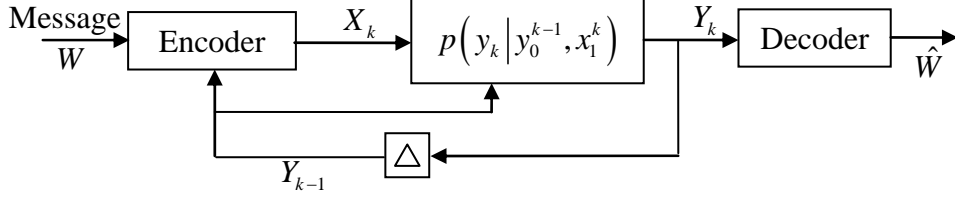


Directed Information and Channel with Feedback

Prapun Suksompong

December 10, 2006

In this article, we consider the problem of sending message via discrete channel with noiseless feedback as shown below.



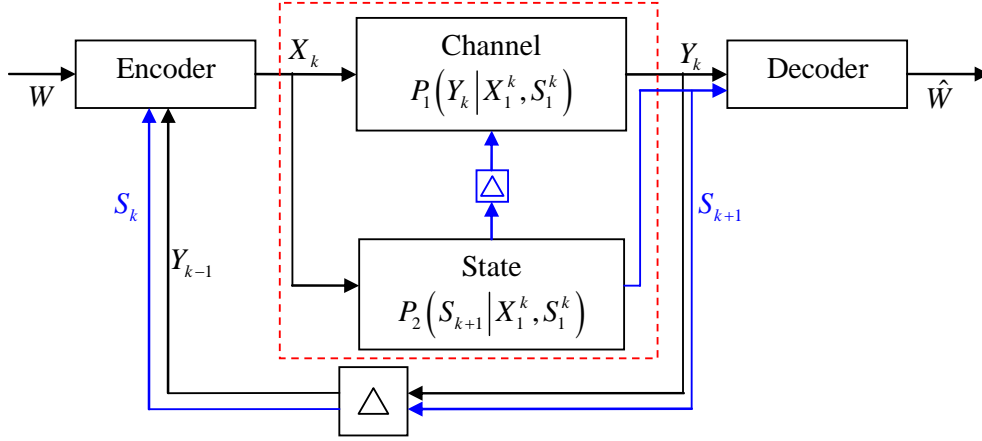
At time k , the encoder produce X_k from W , the previous encoder outputs X_1^{k-1} , and the fed-back channel outputs Y_0^{k-1} . The encoder can be a deterministic function; that is $X_k = f_k(X_1^{k-1}, Y_0^{k-1})$. We also allow random encoder; that is X_k may be governed by the conditional distribution $p(x_k | x_1^{k-1}, y_0^{k-1})$. This X_k is put into channel. The channel output Y_k is generated according to the conditional distribution $p(y_k | x_1^k, y_0^{k-1})$ which not only depends on the current channel input X_k but may also depend on past channel inputs X_1^{k-1} and outputs Y_0^{k-1} .

We modify the classical Shannon's information measures (entropy, mutual information, and their conditional versions) so that they explicitly incorporate feedback. In particular, we shall focus on $I(X^N \rightarrow Y^N)$, a notation introduced by Massey [1990] to capture the directed information flowing from the length N sequence of random variables X^N to the length N sequence of random variables Y^N . In fact, the idea of given direction to information has already been thought of by Marko [Marko 1973] whose paper define a quantity called directed transinformation. The direct information in [Massey 1990] refines this directed transinformation.

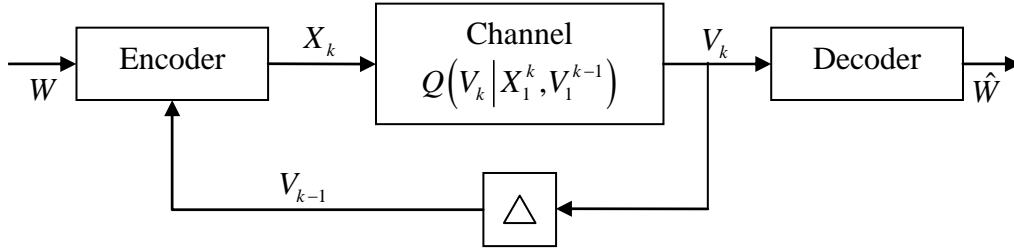
Equivalent System Models

In this section, we show three general models of the system we are interested in. We then prove that they are all equivalent in the sense that we can convert one into another. Hence, in some sense, it is sufficient to analyze only one of them. As usual, k is the time index. X_k and Y_k represent the channel input and output at time k respectively. S_k denotes the channel state at time k . In model 2, we use V_k instead of Y_k to represent the channel output. This notational difference is used so that the proof can be done more smoothly. Model 3 is introduced as an intermediate model to bridge model 1 and model 2. The definitions and proof are almost the same as those in [Chen, Suksompong, and Berger 2004]. What follows is conditioned on the initial state $S_1 = s_1$.

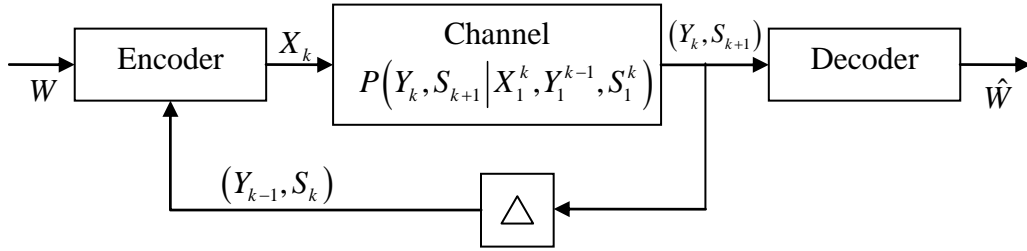
- Model 1:



- Model 2:



- Model 3:



- *Theorem:* Model 1, 2, and 3 are equivalent.

Proof Let \mathcal{M}_i be the set of systems that can be represented by Model i . Then,

(1) $\mathcal{M}_1 \subset \mathcal{M}_3$ because we can let $P(Y_k, S_{k+1} | X_1^k, Y_1^{k-1}, S_1^k) = P_1(Y_k | X_1^k, S_1^k) P_2(S_{k+1} | X_1^k, S_1^k)$.

(2) $\mathcal{M}_3 \subset \mathcal{M}_2$ because we can let $V_k = (Y_k, S_{k+1})$. Then, $V_1^{k-1} = (Y_1^{k-1}, S_1^k)$. Set $Q(V_k | X_1^k, V_1^{k-1}) = P(Y_k, S_{k+1} | X_1^k, Y_1^{k-1}, S_1^k)$.

