

# Directed Information and Conditional Mutual Information

Peter Harremoës

Department of Mathematics

University of Copenhagen

Universitetsparken 5, DK-2100 Copenhagen O

Denmark

Email: moes@math.ku.dk

**Abstract—**We study directed information in Bayesian networks and related structures. Mutual information is split into directed information and residual information. Some basic equations for directed information and residual information are determined.

## I. INTRODUCTION

Shannon introduced mutual information and demonstrated how useful this quantity is to analyze a communication setup [1]. For instance the capacity of a simple information channel equals the maximum mutual information between input and output where the maximum is taken over all input distributions. This result is remarkable in that mutual information is symmetric in the two arguments but an information channel is highly unsymmetric in the sense that the sender has (partly or full) control over the input, but the sender has only indirect control over the distribution of output. Thus, the concept of mutual information does not capture that information flows from the input to the output.

In channels with feedback there is both a flow of information from input to output and a flow from output to input. The concept of directed information has been introduced to analyse channels with feedback, see [2], [3], [4], [5] and . Here we shall go one step further in order to understand what it means that information flows. This question is both of philosophical, theoretical and practical interest. There are many situations where we need more precise concepts in order to understand to what extend information is flowing. We just mention quantum entanglement and secrecy sharing as concepts for which it is not so obvious whether these concepts represent quantities that can flow.

In this paper we shall be more modest and restrict our attention to Bayesian networks, because these networks are important in modelling causality. Thus, for Bayesian networks there is a natural direction of causality and time.

## II. BAYESIAN NETWORKS

By a variable we shall understand any measurable mapping from a probability space into a finite space. The finiteness condition is assumed for simplicity and most of our results can easily be extended to variables that can take infinitely many values. Variables will be denoted with capital letters. If  $X$  and  $Y$  are variables then  $X \cup Y$  shall denote the vector valued variable  $(X, Y)$ . If  $(X_i)$  is a (finite or infinite)

sequence of variables then  $X_m^n$  shall denote the sequence  $(X_m, X_{m+1}, \dots, X_n)$  where  $m \leq n$ . Bold face capital letters will indicate infinite sequences.

If the variables  $X$  and  $Y$  are independent for each value  $z$  of the variable  $Z$  we say that  $X$  and  $Y$  are *conditionally independent given  $Z$*  and write  $(X \perp Y | Z)$ . The relation conditional independence is a *semi graphoid relation*, which means that

$$(X \perp Y | Z) \iff (Y \perp X | Z) \quad (1)$$

and

$$(X \perp Y \cup Z | W) \iff (X \perp Y | W) \text{ and } (X \perp Z | Y \cup W). \quad (2)$$

We introduce *conditional mutual information*  $I(X; Y | Z)$  and note that  $I(X; Y | Z) \geq 0$  with equality if and only if  $(X \perp Y | Z)$ . The relations (1) and (2) for conditional independence have counterparts for conditional independence:

$$I(X; Y | Z) = I(Y; X | Z) \quad (3a)$$

$$I(X; Y \cup Z | W) = I(X; Y | W) + I(X; Z | Y \cup W). \quad (3b)$$

Note that these equations implies the semi graphoid relations. We also get the important relation

$$I(X; Y | Z) = I(X \cup Z; Y) - I(Z; Y) \quad (4)$$

that can be used to define conditional mutual information from mutual information. Now (3a) and (4) together implies (3b).

Semi graphoid relations have proved important in understanding causality [7], but they are difficult to classify. In modelling causality one therefore prefer semi graphoid relations described by Bayesian networks or related graphical models. Here we shall focus on Bayesian networks. A *Bayesian network* is a directed acyclic graph with a variable associated to each node in the graph. If there is a directed path from the variable  $X$  to the variable  $Y$  in the network we shall write  $X \ll Y$  and note that  $\ll$  is a partial ordering of the nodes and the associated variables. The *ancestor set*  $a(Y)$  is defined as the set of variables  $X$  such that there exists a path to a variable in  $Y$ .

