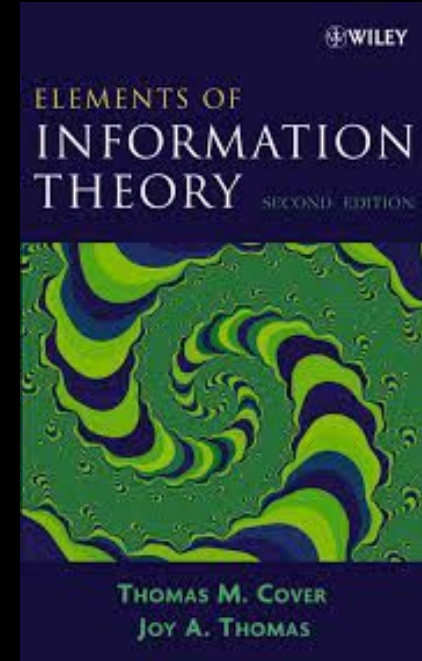"Information is the resolution of uncertainty."

Shannon

Erfan Mirzaei

Information Theory Mini-Course

Sep 2022

# References and Acknowledges

# Session 2

Joint, Conditional, Cross Entropy
Mutual Information
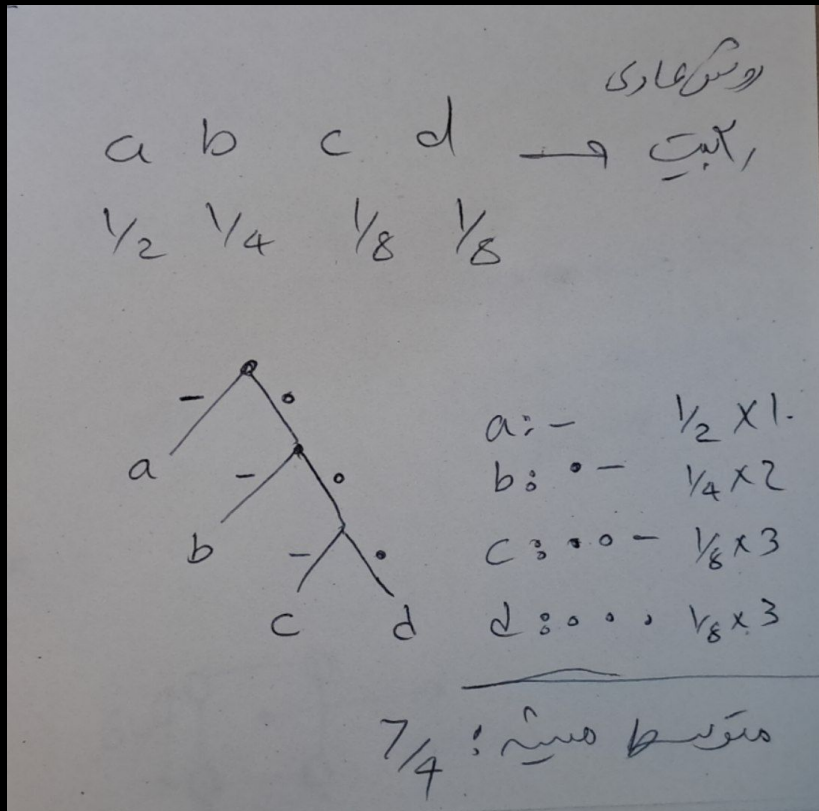
# Recap



- Imagine a 4 sided dice and you want guess the outcome Just by asking Y/N questions:

- A trivial but not efficient way is to ask exhaustively.
  This takes 9/4 questions in average.

- Another approach is that you assume that the toss is fair and ask search for the answer by binary search.
  This takes 2(8/4) questions in average. This is equal to entropy when the toss is fair which is the max entropy among all random variables with the same outcome space size.

- But imagine you know that the probability of numbers is as follows:[½,¼,⅛,⅛]
  With knowing this and choosing good strategy for asking questions you could find the answer with 7/4 questions in average.[Equal to Entropy.]

# Recap



# characters should be sent =
# stones should be thrown =
# questions should be asked


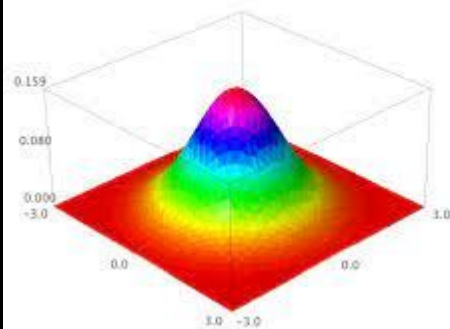This is the beauty of information Theory.