(3) $\mathcal{M}_2 \subset \mathcal{M}_1$ because we can set $S_k = V_{k-1}$, $P_2(S_{k+1} | X_1^k, S_1^k) = Q(V_k | X_1^k, V_1^{k-1})$, and $Y_k = c$ a.s. (Equivalently, $\exists c \forall (x_1^k, s_1^k) P_1(Y_k = c | X_1^k = x_1^k, S_1^k = s_1^k) = 1$.)

Henceforth, we shall focus on Model 2.

Discrete Channel

In this section, we consider the channel part of the system. Special cases of which are also defined.

- **Discrete channel:** $\left(p(y_n | x_1^n, y_1^{n-1}) \right)_{n=1}^N$.
 - [Tatikonda 2000] calls this nonanticipative channel.
- **Definition:** [Cover Thomas 1991] A channel is **memoryless** if

$$p(y_n | x_1^n, y_1^{n-1}) = p(y_n | x_n).$$

Simple equivalent statements are

$$\begin{aligned} &\equiv \text{Given } x_n, (x_1^{n-1}, y_1^{n-1}) \text{ and } y_n \text{ are independent;} \\ &\equiv y_n - x_n - (x_1^{n-1}, y_1^{n-1}) \text{ forms a Markov chain;} \\ &\equiv H(Y_n | X_1^n, Y_1^{n-1}) = H(Y_n | X_n); \\ &\equiv \forall i \ I(Y_i; (X_1^{j-1}, Y_1^{j-1}) | X_j) = 0. \end{aligned}$$

- **Definition** [Massey 1990]: The channel is **used without feedback** if

$$p(x_n | x_1^{n-1}, y_1^{n-1}) = p(x_n | x_1^{n-1}).$$

Equivalent statements are

$$\begin{aligned} &\equiv H(X_n | X_1^{n-1}, Y_1^{n-1}) = H(X_n | X_1^{n-1}); \\ &\equiv I(X_n; Y_1^{n-1} | X_1^{n-1}) = 0; \\ &\equiv Y_1^{n-1} - X_1^{n-1} - X_n \text{ forms a Markov chain;} \\ &\equiv [\text{Ash 1965}] \text{ for } \forall N > n, \ p(y_n | x_1^N, y_1^{n-1}) = p(y_n | x_1^n, y_1^{n-1}); \end{aligned}$$

(Note that this gives exactly the same condition as the definition above when $N = 2$.)

$$\begin{aligned} &\equiv H(Y_n | X_1^N, Y_1^{n-1}) = H(Y_n | X_1^n, Y_1^{n-1}); \\ &\equiv \forall i \ I(Y_i; X_{i+1}^N | X_1^i, Y_1^{i-1}) = 0; \\ &\equiv \forall i \ Y_i - (X_1^i, Y_1^{i-1}) - X_{i+1}^N \text{ forms a Markov chain.} \\ &\equiv I(X_i; Y_j | X_1^{i-1}, Y_1^{j-1}) = 0 \text{ for any } j < i. \end{aligned}$$

We shall show later that these are also equivalent to:

$$\begin{aligned} &\equiv I(X_1^N; Y_1^N) = I(X_1^N \rightarrow Y_1^N) \\ &\equiv I(0 * Y_1^{N-1} \rightarrow X_1^N) = 0. \end{aligned}$$

Interpretation: The choice of the next channel input digit, given all previous input digits, is not further related to the previous channel output digits.

- **Joint distribution:**

$$p(x_1^N, y_1^N) = \prod_{n=1}^N p(x_n, y_n | x_1^{n-1}, y_1^{n-1}) = \prod_{n=1}^N p(x_n | x_1^{n-1}, y_1^{n-1}) p(y_n | x_1^n, y_1^{n-1}).$$

- Discrete **memoryless** channel: $p(x_1^N, y_1^N) = \prod_{n=1}^N p(x_n | x_1^{n-1}, y_1^{n-1}) p(y_n | x_n).$

- Discrete channel used **without feedback**:

$$\begin{aligned} p(x_1^N, y_1^N) &= \prod_{n=1}^N p(x_n | x_1^{n-1}) p(y_n | x_1^n, y_1^{n-1}) = \prod_{n=1}^N p(x_n | x_1^{n-1}) \prod_{n=1}^N p(y_n | x_1^n, y_1^{n-1}) \\ &= p(x_1^N) \prod_{n=1}^N p(y_n | x_1^n, y_1^{n-1}) \end{aligned}$$

or equivalently,

$$p(y_1^N | x_1^N) = \prod_{n=1}^N p(y_n | x_1^n, y_1^{n-1})$$

- Used [Ash], $p(y_n | x_1^N, y_1^{n-1}) = p(y_n | x_1^n, y_1^{n-1})$, we have

$$p(x_1^N, y_1^N) = p(x_1^N) \prod_{n=1}^N p(y_n | x_1^N, y_1^{n-1}) = p(x_1^N) \prod_{n=1}^N p(y_n | x_1^n, y_1^{n-1});$$

hence equivalent to Massey's.

- Discrete **memoryless** channel used **without feedback**:

$$p(x_1^N, y_1^N) = p(x_1^N) \prod_{n=1}^N p(y_n | x_n), \text{ or equivalently } p(y_1^N | x_1^N) = \prod_{n=1}^N p(y_n | x_n).$$

This is **DMC used without feedback**.

In this case,

- $H(X_1^N, Y_1^N) = H(X_1^N) + \sum_{i=1}^N H(Y_i | X_i).$
- $H(Y_1^N | X_1^N) = \sum_{i=1}^N H(Y_i | X_i)$

As an introduction to directed information, we present here one way to look at it. First we shall decompose mutual information into several conditional mutual information using chain rules. We associated each part with an arrow going between X_i and Y_j in a diagram representing the general channel in figure 1. The above specialized channels are then just the general channel with some of the arrows missing. Directed information takes only some of these parts and hence it is a part of mutual information. This idea was presented to Prof. T. Berger and L. Tong during the author's Q exam in June, 2005. This representation of directed information trivializes several of its properties.

- The definitions defined above can be represented using the diagrams.

