

ENEE 627
SPRING 2011
INFORMATION THEORY
DATA PROCESSING

Markov chains

Consider a collection of n rvs, say X_1, \dots, X_n , defined on the same probability triple. For each $i = 1, \dots, n$, the rv X_i is \mathcal{X}_i -valued with \mathcal{X}_i a finite set. We shall write

$$\mathcal{X}^n = \times_{i=1}^n \mathcal{X}_i.$$

The rvs X_1, \dots, X_n are said to form a *Markov chain* if the conditions

$$(1) \quad \begin{aligned} & \mathbb{P}[X_1 = x_1, X_2 = x_2, \dots, X_n = x_n] \\ &= \mathbb{P}[X_1 = x_1] \prod_{k=1}^{n-1} p_{k+1}(x_{k+1}|x_k), \quad (x_1, \dots, x_n) \in \mathcal{X}^n \end{aligned}$$

all hold where for each $k = 1, \dots, n - 1$, we require

$$\begin{aligned} 0 \leq p_{k+1}(x_{k+1}|x_k) \leq 1 \\ \sum_{x_{k+1} \in \mathcal{X}_{k+1}} p_{k+1}(x_{k+1}|x_k) = 1 \end{aligned}, \quad x_k \in \mathcal{X}_k, x_{k+1} \in \mathcal{X}_{k+1}.$$

The Markov chain property of the rvs X_1, \dots, X_n is concisely represented through

$$X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_{n-1} \rightarrow X_n.$$

Simple facts

Here are some simple facts concerning Markov chains as needed in the context of Information Theory.

Fact 0.1 *If $X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_{n-1} \rightarrow X_n$, then for each $k = 2, \dots, n - 1$, the rvs $\{X_i, i = 1, 2, \dots, k - 1\}$ and $\{X_j, j = k + 1, k + 2, \dots, n\}$ are mutually independent given X_k .*

The Markov property is inherited by taking subsets.

Fact 0.2 *If $X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_{n-1} \rightarrow X_n$, then for any subset I of $\{1, \dots, n\}$ with $|I| \geq 2$, the collection of rvs $\{X_i, i \in I\}$ is also a Markov chain, namely if $I = \{i_1, \dots, i_k\}$ with $i_1 < i_2 < \dots < i_k$ for some $k = 2, \dots, n$, then*

$$X_{i_1} \rightarrow X_{i_2} \rightarrow \dots \rightarrow X_{i_k}.$$

Proof. Note that if $X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_{n-1} \rightarrow X_n$, then

$$\begin{aligned} & \mathbb{P}[X_1 = x_1, X_2 = x_2, \dots, X_{n-1} = x_{n-1}] \\ &= \sum_{x_n \in \mathcal{X}_n} \mathbb{P}[X_1 = x_1, X_2 = x_2, \dots, X_n = x_n] \\ &= \sum_{x_n \in \mathcal{X}_n} \mathbb{P}[X_1 = x_1] \prod_{k=1}^{n-1} p_{k+1}(x_{k+1}|x_k) \\ &= \sum_{x_n \in \mathcal{X}_n} \mathbb{P}[X_1 = x_1] \prod_{k=1}^{n-2} p_{k+1}(x_{k+1}|x_k) p_n(x_n|x_{n-1}) \\ &= \mathbb{P}[X_1 = x_1] \prod_{k=1}^{n-2} p_{k+1}(x_{k+1}|x_k) \left(\sum_{x_n \in \mathcal{X}_n} p_n(x_n|x_{n-1}) \right) \\ (2) \quad &= \mathbb{P}[X_1 = x_1] \prod_{k=1}^{n-2} p_{k+1}(x_{k+1}|x_k), \quad (x_1, \dots, x_{n-1}) \in \mathcal{X}^{n-1} \end{aligned}$$

since

$$\sum_{x_n \in \mathcal{X}_n} p_n(x_n|x_{n-1}), \quad x_{n-1} \in \mathcal{X}_{n-1},$$

and it is now plain that $X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_{n-1}$.

Similarly, if $X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_{n-1} \rightarrow X_n$, then

$$\begin{aligned} & \mathbb{P}[X_2 = x_2, \dots, X_{n-1} = x_{n-1}, X_n = x_n] \\ &= \sum_{x_1 \in \mathcal{X}_1} \mathbb{P}[X_1 = x_1, X_2 = x_2, \dots, X_{n-1} = x_{n-1}, X_n = x_n] \\ &= \sum_{x_1 \in \mathcal{X}_1} \mathbb{P}[X_1 = x_1] \prod_{k=1}^{n-1} p_{k+1}(x_{k+1}|x_k) \end{aligned}$$

$$\begin{aligned}
&= \left(\sum_{x_1 \in \mathcal{X}_1} \mathbb{P}[X_1 = x_1] p_2(x_2|x_1) \right) \prod_{k=1}^{n-1} p_{k+1}(x_{k+1}|x_k) \\
(3) \quad &= \mathbb{P}[X_2 = x_2] \prod_{k=2}^{n-1} p_{k+1}(x_{k+1}|x_k), \quad (x_2, \dots, x_n) \in \mathcal{X}_2 \times \dots \times \mathcal{X}_n
\end{aligned}$$

as we note that

$$\mathbb{P}[X_2 = x_2] = \sum_{x_1 \in \mathcal{X}_1} \mathbb{P}[X_1 = x_1] p_2(x_2|x_1).$$

