

Data Processing Theorems and the Second Law of Thermodynamics

Neri Merhav*

Department of Electrical Engineering
Technion - Israel Institute of Technology
Haifa 32000, ISRAEL

Abstract

We draw relationships between the generalized data processing theorems of Zakai and Ziv (1973 and 1975) and the dynamical version of the second law of thermodynamics, a.k.a. the Boltzmann H-Theorem, which asserts that the Shannon entropy, $H(X_t)$, pertaining to a finite-state Markov process $\{X_t\}$, is monotonically non-decreasing as a function of time t , provided that the steady-state distribution of this process is uniform across the state space (which is the case when the process designates an isolated system). It turns out that both the generalized data processing theorems and the Boltzmann H-Theorem can be viewed as special cases of a more general principle concerning the monotonicity (in time) of a certain generalized information measure applied to a Markov process. This gives rise to a new look at the generalized data processing theorem, which suggests to exploit certain degrees of freedom that may lead to better bounds, for a given choice of the convex function that defines the generalized mutual information.

Index Terms: Data processing inequality, convexity, perspective function, H-Theorem, thermodynamics, detailed balance.

*This work was supported by the Israel Science Foundation (ISF) grant no. 208/08.

1 Introduction

In [6], Csiszár considered a generalized notion of the divergence between two probability distributions, a.k.a. the *f*-divergence, by replacing the negative logarithm function, of the classical divergence,

$$D(P_1\|P_2) = \int dx \cdot P_1(x) \left[-\log \frac{P_2(x)}{P_1(x)} \right], \quad (1)$$

with a general convex function¹ Q , i.e.,

$$D_Q(P_1\|P_2) = \int dx \cdot P_1(x) \cdot Q \left(\frac{P_2(x)}{P_1(x)} \right). \quad (2)$$

When the *f*-divergence was applied to the joint distribution (in the role of P_1) and the product of marginals (in the role of P_2) of two random variables, it yielded a generalized notion of mutual information,

$$I^Q(X; Y) = \int dx dy \cdot P(x, y) \cdot Q \left(\frac{P(x)P(y)}{P(x, y)} \right) = \int dx dy \cdot P(x, y) \cdot Q \left(\frac{P(y)}{P(y|x)} \right), \quad (3)$$

which was shown in [6] to obey a data processing inequality, thus extending the well known data processing inequality of the ordinary mutual information (see, e.g., [5, Section 2.8]).

The same ideas were introduced independently by Ziv and Zakai [14], with the primary motivation of using it to obtain sharper distortion bounds for classes of simple codes for joint source-channel coding (e.g., of block length 1), as well as certain situations of signal detection and estimation (see also [1]). The idea was to define both a “rate-distortion function,” $R^Q(d)$ and a “channel capacity,” C^Q , by minimization and maximization (respectively) of the mutual information pertaining to Q , and to derive a lower bound on the distortion d from the data processing inequality

$$R^Q(d) \leq C^Q. \quad (4)$$

In the sequel, this will be referred to as the 1973 version of the generalized data processing theorem. In a somewhat less well known work [15], Zakai and Ziv have substantially further generalized their data processing theorems, so as to apply an even more general information measures, and this will be referred to as the 1975 version. This generalized information measure was in the form

$$I^Q(X; Y) = \int dx dy \cdot P(x, y) \cdot Q \left(\frac{\mu_1(x, y)}{P(x, y)}, \dots, \frac{\mu_k(x, y)}{P(x, y)} \right)$$

¹Originally, this function was denoted by f in [6], hence the name *f*-divergence.

$$= \int dx dy \cdot P(x, y) \cdot Q \left(\frac{\mu_1(y|x)}{P(y|x)}, \dots, \frac{\mu_k(y|x)}{P(y|x)} \right), \quad (5)$$

where Q is now an arbitrary convex function of k variables and $\{\mu_i(x, y)\}$ are arbitrary positive measures (not necessarily probability measures) that are defined consistently with the Markov conditions and where $\mu_i(y|x) = \mu_i(x, y)/P(x)$. It was shown in [15, Theorem 7.1] that the distortion bounds obtained from (5) are tight in the sense that there always exist a convex function Q and measures $\{\mu_i\}$ that would yield the exact distortion pertaining to the optimum communication system, and so, there is no room for improvement of this class of bounds.²