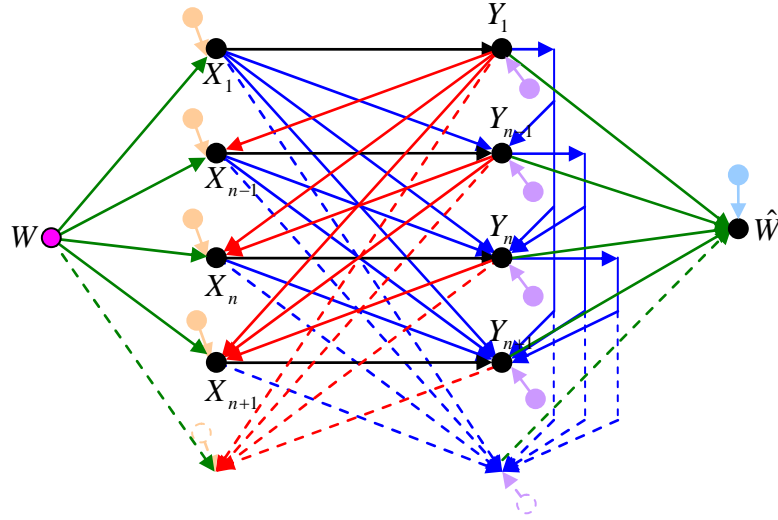


Figure 1 Discrete Channel

Figure 1 shows the discrete channel in its full generality. We can, in fact, decompose it into three parts which are shown in black, blue, and red. Consider, first, the part in black as shown in figure 2.

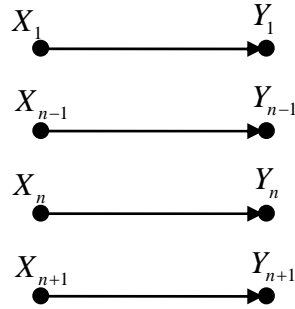


Figure 2 Discrete Memoryless Channel without Feedback

Figure 2 represents the discrete memoryless channel without feedback. We can add memory to channel by adding the blue part, as shown in figure 3. This allows y_k to depends on (x_1^k, y_1^{k-1}) , not just x_k .

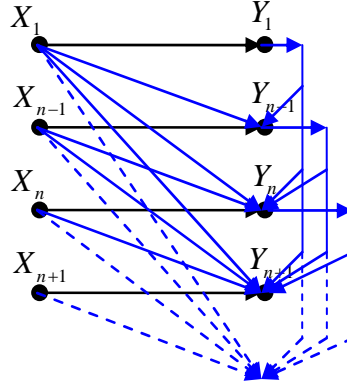


Figure 3 Discrete Channel without Feedback

The last part, with red color in figure 1, represents the feedback. Adding it to figure 3 leads us back the general discrete channel in figure 1.

- The mutual information $I(X_1^N; Y_1^N)$ delivered via the discrete channel described above can also be partitioned into three parts.

First, by applying the chain rule for mutual information twice, we have

$$I(X_1^N; Y_1^N) = \sum_{i=1}^N I(X_i; Y_1^N | X_1^{i-1}) = \sum_{i=1}^N \sum_{j=1}^N I(X_i; Y_j | X_1^{i-1}, Y_1^{j-1}).$$

Now, we separate the terms inside the sum above into three groups by asking whether i is equal to, less than, or greater than j . This gives

- 1) $i = j$: $\sum_{i=1}^N I(X_i; Y_i | X_1^{i-1}, Y_1^{i-1})$. This quantity will be defined as $I(X_1^N \leftrightarrow Y_1^N)$.

We want to say that it relates to the “direct” paths in black above.

- 2) $i < j$: $\sum_{j=1}^N \sum_{i=1}^{j-1} I(X_i; Y_j | X_1^{i-1}, Y_1^{j-1}) = \sum_{j=1}^N I(X_1^{j-1}; Y_j | Y_1^{j-1})$. This quantity will be defined as $I(0 * X_1^{N-1} \rightarrow Y_1^N)$. We want to say that it relates to the channel memory paths in blue above.

- 3) $j < i$: $\sum_{i=1}^N \sum_{j=1}^{i-1} I(X_i; Y_j | X_1^{i-1}, Y_1^{j-1}) = \sum_{i=1}^N I(X_i; Y_1^{i-1} | X_1^{i-1})$. This quantity will be defined as $I(0 * Y_1^{N-1} \rightarrow X_1^N)$. We want to say that it captures the feedback paths in red above.

The directed information only adds up the terms that has $i \leq j$ (all arrows that point from X to Y), that is

$$\begin{aligned} I(X_1^N \rightarrow Y_1^N) &= I(X_1^N \leftrightarrow Y_1^N) + I(0 * X_1^{N-1} \rightarrow Y_1^N) \\ &= I(X_1^N; Y_1^N) - I(0 * Y_1^{N-1} \rightarrow X_1^N) \end{aligned}$$

For the familiar DMC without feedback, we will show that

$$I(X_1^N \rightarrow Y_1^N) = I(X_1^N; Y_1^N) \leq \sum_{i=1}^N I(Y_i; X_i) \leq nC.$$

Directed Information and its properties

- [Kramer 1998] Assume that the sequences X_1^N and Y_1^N are “synchronized”, i.e., that the n^{th} terms in the sequences occur “at the same time”, and that the n^{th} terms occur “before” the $(n+1)^{\text{st}}$ terms.
- [Massey 1990] The [directed information](#) $I(X_1^N \rightarrow Y_1^N)$ from a sequence X_1^N to a sequence Y_1^N is defined by $I(X_1^N \rightarrow Y_1^N) = \sum_{i=1}^N I(X_1^i; Y_i | Y_1^{i-1})$.

• *Remarks:*

- $I(X_1^N; Y_1^N) = \sum_{i=1}^N I(X_1^N; Y_i | Y_1^{i-1})$ by chain rule. We add the information that Y_i tells about X_1^N which Y_1^{i-1} haven't already told.
- For directed information, we use $I(X_1^n; Y_n | Y_1^{n-1})$: To eliminate feedback information, we ignore the information that Y_i may be giving about the future X_{i+1}^N .
- $$I(X_1^N \rightarrow Y_1^N) = \sum_{i=1}^N \mathbb{E} \left[\log \frac{p(X_1^i, Y_i | Y_1^{i-1})}{p(X_1^i | Y_1^{i-1}) p(Y_i | Y_1^{i-1})} \right]$$
$$= \mathbb{E} \left[\log \frac{\prod_{i=1}^N p(Y_i | X_1^i, Y_1^{i-1})}{p(Y_1^N)} \right]$$
- $I(Y_1^N \rightarrow X_1^N) = \sum_{i=1}^N I(Y_i; X_i | X_1^{i-1})$
- $I(X_1^N \rightarrow Y_1^N) \leq I(X_1^N; Y_1^N)$ with equality iff channel is used without feedback

$$\begin{aligned} \text{Proof. } I(X_1^N; Y_i | Y_1^{i-1}) &= H(Y_i | Y_1^{i-1}) - H(Y_i | X_1^N, Y_1^{i-1}) \\ &\geq H(Y_i | Y_1^{i-1}) - H(Y_i | X_1^i, Y_1^{i-1}) \\ &= I(X_1^i; Y_i | Y_1^{i-1}) \end{aligned}$$

Hence,

$$I(X_1^N; Y_1^N) = \sum_{i=1}^N I(X_1^N; Y_i | Y_1^{i-1}) \geq \sum_{i=1}^N I(X_1^i; Y_i | Y_1^{i-1}) = I(X_1^N \rightarrow Y_1^N).$$

Equality occurs iff $\forall i \ H(Y_i | X_1^N, Y_1^{i-1}) = H(Y_i | X_1^i, Y_1^{i-1})$. This is Ash's condition which is equivalent to the “used without feedback” condition.