Just apply (1) and use the conditions

$$\sum_{x_{k+1} \in \mathcal{X}_{k+1}} p_{k+1}(x_{k+1}|x_k) = 1, \quad \begin{array}{l} k = 1, \dots, n-1 \\ x_k \in \mathcal{X}_k \end{array}$$

and we get $X_2 \rightarrow X_2 \rightarrow \dots \rightarrow X_{n-1} \rightarrow X_n$.

Finally, assuming $n \geq 3$, pick $2 \leq \ell \leq n-1$. Similar arguments show that removing X_ℓ does not change the Markov property of the remaining rvs. Thus removing any one of the rvs does not change the Markov property. Iterating this operation k times with the rvs with index in I yields the result. ■

Reversing time does not change the Markov property.

Fact 0.3 *If $X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_{n-1} \rightarrow X_n$, then it is also the case that $X_n \rightarrow X_{n-1} \rightarrow \dots \rightarrow X_2 \rightarrow X_1$.*

Proof. From (1) we note that

$$\begin{aligned}
&\mathbb{P}[X_n = y_1, X_{n-1} = y_2, \dots, X_1 = y_n] \\
&= \mathbb{P}[X_1 = y_n, X_2 = y_{n-1}, \dots, X_{n-1} = y_2, X_n = y_1] \\
(4) \quad &= \mathbb{P}[X_1 = y_n] \prod_{k=1}^{n-1} p_{k+1}(y_k|y_{k-1}), \quad (y_n, \dots, y_1) \in \mathcal{X}^n
\end{aligned}$$

■

Data Processing Inequalities

We begin with the Data Processing Inequality in its standard form.

Lemma 0.1 For any Markov chain $X \rightarrow Y \rightarrow Z$, it is the case that

$$(5) \quad I(X; Z) \leq I(X; Y)$$

and

$$(6) \quad I(X; Z) \leq I(Z; Y).$$

Proof. By the chain rule for mutual informations applied to $I(X; (Y, Z))$ twice, we find

$$(7) \quad I(X; (Y, Z)) = I(X; Y) + I(X; Z|Y)$$

and

$$(8) \quad I(X; (Y, Z)) = I(X; Z) + I(X; Y|X)$$

The Markov property $X \rightarrow Y \rightarrow Z$ implies that X and Z are conditionally independent given Y , whence $I(X; Z|Y) = 0$. Thus,

$$(9) \quad I(X; (Y, Z)) = I(X; Y)$$

and we conclude that

$$(10) \quad I(X; Z) + I(X; Y|X) = I(X; Y).$$

The desired conclusion (5) now follows since $I(X; Y|X) \geq 0$.

By Fact 0.3 we note that $Z \rightarrow Y \rightarrow X$ since $X \rightarrow Y \rightarrow Z$, and applying (5) (this time with $X \leftarrow Z$, $Y \leftarrow Y$ and $Z \leftarrow X$) yields (6). ■

In the context of the Channel Coding Theorem, the following version of the Data Processing Inequality is needed.

Lemma 0.2 For any Markov chain $X \rightarrow U \rightarrow V \rightarrow Y$, we have

$$(11) \quad I(X; Y) \leq I(U; V).$$

Proof. By Fact 0.2 the Markov property $X \rightarrow U \rightarrow V \rightarrow Y$ implies both

$$(12) \quad X \rightarrow V \rightarrow Y$$

and

$$(13) \quad X \rightarrow U \rightarrow V.$$

Applying Lemma 0.1 to (12) and (13) we get

$$(14) \quad I(X; Y) \leq I(X; V)$$

and

$$(15) \quad I(X; V) \leq I(U; V).$$

The conclusion (11) follows by combining (14) and (15). ■

The Markov property and the DMC

Consider the DMC with channel matrix $\mathbf{P} = (p(y|x), x \in \mathcal{X}, y \in \mathcal{Y})$. The message W to be sent is selected from a set of M distinct messages $\mathcal{M} \equiv \{1, 2, \dots, M\}$ with M some positive integer. For each $n = 1, 2, \dots$, consider the (M, n) -code $C_n = (f_n, g_n)$ with encoding function $f_n : \mathcal{M} \rightarrow \mathcal{X}^n$ and decoding function $g_n : \mathcal{Y}^n \rightarrow \mathcal{M}$.