We shall restrict to *synchronious networks* consisting of a number of layers on top of each other such that any arrow in

the graph points from one layer to the next. For such networks the layers shall form a Markov chain where any variables in one layer is independent of all other variables in the same layer given its parents. We shall write  $X \prec Y$  if  $X$  is in a lower layer than  $Y$ , and  $X \approx Y$  if  $X$  and  $Y$  are in the same layer. For each variable there is a Markov kernel from the parents of the variables to the variable and the joint distribution on all variables is determined by these Markov kernels. There exists a synchronization of a network if and only all directed paths from one variable to another are of equal length. By introduction of dummy variables it is easy to synchronize any network and therefore the restriction to synchronous networks gives no loss of generality.

The previous description of a Bayesian network is static in the sense that it relates a graph to some variables with a probabilistic structure but with no time evolution. We shall now introduce *dynamic Bayesian networks*. For each variable  $X$  in a static Bayesian network we associate an double infinite sequence  $\mathbf{X}$  of variables. The joint distribution on all the sequences should be a stationary Markov chain with the Markov kernel being the product of the Markov kernels from the parents of a variable to the variable. The dynamic Bayesian network can be considered as a Bayesian network with infinitely many nodes/variables. The static Bayesian network can then be considered as the stationary distribution of the dynamic Bayesian network.

For a stationary sequence  $\mathbf{X}$  we define the *entropy rate* in the usual way, i.e.

$$\begin{aligned} H(\mathbf{X}) &= \lim \frac{1}{N} H(X_1^N) \\ &= \lim \frac{1}{N} \sum_{n=1}^N H(X_n | X_1^{n-1}) \\ &= H(X_1 | X_{-\infty}^0). \end{aligned}$$

Similarly we define the *conditional mutual information rate* by

$$I(\mathbf{X}; \mathbf{Y} | \mathbf{Z}) = \lim_{N \rightarrow \infty} \frac{1}{N} I(X_1^N; Y_1^N | Z_1^N).$$

For sequences of variables  $\mathbf{X}$ ,  $\mathbf{Y}$  and  $\mathbf{Z}$  in a dynamic Bayesian network we have

$$I(\mathbf{X}; \mathbf{Y} | \mathbf{Z}) = I(X; Y | Z),$$

where  $X, Y$  and  $Z$  are the associated variables in the static Bayesian network. The relation can be used to translate concepts for sequences into concepts for static Bayesian networks. Note that in general

$$I(X_1^N; Y_1^N | Z_1^N) \neq I(X; Y | Z).$$

Let  $X$  denote a sequence. Then we let  $T(\mathbf{X})$  denote the *delayed sequence* such that  $T(\mathbf{X})_n = X_{n-1}$ .

### III. DIRECTED AND RESIDUAL INFORMATION

Consider two sequences of random variables  $X^N$  and  $Y^N$ . Then the *directed information* from  $X^N$  to  $Y^N$  is defined by

$$\begin{aligned} \tilde{I}(X^N \rightarrow Y^N) &= \\ &= \frac{1}{N} \sum_{n=1}^N H(Y_n | Y_1^{n-1}) - H(Y_n | X_1^{n-1} \cup Y_1^{n-1}) \\ &= \frac{1}{N} \sum_{n=1}^N I(Y_n; X_1^{n-1} | Y_1^{n-1}) \end{aligned}$$

If  $\mathbf{X} = (X_n)_{n \in \mathbb{Z}}$  and  $\mathbf{Y} = (Y_n)_{n \in \mathbb{Z}}$  are stationary processes then the *directed information rate* from  $\mathbf{X}$  to  $\mathbf{Y}$  is defined as the limit

$$\tilde{I}(\mathbf{X} \rightarrow \mathbf{Y}) = \lim_{N \rightarrow \infty} \tilde{I}(X^N \rightarrow Y^N).$$