Alternative Proof.

Use $I(X_1, X_2; Y | Z) = I(X_1; Y) + I(X_2; Y | X_1, Z)$. Then,

$$I(X_1^N; Y_i | Y_1^{i-1}) - I(X_1^i; Y_i | Y_1^{i-1}) = I(X_{i+1}^N; Y_i | X_1^i, Y_1^{i-1}) \geq 0$$

This proof give us the next property.

- $I(X_1^N; Y_1^N) - I(X_1^N \rightarrow Y_1^N) = \sum_{i=1}^N I(X_{i+1}^N; Y_i | X_1^i, Y_1^{i-1})$; in fact, the last term is 0; so

$$I(X_1^N; Y_1^N) - I(X_1^N \rightarrow Y_1^N) = \sum_{i=1}^{N-1} I(X_{i+1}^N; Y_i | X_1^i, Y_1^{i-1}).$$

So, $I(X_1^N \rightarrow Y_1^N) \leq I(X_1^N; Y_1^N)$ with equality iff $\forall i \ Y_i - (X_1^i, Y_1^{i-1}) - X_{i+1}^N$ forms a Markov chain.

- If $X_{i+1}^N - X_1^i - Y_1^i$ is a Markov chain (i.e. $p(X_{i+1}^N | X_1^i, Y_1^i) = p(X_{i+1}^N | X_1^i)$), then $I(X_1^N; Y_1^N) = I(X_1^N \rightarrow Y_1^N)$.
- The Markov chain condition above says that the future X 's are not influenced by the past and current Y 's when conditioned on the past and current X 's. In the context of channels this state that if there is no feedback, then the two different mutual information measures are equal. [Tatikonda 2000, p. 81]
- Pearl [1988] call this the “weak union” property of conditional independence.

Proof. $p(X_{i+1}^N | X_1^i, Y_1^i) = p(X_{i+1}^N | X_1^i) \Rightarrow p(X_{i+1}^N | X_1^i, Y_1^i) = p(X_{i+1}^N | X_1^i, Y_1^{i-1}) \Rightarrow Y_i - (X_1^i, Y_1^{i-1}) - X_{i+1}^N$ Markov.

(Recall that $p(Z | V, U_1, U_2) = p(Z | V) \Rightarrow p(Z | V, U_1, U_2) = p(Z | V, U_1) = p(Z | V, U_2) = p(Z | V)$.)

- For DMC, $I(X_1^N \rightarrow Y_1^N) \leq \sum_{i=1}^N I(X_i; Y_i)$ with equality iff Y_i are independent.

Proof. $I(X_1^i; Y_i | Y_1^{i-1}) = H(Y_i | Y_1^{i-1}) - H(Y_i | X_1^i, Y_1^{i-1})$
 $= H(Y_i | Y_1^{i-1}) - H(Y_i | X_i); \text{memoryless}$
 $\leq H(Y_i) - H(Y_i | X_i)$
 $= I(X_i; Y_i)$

Recall that $H(Y_1^N) = \sum_{i=1}^N H(Y_i | Y_1^{i-1}) \leq \sum_{i=1}^N H(Y_i)$ with equality iff Y_i are independent.

- *Definition:* Denote the sequence $(0, Y_1, \dots, Y_{N-1})$ by DY_1^N [Kramer 1998], or $0 * Y_1^{N-1}$ [Massey 1990].
- The letter D represent delay by one time step (with discard of the last component).

- $I(0 * Y_1^{N-1} \rightarrow X_1^N) = \sum_{i=1}^N I(0 * Y_1^{i-1}; X_i | X_1^{i-1}) = \sum_{i=1}^N I(Y_1^{i-1}; X_i | X_1^{i-1})$
 $= \sum_{i=2}^N I(Y_1^{i-1}; X_i | X_1^{i-1})$

- Define $I(X_1^N \leftrightarrow Y_1^N) = I(Y_1^N \leftrightarrow X_1^N) = \sum_{i=1}^N I(Y_i; X_i | Y_1^{i-1}, X_1^{i-1})$.

- $I(X_1^N \rightarrow Y_1^N) = I(0 * X_1^{N-1} \rightarrow Y_1^N) + I(X_1^N \leftrightarrow Y_1^N)$

$$I(Y_1^N \rightarrow X_1^N) = I(0 * Y_1^{N-1} \rightarrow X_1^N) + I(Y_1^N \leftrightarrow X_1^N)$$

Proof. $I(Y_1^i; X_i | X_1^{i-1}) - I(Y_1^{i-1}; X_i | X_1^{i-1}) = I(Y_i; X_i | Y_1^{i-1}, X_1^{i-1})$

$$I(Y_1^N \rightarrow X_1^N) - I(0 * Y_1^{N-1} \rightarrow X_1^N)$$

$$= \sum_{i=1}^N I(Y_i; X_i | X_1^{i-1}) - \sum_{i=1}^N I(Y_1^{i-1}; X_i | X_1^{i-1})$$

$$= \sum_{i=1}^N I(Y_i; X_i | Y_1^{i-1}, X_1^{i-1})$$

- Conservation Law:

$$I(X_1^N; Y_1^N) = I(X_1^N \rightarrow Y_1^N) + I(0 * Y_1^{N-1} \rightarrow X_1^N)$$

$$= I(Y_1^N \rightarrow X_1^N) + I(0 * X_1^{N-1} \rightarrow Y_1^N)$$

- Equivalently, $\sum_{i=2}^N I(Y_1^{i-1}; X_i | X_1^{i-1}) = \sum_{i=1}^{N-1} I(X_{i+1}^N; Y_i | X_1^i, Y_1^{i-1})$.

