entropy.

Claim 1. $H(X) \leq \log_2 n$. Equality holds when $p_i = \frac{1}{n}, \forall i$. If $p_i = 1$ for some i and 0 elsewhere, the entropy would be 0.

Proof. $f(x) = -x \log x$ is concave. (The average is less than the function value of the average). $\sum_{i=1}^n \frac{1}{n} f(x_i) \leq f(\frac{1}{n} \sum x_i)$. Let $x_i = p_i$. $f(x) = x \log_2 \frac{1}{x}$. $f(\frac{1}{n}) \geq \frac{1}{n} \sum (p_i \log_2 \frac{1}{p_i})$. $\frac{1}{n} \log_2 n \geq \frac{1}{n} \sum (p_i \log_2 \frac{1}{p_i})$. Hence we have proven that $H(x) \leq \log_2 n$, with the equality attained by the uniform distribution. ■

Definition (Binary Input Entropy). If B is a bit with bias p , when $H(B) = p \log_2 \frac{1}{p} + (1-p) \log_2 \frac{1}{1-p}$.

The entropy is maximized at $p = \frac{1}{2}$.

Let X, Y be two RVs.

Definition (Conditional entropy). This is like the uncertainty of X over the choices of Y .

$$H(X | Y) = \sum_y H(X | Y = y) \Pr[Y = y]$$

Definition (Joint entropy). This is just like the usual entropy applied to the joint distribution.

$$H(X, Y) = \sum_{x,y} \Pr[(X, Y) = (x, y)] \log_2 \frac{1}{\Pr[(X, Y) = (x, y)]}$$

Fact 1 (Chain rule for entropy). $H(X, Y) = H(X) + H(Y | X)$

Proof.

$$H(X, Y) = \sum_{x, y} \Pr[(X, Y) = (x, y)] \log\left(\frac{1}{\Pr[(X, Y) = (x, y)]}\right) \quad (1)$$

$$= \sum_{x, y} \Pr[X = x] \Pr[Y = y | X = x] \log\left(\frac{1}{\Pr[X = x] \Pr[Y = y | X = x]}\right) \quad (2)$$

$$= \sum_{x, y} \Pr[X = x] \Pr[Y = y | X = x] \left(\log\left(\frac{1}{\Pr[X = x]}\right) + \log\left(\frac{1}{\Pr[Y = y | X = x]}\right) \right) \quad (3)$$

$$= H(X) + H(Y | X) \quad (4)$$

(1) using definition of joint entropy, (2) using definition of conditional probability, (3) splitting the log term using logarithm properties, (4) splitting the sum over x and y, using the definition of entropy. ■

3 Properties of Entropy

Fact 2 (Conditioning cannot increase entropy). $H(X | Y) \leq H(X)$
 "you know more from conditioning".

Proof. Using Jensen's inequality. Let f be a continuous concave function, let p be a distribution of support size n . For any x_1, \dots, x_n

$$\sum_{i=1}^n p_i f(x_i) \leq f\left(\sum_i p_i x_i\right)$$

Recall that $f(x) = \log x$ is concave.

$$H(X | Y) - H(X) = \sum_{x, y} \Pr[Y = y] \Pr[X = x | Y = y] \log\left(\frac{1}{\Pr[X = x | Y = y]}\right) \quad (5)$$

$$- \sum_x \Pr[X = x] \log\left(\frac{1}{\Pr[X = x]}\right) \sum_y \Pr[Y = y | X = x] \quad (6)$$

$$= \sum_{x, y} \Pr[X = x, Y = y] \log \frac{\Pr[X = x]}{\Pr[X = x | Y = y]} \quad (7)$$

$$= \sum_{x, y} \Pr[X = x, Y = y] \log \frac{\Pr[X = x] \Pr[Y = y]}{\Pr[(X, Y) = (x, y)]} \quad (8)$$

$$\leq \log\left(\sum_{x, y} \Pr[X = x, Y = y] \cdot \frac{\Pr[X = x] \Pr[Y = y]}{\Pr[(X, Y) = (x, y)]}\right) \quad (9)$$

$$= \log 1 = 0 \quad (10)$$

The last term of (6) is multiplying $H(X)$ by 1, expressed in terms of conditional probabilities. (7) uses the fact that difference of log is log of ratio. (8) uses definition of conditional probability. (9) uses Jensen's inequality, with the function being the log function. If X, Y are independent, $H(X | Y) = H(X)$. the log goes to 1 in (9). ■