It easy to prove that the limit exists and equals

$$H(Y_1 | Y_{-\infty}^0) - H(Y_1 | X_{-\infty}^0 \cup Y_{-\infty}^0).$$

The mutual information rate between two sequences of random variables  $X^N$  and  $Y^N$  is given by

$$\begin{aligned} & \frac{1}{N} I(X^N; Y^N) \\ &= \frac{1}{N} (H(X^N) + H(Y^N) - H(X^N \cup Y^N)) \\ &= \frac{1}{N} \sum_{n=1}^N \left( \begin{array}{c} H(X_n | X_1^{n-1}) + H(Y_n | Y_1^{n-1}) \\ - H(X_n \cup Y_n | X_1^{n-1} \cup Y_1^{n-1}) \end{array} \right) \\ &= \frac{1}{N} \sum_{n=1}^N \left( \begin{array}{c} H(X_n | X_1^{n-1}) - H(X_n | X_1^{n-1} \cup Y_1^{n-1}) \\ + H(Y_n | Y_1^{n-1}) - H(Y_n | X_1^{n-1} \cup Y_1^{n-1}) \\ + \left( \begin{array}{c} H(X_n | X_1^{n-1} \cup Y_1^{n-1}) \\ + H(Y_n | X_1^{n-1} \cup Y_1^{n-1}) \\ - H(X_n \cup Y_n | X_1^{n-1} \cup Y_1^{n-1}) \end{array} \right) \end{array} \right) \\ &= \tilde{I}(X^N \rightarrow Y^N) + \tilde{I}(Y^N \rightarrow X^N) \\ & \quad + \frac{1}{N} \sum_{n=1}^N I(X_n; Y_n | X_1^{n-1} \cup Y_1^{n-1}). \end{aligned}$$

We see that a mutual information rate is a sum of two directed information terms plus a term that we shall call the *residual information*. Thus

$$\tilde{I}_{res}(X^N; Y^N) = \frac{1}{N} \sum_{n=1}^N I(X_n; Y_n | X_1^{n-1} \cup Y_1^{n-1}).$$

The residual information measures how much the sequences deviates from being a Bayesian network of the form depicted in Figure 1. For stationary sequences we define the residual information as the limit

$$\begin{aligned} \tilde{I}_{res}(\mathbf{X}; \mathbf{Y}) &= \lim_{N \rightarrow \infty} \tilde{I}_{res}(X^N; Y^N) \\ &= I(X_1; Y_1 | X_{-\infty}^0 \cup Y_{-\infty}^0). \end{aligned}$$

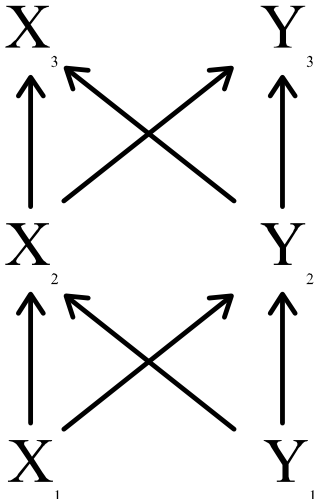


Fig. 1. Two sequences with zero residual information.

Thus we get

$$\frac{1}{N} I(X^N; Y^N) = \tilde{I}(X^N \rightarrow Y^N) + \tilde{I}(Y^N \rightarrow X^N) + \tilde{I}_{res}(X^N; Y^N)$$

and

$$I(\mathbf{X}; \mathbf{Y}) = \tilde{I}(\mathbf{X} \rightarrow \mathbf{Y}) + \tilde{I}(\mathbf{Y} \rightarrow \mathbf{X}) + \tilde{I}_{res}(\mathbf{X}; \mathbf{Y}). \quad (5)$$

*Example 1:* Let  $\mathbf{X}$  and  $\mathbf{Y}$  be sequences of variables in a dynamical Bayesian network corresponding to variables  $X$  and  $Y$  associated with single nodes in the graph. Then

$$I(X; Y) = \begin{cases} \tilde{I}(\mathbf{X} \rightarrow \mathbf{Y}) & \text{if } X \prec Y \\ \tilde{I}(\mathbf{Y} \rightarrow \mathbf{X}) & \text{if } Y \prec X \\ I_{res}(\mathbf{X}; \mathbf{Y}) & \text{if } X \approx Y \end{cases}.$$

With this definition of directed information we have that there may be an information flow from  $X$  to  $Y$  although  $\neg X \ll Y$ . In this sense information may not flow along the direction of causation. We shall discuss this problem later and suggest a solution.