- *Interpretation:* By putting a 0 in the front,

- We shift the delay position from the feedback to the channel. Before, X_i and Y_i are treated as if they happened during the same time step.

- Now, the formula doesn't include the "middle" part which is $I(X_1^N \leftrightarrow Y_1^N)$.

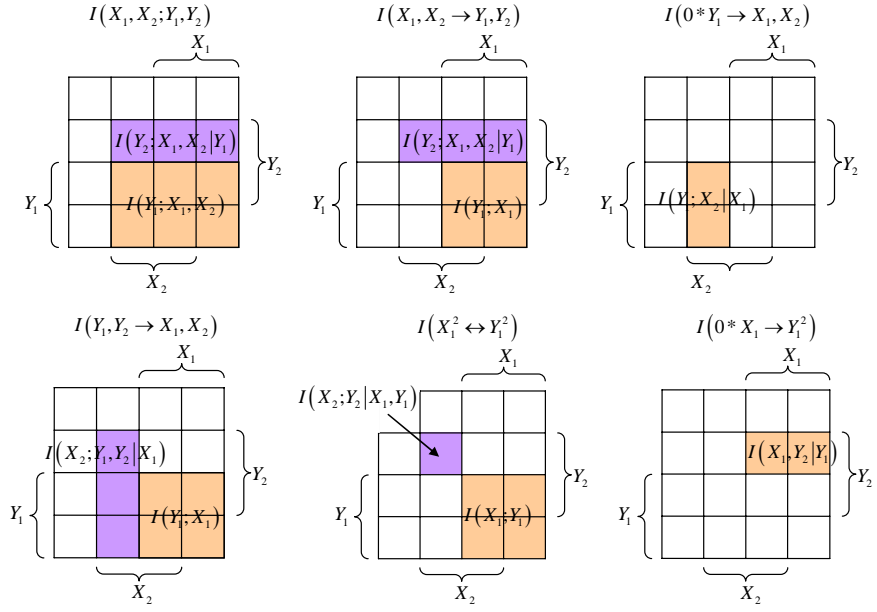
- $I(X_1^N; Y_1^N) = I(X_1^N \rightarrow Y_1^N) + I(Y_1^N \rightarrow X_1^N) - I(X_1^N \leftrightarrow Y_1^N)$

- In general, $I(X_1^N; Y_1^N) \neq I(X_1^N \rightarrow Y_1^N) + I(Y_1^N \rightarrow X_1^N)$.

- This is obvious when $N = 1$ because

$$I(X_1^N; Y_1^N) = I(X_1^N \rightarrow Y_1^N) = I(Y_1^N \rightarrow X_1^N).$$

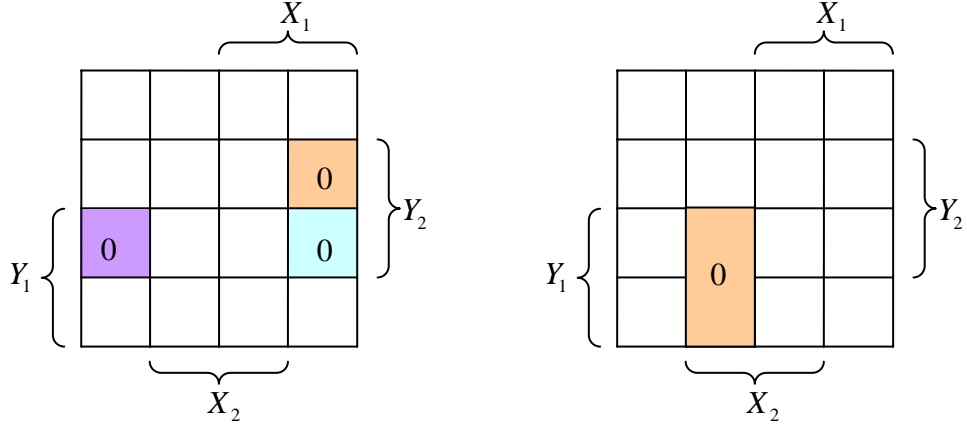
- *Example: (I-measure) $N = 2$*



- DMC: $I(Y_2; (X_1, Y_1) | X_2) = 0$ which equivalent to $I(Y_2; X_1 | X_2, Y_1) = 0$, $I(Y_2; Y_1 | X_1, X_2) = 0$, and $I(X_1; Y_1; Y_2 | X_2) = 0$.

DMC: $I(Y_2; (X_1, Y_1) | X_2) = 0$

No feedback: $I(X_2; Y_1 | X_1) = 0$



We can clearly see now that when there is no feedback, $I(X_1^2; Y_1^2) = I(X_1^2 \rightarrow Y_1^2)$.

- We summarize the definitions and properties involving the directed information below:

- $I(X_1^N; Y_1^N) = \sum_{i=1}^N I(X_1^N; Y_i | Y_1^{i-1})$
- $I(X_1^N \rightarrow Y_1^N) = \sum_{i=1}^N I(X_1^i; Y_i | Y_1^{i-1}) = H(Y_1^N) - \sum_{i=1}^N H(Y_i | X_1^i, Y_1^{i-1})$
- $I(0 * Y_1^{N-1} \rightarrow X_1^N) = \sum_{i=2}^N I(Y_1^{i-1}; X_i | X_1^{i-1}) = \sum_{i=1}^{N-1} I(X_{i+1}^N; Y_i | X_1^i, Y_1^{i-1})$
 $= I(X_2; Y_1 | X_1) + I(X_3; Y_1^2 | X_1^2) + \dots$

$$I(0 * X_1^{N-1} \rightarrow Y_1^N) = \sum_{i=1}^N I(Y_i; X_1^{i-1} | Y_1^{i-1})$$

$$= I(X_1; Y_2 | Y_1) + I(X_1^2; Y_3 | Y_1^2) + \dots$$

- $I(X_1^N \leftrightarrow Y_1^N) = I(Y_1^N \leftrightarrow X_1^N) = \sum_{i=1}^N I(Y_i; X_i | Y_1^{i-1}, X_1^{i-1})$
- $I(X_1^N \rightarrow Y_1^N) = I(0 * X_1^{N-1} \rightarrow Y_1^N) + I(X_1^N \leftrightarrow Y_1^N)$
 $I(Y_1^N \rightarrow X_1^N) = I(0 * Y_1^{N-1} \rightarrow X_1^N) + I(Y_1^N \leftrightarrow X_1^N)$
- $I(X_1^N; Y_1^N) = I(X_1^N \rightarrow Y_1^N) + I(Y_1^N \rightarrow X_1^N) - I(X_1^N \leftrightarrow Y_1^N)$
 $= I(0 * X_1^{N-1} \rightarrow Y_1^N) + I(0 * Y_1^{N-1} \rightarrow X_1^N) + I(X_1^N \leftrightarrow Y_1^N)$

- [Tatikonda 2000] If the process $\{p(x_1^N, y_1^N)\}_{N=1}^\infty$ is information stable [p. 89], then $\lim_{N \rightarrow \infty} \frac{1}{N} I(X_1^N \rightarrow Y_1^N)$ exists and we can work directly with $I(X_1^N \rightarrow Y_1^N)$ (instead of liminf in probability as defined in [Tatikonda 2000 p. 89] and [Verdú 1994]).