By setting $\mu_i(y|x) = P(y|x_i)$, $i = 1, 2, \dots, k-1$, where $\{x_i\}$ are $k-1$ particular letters in the alphabet of X , and $\mu_k(y|x) = P(y)$, they defined yet another generalized information measure that satisfies the data processing theorem as

$$\mathbf{E} \left\{ Q \left(\frac{P(Y|X_1)}{P(Y|X)}, \dots, \frac{P(Y|X_{k-1})}{P(Y|X)}, \frac{P(Y)}{P(Y|X)} \right) \right\}, \quad (6)$$

where the expectation is taken w.r.t. the joint distribution

$$P(x_1, \dots, x_{k-1}, x, y) = P(x)P(y|x)P(x_1)P(x_2) \cdots P(x_{k-1}).$$

In both [14] and [15], there are many examples how these data processing inequalities can be used to improve on earlier distortion bounds.

The data processing theorems of Csiszár and Zakai and Ziv form one aspect of this work. The other aspect, which may seem unrelated at first glance (but will nevertheless be shown here to be strongly related) is the second law of thermodynamics, or more precisely, *Boltzmann's H-theorem*. The second law of thermodynamics tells that in an isolated physical system (i.e., when no energy flows in or out), the entropy cannot decrease over time. Since one of the basic postulates of statistical physics tells that all states of the system, which have the same energy, also have the same probability in equilibrium, it follows that the stationary (equilibrium) distribution of these states must be uniform, because all accessible states must have the same energy when the system is isolated. Indeed, if the state of this system is designated by a Markov process, $\{X_t\}$ with a uniform stationary state distribution, the Boltzmann H-theorem tells that the Shannon entropy of X_t , $H(X_t)$, cannot decrease with t , which is a restatement of the second law.

²This result is non-constructive, however, in the sense that this choice of Q and $\{\mu_i\}$ depends on the optimum encoder and decoder.

We show, in this paper, that the generalized data processing theorems of [6], [14], and [15] on the one hand, and the Boltzmann H–theorem, on the other hand, are all special cases of a more general principle, which asserts that a certain generalized information measure, applied to the underlying Markov process must be a monotonic function of time. This unified framework provides a new perspective on the generalized data processing theorem. Beyond the fact that this new perspective may be interesting on its own right, it naturally suggests to exploit certain degrees of freedom of the Ziv–Zakai generalized mutual information that may lead to better bounds, for a given choice of the convex function that defines this generalized mutual information. These additional degrees of freedom may be important, because the variety of convex functions $\{Q\}$ which are convenient to work with, is rather limited. The fact that better bounds may indeed be obtained is demonstrated by an example.

The outline of the remaining part of this paper is as follows. In Section 2, we provide some background on Markov processes with a slight physical flavor, which will include the notion of detailed balance, global balance, as well as known results like the Boltzmann H–theorem, and its generalizations to information measures other than the entropy. In Section 3, we relate the generalized version of the Boltzmann H–theorem and the generalized data processing theorems and formalize the uniform framework that supports both. This is done, first for the 1973 version [14] of the Ziv–Zakai data processing theorem (along with an example), and then for the 1975 version by Zakai and Ziv [15]. Finally, in Section 4, we summarize and conclude.

2 Background

2.1 Detailed Balance and Global Balance

Many dynamical models of a physical system describe the microscopic state (or microstate, for short) of this system as a Markov process, $\{X_t\}$, either in discrete time or in continuous time. In this section, we discuss a few properties of these processes as well as the evolution of information measures associated with them, like entropy, divergence and more.

We begin with an isolated system in continuous time, which is not necessarily assumed to have reached yet its stationary distribution pertaining to equilibrium. Let us suppose that the state X_t

may take on values in a finite set \mathcal{X} . For $x, x' \in \mathcal{X}$, let us define the state transition rates

$$W_{xx'} = \lim_{\delta \rightarrow 0} \frac{\Pr\{X_{t+\delta} = x' | X_t = x\}}{\delta} \quad x' \neq x \quad (7)$$

which means, in other words,

$$\Pr\{X_{t+\delta} = x' | X_t = x\} = W_{xx'} \cdot \delta + o(\delta). \quad (8)$$