#### IV. RELATIVITY OF INFORMATION FLOWS

Both the definition of directed information and residual information depends on the synchronization of the sequences. We shall now see how the definitions rely on the synchronization. We have

$$\begin{aligned} \tilde{I}(T(\mathbf{X}) \rightarrow \mathbf{Y}) &= H(Y_1 | Y_{-\infty}^0) - H(Y_1 | T(X)_{-\infty}^0 \cup Y_{\infty}^0) \\ &= H(Y_1 | Y_{-\infty}^0) - H(Y_1 | X_{-\infty}^{-1} \cup Y_{\infty}^0) \\ &\leq H(Y_1 | Y_{-\infty}^0) - H(Y_1 | X_{-\infty}^0 \cup Y_{\infty}^0) \\ &= \tilde{I}(\mathbf{X} \rightarrow \mathbf{Y}). \end{aligned}$$

The identity

$$\tilde{I}(\mathbf{X} \rightarrow T(\mathbf{Y})) = \tilde{I}(T^{-1}(\mathbf{X}) \rightarrow \mathbf{Y})$$

is obvious. Thus  $\tilde{I}(T^n(\mathbf{X}) \rightarrow \mathbf{Y})$  is a decreasing function of  $n$  and  $\tilde{I}(\mathbf{X} \rightarrow T^n(\mathbf{Y}))$  is an increasing function of  $n$ . Therefore the numbers for which both  $\tilde{I}(T^n(\mathbf{X}) \rightarrow \mathbf{Y})$  and  $\tilde{I}(\mathbf{X} \rightarrow T^n(\mathbf{Y}))$  are positive is an interval. Only in this interval communication is possible in both directions and in practice this idea is often used to check synchronization. In the literature directed information is defined by

$$I(\mathbf{X} \rightarrow \mathbf{Y}) = \tilde{I}(\mathbf{X} \rightarrow T^{-1}(\mathbf{Y})),$$

see [5] and [8]. With this definition mutual information will not split up into two flows and a residual term. Instead one has to involve the operator  $T$  to get formulas relating mutual information and directed information.

Let  $k$  be a positive integer and let  $\mathbf{X}$  be a stationary sequence. Then  $k\mathbf{X}$  shall denote "block sequence" where the  $n$ 'th variable is block  $X_{(n-1)k+1}, X_{(n-1)k+2}, \dots, X_{nk}$ . We get

$$\begin{aligned} \tilde{I}(k\mathbf{X} \rightarrow k\mathbf{Y}) &= H(Y_1^k | Y_{-\infty}^0) - H(Y_1^k | X_{-\infty}^0 \cup Y_{\infty}^0) \\ &= \sum_{n=1}^k H(Y_n | Y_{-\infty}^{n-1}) - H(Y_n | X_{-\infty}^0 \cup Y_{\infty}^{n-1}) \\ &\leq \sum_{n=1}^k H(Y_n | Y_{-\infty}^{n-1}) - H(Y_n | X_{-\infty}^{n-1} \cup Y_{\infty}^{n-1}) \\ &= k\tilde{I}(\mathbf{X} \rightarrow \mathbf{Y}). \end{aligned}$$

Similarly we get

$$\tilde{I}_{res}(k\mathbf{X}; k\mathbf{Y}) \geq kI_{res}(\mathbf{X}; \mathbf{Y}).$$

The conclusion is that one can explain less mutual information as flow and more as residual information if one considers a block sequence instead of the original sequence.