- **Definition: [Kramer's] causal conditioning**

- We see above that the channel outputs are given by

$$p(y_1^N \| x_1^N) = \prod_{n=1}^N p(y_n | x_1^n, y_1^{n-1}).$$

This lead us to define

$$H(Y_1^N \| X_1^N) = -\mathbb{E} \left[p(Y_1^N \| X_1^N) \right] = \sum_{n=1}^N H(Y_n | X_1^n, Y_1^{n-1}).$$

- Again, $H(Y_1^N \| X_1^N)$ differs from the conditional entropy $H(Y_1^N | X_1^N)$ only in that X_1^n replaces X_1^{n-1} .
- The term “causal” reflect the conditioning on past and present values of the sequence X_1^N only.
- It differs from “free information” [Marko 1973] only in that X_n is included in the conditioning.
- $p(x_1^N \|^\circ y_1^N) = \prod_{n=1}^N p(x_n | x_1^{n-1}, y_1^{n-1});$

$$H(X_1^N \|^\circ Y_1^N) = -\mathbb{E} \left[\log p(X_1^N \|^\circ Y_1^N) \right] = \sum_{n=1}^N H(X_n | X_1^{n-1}, Y_1^{n-1}).$$

Note the **asymmetry** in the definitions above.

- $p(x_1^N, y_1^N) = p(x_1^N \|^\circ y_1^N) p(y_1^N \| x_1^N).$
 $H(X_1^N, Y_1^N) = H(X_1^N \|^\circ Y_1^N) + H(Y_1^N \| X_1^N).$
- $\tilde{p}(x_1^N, y_1^N) = p(x_1^N \|^\circ y_1^N) p(y_1^N) = p(y_1^N) \prod_{n=1}^N p(x_n | x_1^{n-1}, y_1^{n-1})$

- $p(x_1^N \| y_1^N) = \prod_{n=1}^N p(x_n | x_1^{n-1}, y_1^n)$

$$H(X_1^N \| Y_1^N) = -\mathbb{E}[\log p(X_1^N \| Y_1^N)] = \sum_{n=1}^N H(X_n | X_1^{n-1}, Y_1^n).$$

- $H(X_1^N \| Y_1^N) \leq H(X_1^N \|^\circ Y_1^N)$

Proof. Conditioning only reduces entropy.

- $H(X_1^N \|^\circ Y_1^N) - H(X_1^N \| Y_1^N) = I(X_1^N \leftrightarrow Y_1^N)$

Proof.
$$\begin{aligned} H(X_1^N \|^\circ Y_1^N) - H(X_1^N \| Y_1^N) &= \sum_{n=1}^N (H(X_n | X_1^{n-1}, Y_1^{n-1}) - H(X_n | X_1^{n-1}, Y_1^n)) \\ &= \sum_{n=1}^N (I(X_n; Y_n | X_1^{n-1}, Y_1^{n-1})) \\ &= I(X_1^N \leftrightarrow Y_1^N) \end{aligned}$$

- $I(X_1^N \rightarrow Y_1^N) = H(Y_1^N) - H(Y_1^N \| X_1^N).$

Proof.
$$\begin{aligned} I(X_1^N \rightarrow Y_1^N) &= \sum_{i=1}^N I(X_1^i; Y_i | Y_1^{i-1}) = H(Y_1^N) - \sum_{i=1}^N H(Y_i | X_1^i, Y_1^{i-1}) \\ &= H(Y_1^N) - H(Y_1^N \| X_1^N) \end{aligned}$$

Alternative Proof. More directly,

$$I(X_1^N \rightarrow Y_1^N) = \mathbb{E} \left[\log \frac{\prod_{i=1}^N p(Y_i | X_1^i, Y_1^{i-1})}{p(Y_1^N)} \right] = \mathbb{E} \left[\log \frac{p(Y_1^N \| X_1^N)}{p(Y_1^N)} \right].$$

- $$\begin{aligned} I(X_1^N \rightarrow Y_1^N) &= \mathbb{E} \left[\log \frac{p(Y_1^N \| X_1^N)}{p(Y_1^N)} \right] = \mathbb{E} \left[\log \frac{p(X_1^N, Y_1^N)}{p(X_1^N \|^\circ Y_1^N) p(Y_1^N)} \right] \\ &= D(p(x_1^N, y_1^N) \| p(x_1^N \|^\circ y_1^N) p(y_1^N)) \end{aligned}$$

- $I(0 * Y_1^{N-1} \rightarrow X_1^N) = \sum_{i=1}^N I(Y_1^{i-1}; X_i | X_1^{i-1}) = H(X_1^N) - H(X_1^N \|^\circ Y_1^N)$

Proof.
$$I(0 * Y_1^{N-1} \rightarrow X_1^N) = \sum_{i=1}^N I(Y_1^{i-1}; X_i | X_1^{i-1}).$$

$$I(Y_1^{i-1}; X_i | X_1^{i-1}) = H(X_i | X_1^{i-1}) - H(X_i | Y_1^{i-1}, X_1^{i-1}).$$

- $I(X_1^N \rightarrow Y_1^N) = H(X_1^N \|^\circ Y_1^N) - H(X_1^N | Y_1^N)$

Proof.
$$\begin{aligned} I(X_1^N \rightarrow Y_1^N) &= H(Y_1^N) - H(Y_1^N \| X_1^N) \\ &= H(Y_1^N) - (H(X_1^N, Y_1^N) - H(X_1^N \|^\circ Y_1^N)) \\ &= H(X_1^N \|^\circ Y_1^N) - H(X_1^N | Y_1^N) \end{aligned}$$

- *Definition:* $I(X_1^N \rightarrow Y_1^N \| Z_1^N) = H(Y_1^N \| Z_1^N) - H(Y_1^N \| X_1^N, Z_1^N)$.