Definition (Mutual information). how much info does X reveal about Y :

$$I(X; Y) = H(X) - H(X | Y) = H(Y) - H(Y | X) = I(Y; X)$$

Why is it symmetric? Using the chain rule, $H(X) + H(Y | X) = H(Y) + H(X | Y) = H(X, Y)$. If $X = Y$, the mutual information would just be the entropy of X . (X reveals everything about Y) If X and Y are independent, the mutual information would be 0. (X reveals nothing about Y).

Definition (Conditional mutual information). $I(X; Y | Z) = H(X | Z) - H(X | Y, Z)$

How does $I(X; Y | Z)$ compare to $I(X; Y)$? Both cases could happen.

Claim 2. For certain X, Y, Z , we can have $I(X; Y | Z) \leq I(X; Y)$. Consider $X = Y = Z$. Recall that $I(X; Y) = H(X)$. However, if we knew Z , there's no information left for X to reveal.

$$\begin{aligned} I(X; Y | Z) &= H(X | Z) - H(X | Y, Z) = 0 - 0 = 0 \\ I(X; Y) &= H(X) - H(X | Y) = H(X) - 0 = H(X) \end{aligned}$$

Claim 3. For certain X, Y, Z , we can have $I(X; Y | Z) \geq I(X; Y)$. Consider $X = Y + Z \pmod 2$, where X, Y are uniform in $\{0, 1\}$. Note that due to the modular arithmetics thing, $X + Y = Z \pmod 2$. Therefore, any two are p.w. independent but they aren't mutually independent.

$$\begin{aligned} I(X; Y | Z) &= H(X | Z) - H(X | Y, Z) = H(X) - 0 = 1 - 0 = 1 \\ I(x; y) &= h(x) - h(x | y) = H(X) - H(X) = 1 - 1 = 0 \end{aligned}$$

Intuitively, Y only reveals useful information about X after conditioning on Z .

Fact 3 (Chain rule for mutual info). $I(X, Y; Z) = I(X; Z) + I(Y; Z | X)$.

It's the information that X reveals about Z plus the information that Y reveals about Z , given that we know X .

$$\begin{aligned} I(X, Y; Z) &= H(X, Y) - H(X, Y | Z) && \text{definition} \\ &= H(X) + H(Y | X) - H(X | Z) - H(Y | X, Z) && \text{chain rule} \\ &= I(X; Z) + I(Y; Z | X) && \text{definition} \end{aligned}$$

By induction, this generalizes to n variables:

$$I(X_1, \dots, X_n; Z) = \sum_i (X_i; Z | X_1, \dots, X_{i-1})$$

Theorem 1 (Fano's Inequality). For any estimator $X' : X \rightarrow Y \rightarrow X'$ i.e., given random variable Y , X' is independent of X , if $P_e = \Pr[X' \neq X]$, we have

$$H(X | Y) \leq H(P_e) + P_e \cdot \log(|X| - 1)$$

where $|X|$ is the support size of X .

$H(P_e)$ is the bit entropy function.

Here $X \rightarrow Y \rightarrow X'$ is a Markov Chain, meaning X' and X are independent given Y . Think of X as the past, Y as the present, and X' as the future. "Past and future are conditionally independent given the present".

In the context of Fano's inequality, X' is one's 'best guess' about X given Y .

Consider the following special case example: $X \in \{0, 1\}$. Therefore, we are only dealing with two bits going through the Y channel. $H(X | Y) \leq H(P_e)$. Suppose that $Y = X$, then $X' = Y = X$, so $P_e = 0$. In this case, Fano's gives us that $H(P_e) = 0$ (checks out). Suppose that Y is independent of X . In this case, the smallest possible probability error is when we output a bit uniformly independently, which is $\frac{1}{2}$. Now, $H(P_e) = 1$, and $H(X | Y) = 1$.