#### V. CONDITIONAL FLOWS

Conditional versions of directed information have been studied in [4], [5], [8] and [6]. Note that the definition we shall use is a little different from what is found in the literature. We define the *conditional directed information* by

$$\tilde{I}(X^N \rightarrow Y^N \| Z^N) = \tilde{I}(X^N \cup Z^N \rightarrow Y^N) - \tilde{I}(Z^N \rightarrow Y^N) \quad (6)$$

so that the information flowing from  $X^N$  to  $Y^N$  given  $Z^N$  is the information flowing from both variables minus the amount of information flowing  $Z^N$ . Thus

$$\begin{aligned} \tilde{I}(X^N \rightarrow Y^N \| Z^N) &= \frac{1}{N} \sum_{n=1}^N \left( H(Y_n | Y_1^{n-1}) - H(Y_n | X_1^{n-1} \cup Y_1^{n-1} \cup Z_1^{n-1}) \right) \\ &\quad - \frac{1}{N} \sum_{n=1}^N \left( H(Y_n | Y_1^{n-1}) - H(Y_n | X_1^{n-1} \cup Z_1^{n-1}) \right) \\ &= \frac{1}{N} \sum_{n=1}^N \left( H(Y_n | Y_1^{n-1} \cup Z_1^{n-1}) - H(Y_n | X_1^{n-1} \cup Y_1^{n-1} \cup Z_1^{n-1}) \right). \end{aligned}$$

Note that  $\tilde{I}(X^N \rightarrow Y^N \| Z^N) \geq 0$ .

For stationary sequences we define

$$\tilde{I}(\mathbf{X} \rightarrow \mathbf{Y} \| \mathbf{Z}) = \lim_{N \rightarrow \infty} \tilde{I}(X^N \rightarrow Y^N \| Z^N).$$

Then

$$\begin{aligned} \tilde{I}(\mathbf{X} \rightarrow \mathbf{Y} \| \mathbf{Z}) &= H(Y_1 | Y_{-\infty}^0 \cup Z_{-\infty}^0) \\ &\quad - H(Y_1 | X_{-\infty}^0 \cup Y_{-\infty}^0 \cup Z_{-\infty}^0). \end{aligned}$$

We observe that

$$\begin{aligned} \tilde{I}(\mathbf{X} \cup \mathbf{Y} \rightarrow \mathbf{Z} \| \mathbf{W}) &= H(Z_1 | Z_{-\infty}^0 \cup W_{-\infty}^0) \\ &\quad - H(Z_1 | X_{-\infty}^0 \cup Y_{-\infty}^0 \cup Z_{-\infty}^0 \cup W_{-\infty}^0) \\ &= H(Z_1 | Z_{-\infty}^0 \cup W_{-\infty}^0) - H(Z_1 | X_{-\infty}^0 \cup Z_{-\infty}^0 \cup W_{-\infty}^0) \\ &\quad + H(Z_1 | X_{-\infty}^0 \cup Z_{-\infty}^0 \cup W_{-\infty}^0) \\ &\quad - H(Z_1 | X_{-\infty}^0 \cup Y_{-\infty}^0 \cup Z_{-\infty}^0 \cup W_{-\infty}^0) \\ &= \tilde{I}(\mathbf{X} \rightarrow \mathbf{Z} \| \mathbf{W}) + \tilde{I}(\mathbf{Y} \rightarrow \mathbf{Z} \| \mathbf{X} \cup \mathbf{W}). \end{aligned}$$

We are now able to write a conditional mutual information in terms of conditional flows and residual information. We have

$$\begin{aligned} I(\mathbf{X}; \mathbf{Y} | \mathbf{Z}) &= \lim_{N \rightarrow \infty} \frac{1}{N} I(X^N; Y^N | Z^N) \\ &= \lim_{N \rightarrow \infty} \left( \frac{1}{N} I(X^N \cup Z^N; Y^N) \right. \\ &\quad \left. - \frac{1}{N} I(Z^N; Y^N) \right) \\ &= I(\mathbf{X} \cup \mathbf{Z}; \mathbf{Y}) - I(\mathbf{Z}; \mathbf{Y}). \end{aligned}$$