Denoting

$$P_t(x) = \Pr\{X_t = x\}, \quad (9)$$

it is easy to see that

$$P_{t+dt}(x) = \sum_{x' \neq x} P_t(x') W_{x'x} dt + P_t(x) \left(1 - \sum_{x' \neq x} W_{xx'} dt \right), \quad (10)$$

where the first sum describes the probabilities of all possible transitions from other states to state x and the second term describes the probability of not leaving state x . Subtracting $P_t(x)$ from both sides and dividing by dt , we immediately obtain the following set of differential equations:

$$\frac{dP_t(x)}{dt} = \sum_{x'} [P_t(x') W_{x'x} - P_t(x) W_{xx'}], \quad x \in \mathcal{X}, \quad (11)$$

where W_{xx} is defined in an arbitrary manner, e.g., $W_{xx} \equiv 0$ for all $x \in \mathcal{X}$. In the physics terminology (see, e.g., [10],[12]), these equations are called the *master equations*.³ When the process reaches stationarity, i.e., for all $x \in \mathcal{X}$, $P_t(x)$ converge to some $P(x)$ that is time-invariant, then

$$\sum_{x'} [P(x') W_{x'x} - P(x) W_{xx'}] = 0, \quad \forall x \in \mathcal{X}. \quad (12)$$

This situation is called *global balance* or *steady state*. When the physical system under discussion is isolated, namely, no energy flows into the system or out, the steady state distribution must be uniform across all states, because all accessible states must be of the same energy and the equilibrium probability of each state depends solely on its energy. Thus, in the case of an isolated system, $P(x) = 1/|\mathcal{X}|$ for all $x \in \mathcal{X}$. From quantum mechanical considerations, as well as considerations pertaining to time reversibility in the microscopic level,⁴ it is customary to assume $W_{xx'} = W_{x'x}$

³Note that the master equations apply in discrete time too, provided that the derivative at the l.h.s. is replaced by a simple difference, $P_{t+1}(x) - P_t(x)$, and $\{W_{xx'}\}$ are replaced one-step state transition probabilities.

⁴Consider, for example, an isolated system of moving particles of mass m and position vectors $\{\mathbf{r}_i(t)\}$, obeying the differential equations $m d^2 \mathbf{r}_i(t) / dt^2 = \sum_{j \neq i} F(\mathbf{r}_j(t) - \mathbf{r}_i(t))$, $i = 1, 2, \dots, n$, ($F(\mathbf{r}_j(t) - \mathbf{r}_i(t))$ being mutual interaction forces), which remain valid if the time variable t is replaced by $-t$ since $d^2 \mathbf{r}_i(t) / dt^2 = d^2 \mathbf{r}_i(-t) / d(-t)^2$.

for all pairs $\{x, x'\}$. We then observe that, not only do $\sum_{x'} [P(x')W_{x'x} - P(x)W_{xx'}]$ all vanish, but moreover, each individual term in this sum vanishes, as

$$P(x')W_{x'x} - P(x)W_{xx'} = \frac{1}{|\mathcal{X}|}(W_{x'x} - W_{xx'}) = 0. \quad (13)$$

This property is called *detailed balance*, which is stronger than global balance, and it means equilibrium, which is stronger than steady state. While both steady-state and equilibrium refer to situations of time-invariant state probabilities $\{P(x)\}$, a steady-state still allows cyclic “flows of probability.” For example, a Markov process with cyclic deterministic transitions $1 \rightarrow 2 \rightarrow 3 \rightarrow 1 \rightarrow 2 \rightarrow 3 \rightarrow \dots$ is in steady state provided that the probability distribution of the initial state is uniform $(1/3, 1/3, 1/3)$, however, the cyclic flow among the states is in one direction. On the other hand, in detailed balance ($W_{xx'} = W_{x'x}$ for an isolated system), which is equilibrium, there is no net flow in any cycle of states. All the net cyclic probability fluxes vanish, and therefore, time reversal would not change the probability law, that is, $\{X_{-t}\}$ has the same probability law as $\{X_t\}$ (see [9, Sect. 1.2]). For example, if $\{Y_t\}$ is a Bernoulli process, taking values equiprobably in $\{-1, +1\}$, then X_t defined recursively by

$$X_{t+1} = (X_t + Y_t) \bmod K, \quad (14)$$

has a symmetric state-transition probability matrix W , a uniform stationary state distribution, and it satisfies detailed balance.