# Joint Entropy

Have two random variables X and Y with joint probability distribution P(X = x, Y = y) = P(x, y)



$$p_{X,Y}(x, y) = \mathrm{P}(X = x \text{ and } Y = y)$$

$$H(X, Y) = -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(x, y),$$

$$H(X, Y) = -E \log p(X, Y).$$

|   | X |   |   |
|---|---|---|---|
|   | 1 | 2 | 3 |
| 1 | 0 | 1/6 | 1/6 |
| Y 2 | 1/6 | 0 | 1/6 |
| 3 | 1/6 | 1/6 | 0 |

# Conditional Entropy

Have two random variables X and Y with joint probability distribution $P(X = x, Y = y) = P(x, y)$

$$
\begin{aligned}
H(Y|X) &= \sum_{x \in \mathcal{X}} p(x) H(Y|X = x) \\
&= -\sum_{x \in \mathcal{X}} p(x) \sum_{y \in \mathcal{Y}} p(y|x) \log p(y|x) \\
&= -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(y|x) \\
&= -E \log p(Y|X).
\end{aligned}
$$

# First Chain rule

The naturalness of definition of joint and conditional entropy is exhibited:

(Chain rule)

$$H(X, Y) = H(X) + H(Y|X).$$

$$P(A \mid B) = \frac{P(A \text{ and } B)}{P(B)}$$

**Corollary**

$$H(X, Y|Z) = H(X|Z) + H(Y|X, Z). \qquad (2.21)$$

- $H(Y, X) = H(X) + H(X \mid Y) = H(Y) + H(Y|X)$

# Proof

$$H(X, Y) = -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(x, y) \tag{2.15}$$

$$= -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(x) p(y|x) \tag{2.16}$$

$$= -\sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(x) - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(y|x) \tag{2.17}$$

$$= -\sum_{x \in \mathcal{X}} p(x) \log p(x) - \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log p(y|x) \tag{2.18}$$

$$= H(X) + H(Y|X). \tag{2.19}$$

Equivalently, we can write

$$\log p(X, Y) = \log p(X) + \log p(Y|X) \tag{2.20}$$

and take the expectation of both sides of the equation to obtain the theorem. □
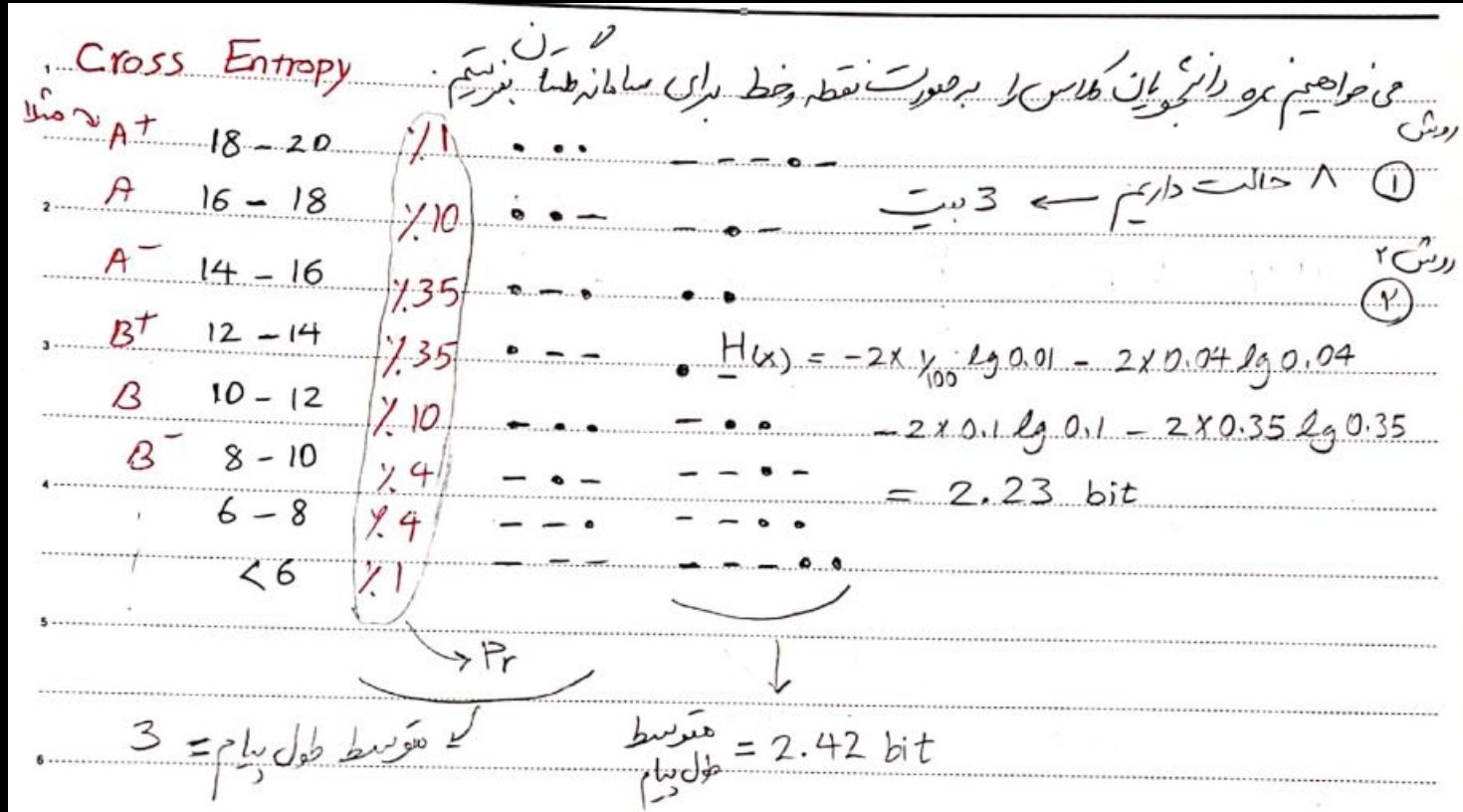
# What is the intuition?