Each of these terms can be written as a sum of two flow and residual information. Thus,

$$\begin{aligned} I(\mathbf{X}; \mathbf{Y} | \mathbf{Z}) &= \\ &\tilde{I}(\mathbf{X} \cup \mathbf{Z} \rightarrow \mathbf{Y}) + \tilde{I}(\mathbf{Y} \rightarrow \mathbf{X} \cup \mathbf{Z}) + \tilde{I}_{res}(\mathbf{X} \cup \mathbf{Z}; \mathbf{Y}) \\ &\quad - \tilde{I}(\mathbf{Z} \rightarrow \mathbf{Y}) - \tilde{I}(\mathbf{Y} \rightarrow \mathbf{Z}) - \tilde{I}_{res}(\mathbf{Z}; \mathbf{Y}). \end{aligned}$$

and we get the formula

$$\begin{aligned} I(\mathbf{X}; \mathbf{Y} | \mathbf{Z}) &= \\ &\tilde{I}(\mathbf{X} \rightarrow \mathbf{Y} \| \mathbf{Z}) + \tilde{I}(\mathbf{Y} \rightarrow \mathbf{X} \cup \mathbf{Z}) - \tilde{I}(\mathbf{Y} \rightarrow \mathbf{Z}) \\ &\quad + \tilde{I}_{res}(\mathbf{X} \cup \mathbf{Z}; \mathbf{Y}) - \tilde{I}_{res}(\mathbf{Z}; \mathbf{Y}). \end{aligned}$$

The formula is closely related to [8, Prop. 3], but by introducing residual information we have formulas that do not involve the delay operator. It seems tempting to define  $\tilde{I}_{res}(\mathbf{X} \cup \mathbf{Z}; \mathbf{Y}) - \tilde{I}_{res}(\mathbf{Z}; \mathbf{Y})$  as the residual mutual information of  $X$  and  $Y$  given  $Z$ , but it may lead to a negative quantity.

The basic equations for conditional mutual information can now be stated in terms of directed and residual information. Equation 3a states that

$$\begin{aligned} &\tilde{I}(\mathbf{X} \rightarrow \mathbf{Y} \| \mathbf{Z}) + \tilde{I}(\mathbf{Y} \rightarrow \mathbf{X} \cup \mathbf{Z}) - \tilde{I}(\mathbf{Y} \rightarrow \mathbf{Z}) \\ &\quad + \tilde{I}_{res}(\mathbf{X} \cup \mathbf{Z}; \mathbf{Y}) - \tilde{I}_{res}(\mathbf{Z}; \mathbf{Y}) \\ &= \tilde{I}(\mathbf{Y} \rightarrow \mathbf{X} \| \mathbf{Z}) + \tilde{I}(\mathbf{X} \rightarrow \mathbf{Y} \cup \mathbf{Z}) - \tilde{I}(\mathbf{X} \rightarrow \mathbf{Z}) \\ &\quad + \tilde{I}_{res}(\mathbf{Y} \cup \mathbf{Z}; \mathbf{X}) - \tilde{I}_{res}(\mathbf{Z}; \mathbf{X}) \end{aligned}$$

The terms can be reorganized so that we get

$$\begin{aligned} &\tilde{I}(\mathbf{X} \rightarrow \mathbf{Y} \cup \mathbf{Z}) - \tilde{I}(\mathbf{X} \rightarrow \mathbf{Y} \| \mathbf{Z}) - \tilde{I}(\mathbf{X} \rightarrow \mathbf{Z}) \quad (7) \\ &\quad + \tilde{I}_{res}(\mathbf{X}; \mathbf{Y} \cup \mathbf{Z}) - \tilde{I}_{res}(\mathbf{Z}; \mathbf{X}) \\ &= \tilde{I}(\mathbf{Y} \rightarrow \mathbf{X} \cup \mathbf{Z}) - \tilde{I}(\mathbf{Y} \rightarrow \mathbf{X} \| \mathbf{Z}) - \tilde{I}(\mathbf{Y} \rightarrow \mathbf{Z}) \\ &\quad + \tilde{I}_{res}(\mathbf{Y}; \mathbf{X} \cup \mathbf{Z}) - \tilde{I}_{res}(\mathbf{Z}; \mathbf{Y}). \end{aligned}$$

We see that relation 7 together with Equation 5 and the definitions 6 and 4 implies the identities 3a and 3b.