To prove Fano's inequality, we need the data processing inequality.

Theorem 2 (Data Processing Inequality). *Suppose $X \rightarrow Y \rightarrow Z$ is a Markov Chain. Then $I(X; Y) \geq I(X; Z)$*

The moral takeaway is "no clever combination of the data can improve the estimation". Here, Z is the clever combination.

Proof (Data Processing).

$$I(X; Y, Z) = I(X; Z) + I(X; Y | Z) = I(X; Y) + I(X; Z | Y) \text{ (expand in two ways)}$$

Using chain rule, taking Y and Z out first respectively. Therefore, it suffices to show that $I(X; Z | Y) = 0$, since mutual information is nonnegative. ($H(X) \geq H(X | Y)$, using the lemma that conditioning only reduces entropy).

$$I(X; Z | Y) = H(X | Y) - H(X | Y, Z)$$

using the definition. But given Y , then X and Z are independent. So $H(X | Y, Z) = H(X | Y)$. A corollary is $H(X | Y) \leq H(X | Z)$, writing the original form into $H(X) - H(X | Y) \geq H(X) - H(X | Z)$. ■

Proof (Fano). Let $E = 1$ if $X' \neq X$, and $E = 0$ otherwise.

$$H(E, X | X') = H(X | X') + H(E | X, X') = H(X | X') \tag{11}$$

$$H(E, X | X') = H(E | X') + H(X | E, X') \leq H(P_e) + H(X | E, X') \tag{12}$$

(11) using chain rule, and the fact that E completely depends on X, X' . The first equality in (12) is just expanding in two ways, and the inequality is using the fact that conditioning cannot increase entropy. $H(P_e) = H(E) \geq H(E | X')$.

$$\begin{aligned} H(X | E, X') &= \Pr[E = 0]H(X | X', E = 0) + \Pr[E = 1]H(X | X', E = 1) \text{ (conditional prob)} \\ &\leq (1 - P_e) \cdot 0 + P_e \cdot \log_2(|X| - 1) \end{aligned}$$

Combining the two expressions, we get $H(X | X') \leq H(P_e) + P_e \cdot \log_2(|X| - 1)$ By data processing, $H(X | Y) \leq H(X | X') \leq H(P_e) + P_e \cdot \log_2(|X| - 1)$, since we are dealing with a Markov chain by assumption. ■

Statistically, Fano's inequality gives us a lower bound of the error probability. It measures how well we could predict X , in terms of the entropy of X given Y .

By the following example, we can see that Fano's inequality is actually tight: Suppose the distribution p of X satisfies $p_1 \geq p_2 \geq \dots \geq p_n$. Suppose Y is a constant distribution. So $I(X; Y) = H(X) - H(X | Y) = 0$. The Best predictor X' of X is $X = 1$. $P_e = \Pr[X' \neq X] = 1 - p_1$. $H(X | Y) \leq H(p_1) + (1 - p_1) \log_2(n - 1)$, as predicted by Fano's inequality. However, $H(X) = H(X | Y)$ and if $p_2 = p_3 = \dots = p_n = \frac{1-p_1}{n-1}$, the inequality is tight:

$$X \sim (p_1, \frac{1-p_1}{n-1}, \dots, \frac{1-p_1}{n-1}) \quad (13)$$

$$H(X) = \sum_i p_i \log\left(\frac{1}{p_i}\right) \quad (14)$$

$$= p_1 \log\left(\frac{1}{p_1}\right) + \sum_{i>1} \frac{1-p_1}{n-1} \log \frac{n-1}{1-p_1} \quad (15)$$

$$= p_1 \log \frac{1}{p_1} + \frac{1}{(n-1)} \sum_{i>1} (1-p_1) \log\left(\frac{1}{1-p_1}\right) + (1-p_1) \log(n-1) \quad (16)$$

$$= H(p_1) + (1-p_1) \log(n-1) \quad (17)$$

Y is a constant. So $H(X | Y) = H(X)$. (15), (16) uses log properties to foil the p_1 term out. Then the first two terms of (16) is the binary entropy function, and the last term is what Fano's gives. This chain of equalities shows that the equality in Fano's is attainable.