When you have uncertainty about two random variables you can think about it in a sequential manner:



- $H(Y \mid X) <= H(Y)$

  When we know X, we can not have more uncertainty about

  Y than we know nothing. [Or revealing doesn't add more info]

# Cross Entropy

# Cross Entropy

# Cross Entropy



Cross Entropy

$P$ : توزیع واقعی

$q$ : توزیع مفروض

$$H(p,q) = -\sum_{x \in X} p(x) \, lg \, q(x)$$

اگر توزیع مفروض با توزیع واقعی برابر باشد

Cross Entropy = Entropy

$H(P) = H(X) \longrightarrow 2.23 \; bit$

$H(P,q) \longrightarrow 4.58 \; bit$

خطای ما در تخمین احتمالات (یعنی q به جای P)

باعث شده که نویزی داشته باشیم representation

و این باعث افزایش یقین (اطلاعات) شده.

$H(P)$

# Cross Entropy Loss

$$H_p(q) = -\frac{1}{N} \sum_{i=1}^{N} y_i \cdot log(p(y_i)) + (1 - y_i) \cdot log(1 - p(y_i))$$

Binary Cross-Entropy / Log Loss

"Minimizing this KL divergence corresponds exactly to minimizing the cross-entropy between the distributions."
Page 132, Deep Learning, 2016                    Reference

Any loss consisting of a negative log-likelihood is a cross-entropy between the empirical distribution defined by the training set and the probability distribution defined by model. For example, mean squared error is the cross-entropy between the empirical distribution and a Gaussian model.

# Relative Entropy OR KL Divergence



$$H(P,q) - H(P) = -\sum P(x) \lg q(x) + \sum P(x) \lg P(x)$$

$$= \sum P(x) \lg \frac{P(x)}{q(x)} = D(P\|q) \quad \text{Relative Entropy}$$

$$D(P\|q) \neq D(q\|P)$$

Kullback-Liebler Divergence (KL divergence)

- How much p and q as two probability distribution is similar to each other?

- How much we had more uncertainty just because we don't know the probability distribution. We sent longer messages. And revealing the outcome was more surprising to us and gave us more information.