## VI. LOWER BOUNDS ON DIRECTED INFORMATION

As we have noticed that directed information is relative to the synchronization of the sequences. Nevertheless it is sometimes possible to provide lower bound to the directed information, which are independent of the synchronization. Here we shall just provide a simple example that will illustrate the idea.

Consider (sets of) variables  $X, Y, Z$  and  $W$  in a Bayesian network satisfying that statistical independence holds if and only the  $d$ -separation criteria is fulfilled. Assume that  $X$  and  $Y$  are independent given  $Z \cup W$ . Then we have the lower bound

$$I(W \rightarrow X \cup Y \cup Z) \geq I(X; Y | Z). \quad (8)$$

To see this we put  $W' = W \cap a(X \cup Y \cup Z)$  and note  $X$  and  $Y$  are independent of  $Z \cup W'$  [9]. Now,

$$\begin{aligned} \tilde{I}(W \rightarrow X \cup Y \cup Z) &= \tilde{I}(W' \rightarrow X \cup Y \cup Z) \\ &\quad + \tilde{I}(\mathbb{C}W \rightarrow X \cup Y \cup Z | W') \\ &\geq \tilde{I}(W' \rightarrow X \cup Y \cup Z) \\ &\geq I(W; X \cup Y \cup Z) \\ &\geq I(W; X \cup Y | Z) \\ &\geq I(X; Y | Z). \end{aligned}$$

The last inequality is easily proved by using the Venn diagram method [10].

The bound (8) can be written as

$$\tilde{I}(W \rightarrow X \cup Y \cup Z) \geq I(X; Y | Z) - I(X; Y | Z \cup W).$$

Under weak conditions this lower bound holds even when  $X$  and  $Y$  are not independent given  $Z \cup W$ . More refined lower bounds are an area for future investigations.

## VII. ACKNOWLEDGEMENT

The author want to thank Jim Massey for useful discussions.

## REFERENCES

- [1] C. E. Shannon, "A mathematical theory of communication," *Bell Syst. Tech. J.*, vol. 27, pp. 379–423 and 623–656, 1948.
- [2] H. Marko, "The bidirectional communication theory - a generalization of information theory," *IEEE Trans. Commun.*, vol. COM-21, pp. 1345–1351, Dec. 1973.
- [3] J. L. Massey, "Causality, feedback and directed information," in *Proc. 1990 Int. Symp. on Info. Th. and its Appls., Hawaii, USA.*, pp. 303–305, 1990.
- [4] G. Kramer, *Directed Information for Channels with Feedback*, vol. 11 of *ETH Series in Inform. Proc.* Konstanz: Hartung-Gorre, 1998.
- [5] G. Kramer, "Capacity results for the discrete memoryless network," *IEEE Trans. Inform. Theory*, vol. IT-49, pp. 4–21, Jan. 2003.

- [6] R. Venkataramanan and S. S. Pradhan, "Directed information for communication problems with common side information and delayed feedback/feedforward," in *Proc. Allerton conference on communication, control and computing*, Sept. 2005.
- [7] J. Pearl, *Probabilistic Reasoning in Intelligent Systems*. San Mateo, California: Morgan Kaufmann Publ., 1988.
- [8] J. L. Massey and P. C. Massey, "Conservation of mutual and directed information," in *Proceedings of the 2005 IEEE International Symposium on Information Theory*, (Adelaide, South Australia, Australia), pp. 157–158, Sept. 2005.
- [9] P. Harremoës, *Time and Conditional Independence*, vol. 255 of *IMFUFA-tekst*. IMFUFA Roskilde University, 1993. Original in Danish entitled Tid og Betinget Uafhængighed. English translation partially available at <http://www.math.ku.dk/moes/index.html>.
- [10] R. W. Yeung, *A First Course in Information Theory*. New York: Kluwer, 2002.