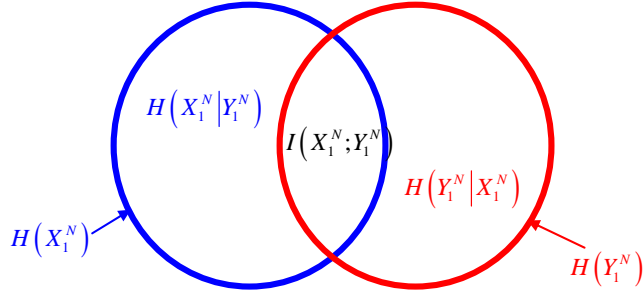
- $I(X_1^N \rightarrow Y_1^N \| Z_1^N) = \sum_{n=1}^N I(Y_n; X_1^n | Y_1^{n-1}, Z_1^n)$

Proof.
$$\begin{aligned} I(X_1^N \rightarrow Y_1^N \| Z_1^N) &= \sum_{n=1}^N H(Y_n | Y_1^{n-1}, Z_1^n) - H(Y_n | Y_1^{n-1}, X_1^n, Z_1^n) \\ &= \sum_{n=1}^N H(Y_n | Y_1^{n-1}, Z_1^n) - H(Y_n | Y_1^{n-1}, X_1^n, Z_1^n) \\ &= \sum_{n=1}^N I(Y_n; X_1^n | Y_1^{n-1}, Z_1^n) \end{aligned}$$

- Let $Z_1^N = 0 * X_1^{N-1}$, then $I(X_1^N \rightarrow Y_1^N \| 0 * X_1^{N-1}) = I(X_1^N \leftrightarrow Y_1^N)$.

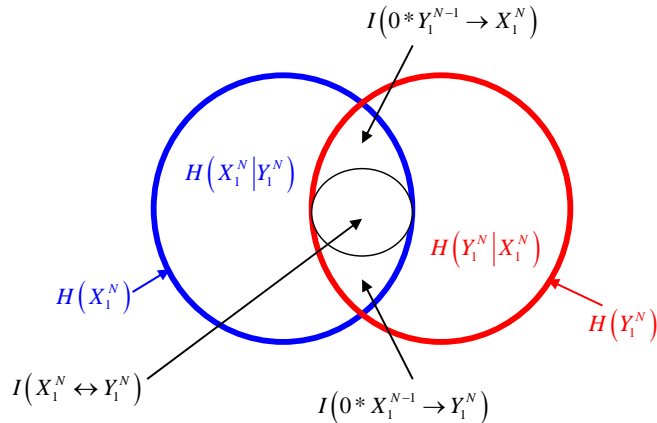
Proof.
$$\begin{aligned} I(X_1^N \rightarrow Y_1^N \| 0 * X_1^{N-1}) &= \sum_{n=1}^N I(Y_n; X_1^n | Y_1^{n-1}, X_1^{n-1}) \\ &= \sum_{n=1}^N I(Y_n; X_n | Y_1^{n-1}, X_1^{n-1}) \end{aligned}$$

- Similarly, by symmetry, we have $I(Y_1^N \rightarrow X_1^N \| 0 * Y_1^{N-1}) = I(X_1^N \leftrightarrow Y_1^N)$.
- We conclude this section with diagrams. All the diagrams below are the same diagram. First, we have the familiar diagram:

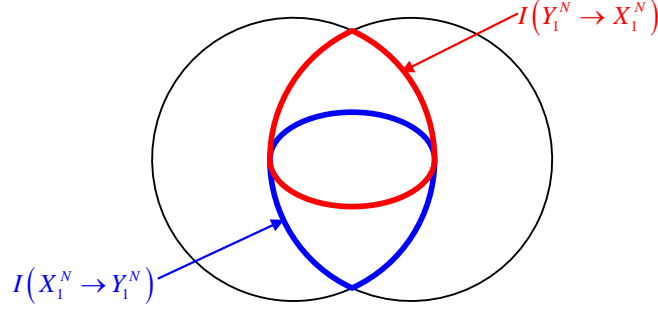


Next, the mutual information $I(X_1^N; Y_1^N)$ is partitioned into three subsets:

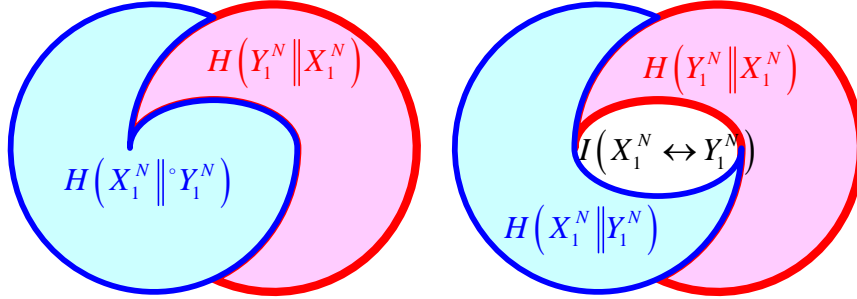
$$I(X_1^N; Y_1^N) = I(0 * X_1^{N-1} \rightarrow Y_1^N) + I(0 * Y_1^{N-1} \rightarrow X_1^N) + I(X_1^N \leftrightarrow Y_1^N).$$



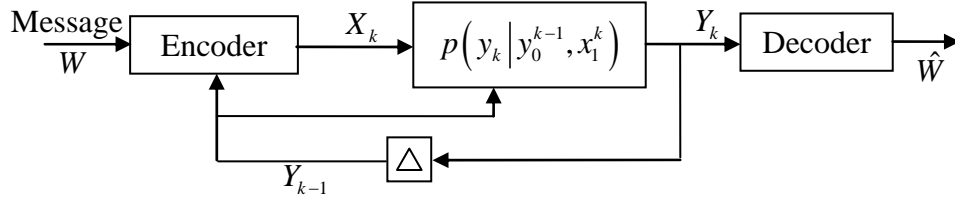
$I(X_1^N \rightarrow Y_1^N)$ and $I(Y_1^N \rightarrow X_1^N)$ are parts of $I(X_1^N; Y_1^N)$:



Finally, $H(X_1^N, Y_1^N) = H(X_1^N \| Y_1^N) + H(Y_1^N \| X_1^N)$
 $= H(X_1^N \| Y_1^N) + H(Y_1^N \| X_1^N) + I(X_1^N \leftrightarrow Y_1^N)$



Converse Channel Coding Theorem



- **Assume** that the system is causal; that is $p(y_i | x_1^i, y_1^{i-1}, w) = p(y_i | x_1^i, y_1^{i-1})$.

The idea is that the source output sequence should be thought of a specified prior to the process of sending sequences over channels and the channel should be aware of such sequences only via its past inputs and outputs and its current input.