# Mutual Information

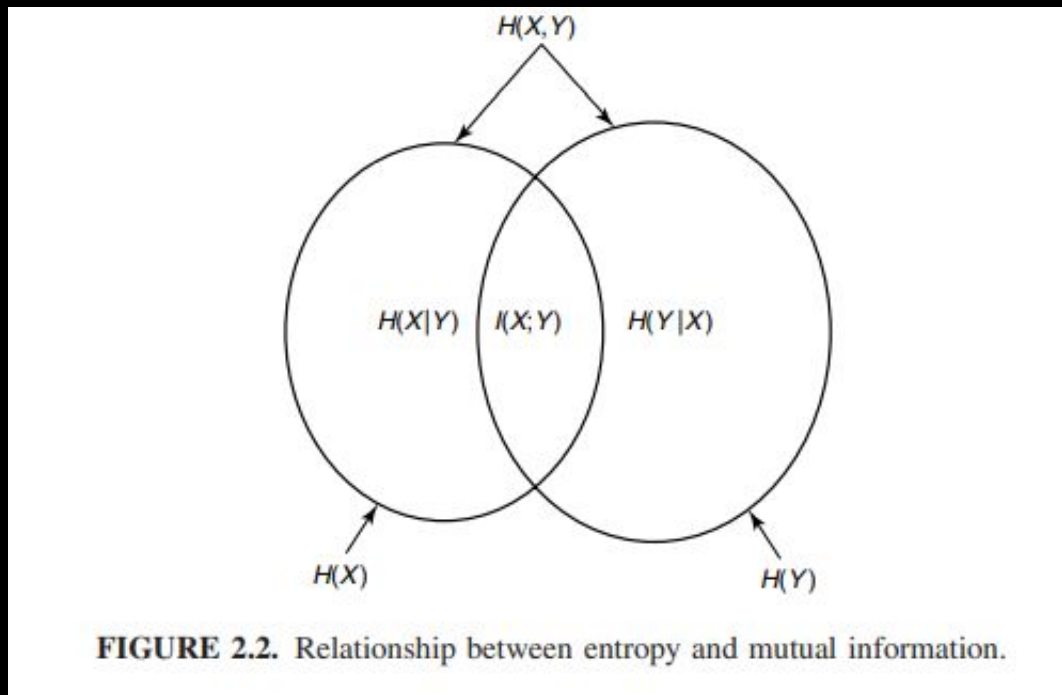Have two random variables X and Y with joint probability distribution $P(X = x, Y = y) = P(x, y)$

- Relative Entropy or KL-Divergence between P(x,y) and P(x)*P(y)

$$I(X; Y) = \sum_{x \in \mathcal{X}} \sum_{y \in \mathcal{Y}} p(x, y) \log \frac{p(x, y)}{p(x) p(y)}$$
$$= D(p(x, y) \| p(x) p(y))$$
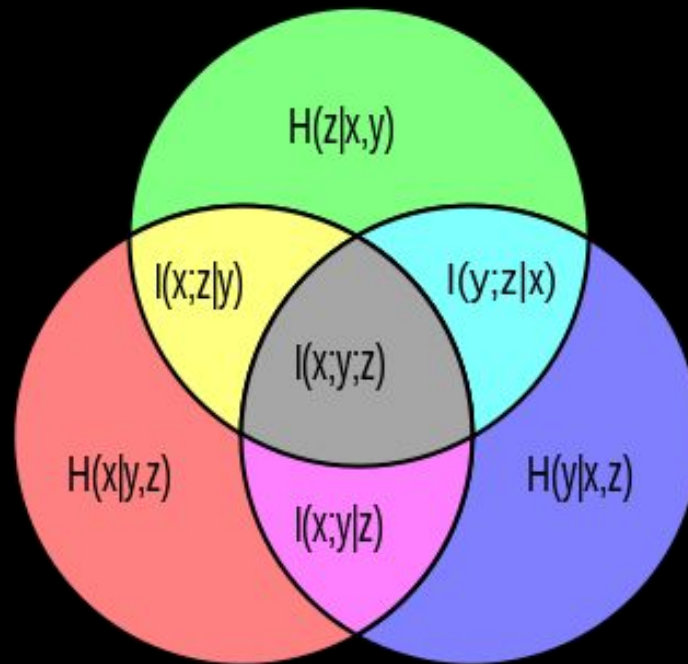$$= E_{p(x,y)} \log \frac{p(X, Y)}{p(X) p(Y)}.$$

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

- $I(X ; Y) = H(X) - H(X \mid Y) = H(Y) - H(Y \mid X) = H(X) + H(Y) - H(X, Y)$

# Mutual Information and Entropy



**FIGURE 2.2.** Relationship between entropy and mutual information.

# Mutual Information and Entropy

# Chain Rules

$$H(X_1, X_2, \ldots, X_n) = H(X_1) + H(X_2|X_1) + \cdots + H(X_n|X_{n-1}, \ldots, X_1)$$

$$\text{(2.52)}$$

$$= \sum_{i=1}^{n} H(X_i|X_{i-1}, \ldots, X_1). \quad \square \qquad \text{(2.53)}$$

$$I(X_1, X_2, \ldots, X_n; Y)$$

$$= H(X_1, X_2, \ldots, X_n) - H(X_1, X_2, \ldots, X_n|Y) \qquad \text{(2.63)}$$

$$= \sum_{i=1}^{n} H(X_i|X_{i-1}, \ldots, X_1) - \sum_{i=1}^{n} H(X_i|X_{i-1}, \ldots, X_1, Y)$$

$$= \sum_{i=1}^{n} I(X_i; Y|X_1, X_2, \ldots, X_{i-1}). \quad \square \qquad \text{(2.64)}$$

# Thanks for your attention

Erfan Mirzaei
erfunmirzaei@gmail.com