The string of symbols \mathbf{X}^n to be sent over the channel is specified by

$$\mathbf{X}^n = f_n(W),$$

and upon receiving the string of symbols \mathbf{Y}^n , the estimate \widehat{W} is generated according to

$$\widehat{W} = g_n(\mathbf{Y}^n).$$

As usual the DMC assumption is encapsulated through

$$\mathbb{P}[\mathbf{Y}^n = \mathbf{y}^n | \mathbf{X}^n = \mathbf{x}^n] = \prod_{k=1}^n p(y_k | x_k), \quad \mathbf{x}^n \in \mathcal{X}^n, \mathbf{y}^n \in \mathcal{Y}^n.$$

Lemma 0.3 *With the usual notation, for each $n = 1, 2, \dots$, we have*

$$(16) \quad W \rightarrow \mathbf{X}^n \rightarrow \mathbf{Y}^n \rightarrow \widehat{W}$$

provided W is selected independently of the operation of the DMC.

Proof. Select w in \mathcal{M} , \mathbf{x}^n in \mathcal{X}^n , \mathbf{y}^n in \mathcal{Y}^n and v in \mathcal{M} . Note that

$$\begin{aligned}
& \mathbb{P} \left[W = w, \mathbf{X}^n = \mathbf{x}^n, \mathbf{Y}^n = \mathbf{y}^n, \widehat{W} = v \right] \\
&= \mathbb{P} [W = w, \mathbf{X}^n = \mathbf{x}^n, \mathbf{Y}^n = \mathbf{y}^n, g_n(\mathbf{Y}^n) = v] \\
&= \mathbb{P} [W = w, \mathbf{X}^n = \mathbf{x}^n, \mathbf{Y}^n = \mathbf{y}^n, g_n(\mathbf{y}^n) = v] \\
&= \mathbb{P} [W = w, \mathbf{X}^n = \mathbf{x}^n, \mathbf{Y}^n = \mathbf{y}^n] \cdot \delta(g_n(\mathbf{y}^n), v) \\
&= \mathbb{P} [W = w, \mathbf{X}^n = \mathbf{x}^n] \cdot \prod_{k=1}^n p(y_k|x_k) \cdot \delta(g_n(\mathbf{y}^n), v) \\
&= \mathbb{P} [W = w, f_n(W) = \mathbf{x}^n] \cdot \prod_{k=1}^n p(y_k|x_k) \cdot \delta(g_n(\mathbf{y}^n), v) \\
(17) \quad &= \mathbb{P} [W = w] \cdot \delta(f_n(w), \mathbf{x}^n) \cdot \prod_{k=1}^n p(y_k|x_k) \cdot \delta(g_n(\mathbf{y}^n), v).
\end{aligned}$$

■

Basic arguments in the converse of the CCT

The converse to the Channel Coding Theorem results from the following chain of arguments:

Assume W to be uniformly distributed over the message set \mathcal{M} , and independent of the operation of the DMC. Thus,

$$\begin{aligned}
& \log_2 M \\
&= H(W) \\
&= H(W|\widehat{W}) + I(W; \widehat{W}) \\
&\leq 1 + \log_2 M \cdot \mathbb{P} [\widehat{W} \neq W] + I(W; \widehat{W}) \quad (\text{Fano's Inequality}) \\
&\leq 1 + \log_2 M \cdot \mathbb{P} [\widehat{W} \neq W] + I(\mathbf{X}^n; \mathbf{Y}^n) \quad (\text{Data Processing}) \\
&= 1 + \log_2 M \cdot \mathbb{P} [\widehat{W} \neq W] + H(\mathbf{Y}^n) - H(\mathbf{Y}^n|\mathbf{X}^n) \\
&= 1 + \log_2 M \cdot \mathbb{P} [\widehat{W} \neq W] + H(\mathbf{Y}^n) - \sum_{i=1}^n H(Y_i|\mathbf{X}^n, \mathbf{Y}^{i-1}) \\
&\quad (\text{Chain rule for conditional entropy})
\end{aligned}$$

$$\begin{aligned} &= 1 + \log_2 M \cdot \mathbb{P} \left[\widehat{W} \neq W \right] + H(\mathbf{Y}^n) - \sum_{i=1}^n H(Y_i|X_i) \quad (\text{DMC}) \\ &\leq 1 + \log_2 M \cdot \mathbb{P} \left[\widehat{W} \neq W \right] + \sum_{i=1}^n H(Y_i) - \sum_{i=1}^n H(Y_i|X_i) \\ &= 1 + \log_2 M \cdot \mathbb{P} \left[\widehat{W} \neq W \right] + \sum_{i=1}^n I(X_i; Y_i) \\ &\leq 1 + \log_2 M \cdot \mathbb{P} \left[\widehat{W} \neq W \right] + nC \quad (\text{Definition of channel capacity for the DMC}) \end{aligned}$$

In short,

$$(18) \quad \frac{\log_2 M}{n} \leq \frac{1}{n} + \frac{\log_2 M}{n} \cdot \mathbb{P} \left[\widehat{W} \neq W \right] + C$$