- So, $p(w, x_1^N, y_1^N) = p(w) \prod_{n=1}^N p(x_n | x_1^{n-1}, y_1^{n-1}, w) p(y_n | x_1^n, y_1^{n-1})$.
- $I(W; Y_0^N) \leq I(X_1^N \rightarrow Y_1^N | Y_0)$ which is $\leq I(X_1^N; Y_1^N | Y_0)$.

Proof. $I(W; Y_0^N) = H(Y_0^N) - H(Y_0^N | W) = \sum_{k=1}^N (H(Y_k | Y_0^{k-1}) - H(Y_k | Y_0^{k-1}, W))$.

Because X_k is produced by (Y_0^{k-1}, W) , we have X_1^k is produced by $f(Y_0^{k-1}, W)$, and therefore

$$H(Y_k | Y_0^{k-1}, W) = H(Y_k | Y_0^{k-1}, X_1^k, W) \stackrel{(a)}{=} H(Y_k | Y_0^{k-1}, X_1^k),$$

where (a) comes from the causality assumption.

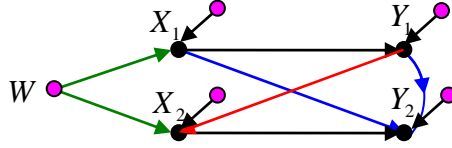
A more general X_k gives

$$H(Y_k | Y_0^{k-1}, W) \geq H(Y_k | Y_0^{k-1}, X_1^k, W) \stackrel{(a)}{=} H(Y_k | Y_0^{k-1}, X_1^k).$$

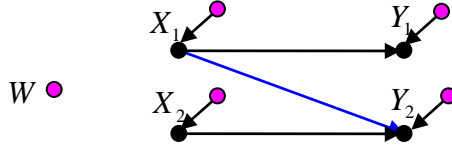
In any case,

$$I(W; Y_0^N) \leq \sum_{k=1}^N \left(H(Y_k | Y_0^{k-1}) - H(Y_k | Y_0^{k-1}, X_1^k) \right) = \sum_{k=1}^N I(Y_k; X_1^k | Y_0^{k-1}).$$

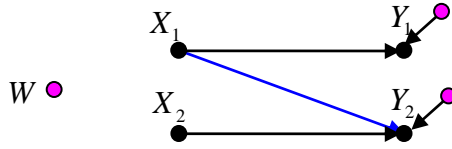
- *Remark:* Let's reconsider the equality of $H(Y_k | Y_0^{k-1}, W)$ and $H(Y_k | Y_0^{k-1}, X_1^k, W)$ using functional dependence graphs as in [Kramer 1998]. Suppose $k = 2$, then the relevant graph is shown below:



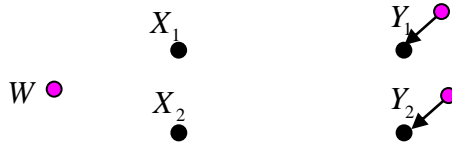
We then remove the arrows coming out of (W, Y_1) . The resulted graph shows that (W, Y_1) does not d -separate Y_2 from X_1^2 . Since all secondary random variables have incoming branches, we also know that (W, Y_1) does not fd -separate Y_2 from X_1^2 .



Now, instead of generating X_n by conditional distribution $p(x_n | x_1^{n-1}, y_1^{n-1})$, we use deterministic encoder; that is $X_n = f_n(X_1^{n-1}, Y_1^{n-1})$. Then, after removing the arrows coming out of (W, Y_1) , the secondary random variables X_1 and X_2 have no incoming branches.



Hence, we can further delete the arrows coming out of X_1 and X_2 .



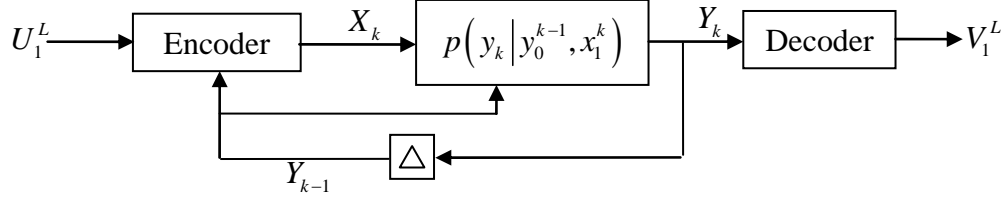
We conclude that (W, Y_1) fd -separate Y_2 from X_1^2 in the case of deterministic encoder.

- Feedback does not increase the capacity of DMC.

Proof. For DMC, we know that $I(X_1^N \rightarrow Y_1^N | Y_0) \leq \sum_{i=1}^N I(X_i; Y_i | Y_0)$. Hence,

$$I(X_1^N \rightarrow Y_1^N | Y_0) \leq \sum_{i=1}^N I(X_i; Y_i) \leq NC_{\text{DMC w/o feedback}}.$$

- This first result in the information theory of feedback channel is due to Shannon [1948].
- In fact, if we replace W by U_1^L and replace \hat{W} by V_1^L .



Then,

- For any discrete channel, we have $I(U_1^L; Y_0^N) \leq I(X_1^N \rightarrow Y_1^N | Y_0)$.
- Furthermore, if the channel is memoryless, we have

$$I(X_1^N \rightarrow Y_1^N | Y_0) \leq \sum_{i=1}^N I(X_i; Y_i | Y_0).$$

References

- C. Shannon, "The zero error capacity of a noisy channel," 1956.
- Ash, Robert B., Information theory. Corrected reprint of the 1965 original. Dover Publications, Inc., New York, 1990. xii+339 pp.
- J. L. Massey, "Causality, Feedback and Directed Information," pp. 303-305 in Proc. 1990 Int. Symp. on Info. Th. & its Appls., Hawaii, USA, Nov. 27-30, 1990.
- H. Marko, "The Bidirectional Communication Theory A Generalization of Information Theory", IEEE Trans. Commun., vol. COM-21, pp. 1345-1351, Dec. 1973.
- S. Tatikonda, "Control Under Communication Constraints." MIT Ph.D. thesis, August 2000.
- G. Kramer, Directed Information for Channels with Feedback, ETH Series in Inform. Proc., vol. 11. Konstanz: Hartung--Gorre, 1998. (PhD Thesis)
 - G. Kramer, "Capacity Results for the Discrete Memoryless Network," IEEE Trans. Inform. Th., vol. IT-49, pp. 4-21, Jan. 2003.
- J. Chen, P. Suksompong and T. Berger, "Communication through a Finite-State Machine with Markov Property", CISS 2004