2.2 Monotonicity of Information Measures

Returning to the case where the process $\{X_t\}$ pertaining to our isolated system has not necessarily reached equilibrium, let us take a look at the entropy of the state

$$H(X_t) = - \sum_{x \in \mathcal{X}} P_t(x) \log P_t(x). \quad (15)$$

The Boltzmann H-theorem (see, e.g., [3, Chap. 7], [8, Sect. 3.5], [10, pp. 171–173] [12, pp. 624–626]) asserts that $H(X_t)$ is monotonically non-decreasing. This result is a restatement of the second law of thermodynamics, which tells that the entropy of an isolated system cannot decrease with time. To see why this is true, we next show that detailed balance implies

$$\frac{dH(X_t)}{dt} \geq 0, \quad (16)$$

where for convenience, we denote $dP_t(x)/dt$ by $\dot{P}_t(x)$. Now,

$$\begin{aligned}
\frac{dH(X_t)}{dt} &= - \sum_x [\dot{P}_t(x) \log P_t(x) + \dot{P}_t(x)] \\
&= - \sum_x \dot{P}_t(x) \log P_t(x) \\
&= - \sum_x \sum_{x'} W_{x'x} [P_t(x') - P_t(x)] \log P_t(x) \\
&= - \frac{1}{2} \sum_{x,x'} W_{x'x} [P_t(x') - P_t(x)] \log P_t(x) - \\
&\quad \frac{1}{2} \sum_{x,x'} W_{x'x} [P_t(x) - P_t(x')] \log P_t(x') \\
&= \frac{1}{2} \sum_{x,x'} W_{x'x} [P_t(x') - P_t(x)] \cdot [\log P_t(x') - \log P_t(x)] \\
&\geq 0,
\end{aligned} \tag{17}$$

where in the second line we used the fact that $\sum_x \dot{P}_t(x) = 0$, in the third line we used detailed balance ($W_{xx'} = W_{x'x}$), and the last inequality is due to the increasing monotonicity of the logarithmic function: the product $[P_t(x') - P_t(x)] \cdot [\log P_t(x') - \log P_t(x)]$ cannot be negative for any pair (x, x') , as the two factors of this product are either both negative, both zero, or both positive. Thus, $H(X_t)$ cannot decrease with time.

The H-theorem has a discrete-time analogue: If a finite-state Markov process has a symmetric transition probability matrix (which is the discrete-time counterpart of the above detailed balance property), which means that the stationary state distribution is uniform, then $H(X_t)$ is a monotonically non-decreasing sequence.

A well-known paradox, in this context, is associated with the notion of the *arrow of time*. On the one hand, we are talking about time-reversible processes, obeying detailed balance, but on the other hand, the increase of entropy suggests that there is asymmetry between the two possible directions that the time axis can be exhausted, the forward direction and the backward direction. If we go back in time, the entropy would decrease. So is there an arrow of time? This paradox was resolved, by Boltzmann himself, once he made the clear distinction between equilibrium and non-equilibrium situations: The notion of time reversibility is associated with equilibrium, where the process $\{X_t\}$ is stationary. On the other hand, the increase of entropy is a result that belongs to the non-stationary regime, where the process is on its way to stationarity and equilibrium. In

the latter case, the system has been initially prepared in a non-equilibrium situation. Of course, when the process is stationary, $H(X_t)$ is fixed and there is no contradiction.

So far we discussed the property of detailed balance only for an isolated system, where the stationary state distribution is the uniform distribution. How is the property of detailed balance defined when the stationary distribution is non-uniform? For a general Markov process, whose steady state-distribution is not necessarily uniform, the condition of detailed balance, which means time-reversibility [9], reads

$$P(x)W_{xx'} = P(x')W_{x'x}, \quad (18)$$

in the continuous-time case. In the discrete-time case (where t takes on positive integer values only), it is defined by a similar equation, except that $W_{xx'}$ and $W_{x'x}$ are replaced by the corresponding one-step state transition probabilities, i.e.,

$$P(x)P(x'|x) = P(x')P(x|x'), \quad (19)$$

where

$$P(x'|x) \triangleq \Pr\{X_{t+1} = x' | X_t = x\}. \quad (20)$$

The physical interpretation is that now our system is (a small) part of a much larger isolated system, which obeys detailed balance w.r.t. the uniform equilibrium distribution, as before. A well known example of a process that obeys detailed balance in its more general form is the M/M/1 queue with an arrival rate λ and service rate μ ($\lambda < \mu$). Here, since all states are arranged along a line, with bidirectional transitions between neighboring states only (see Fig. 1), there cannot be any cyclic probability flux. The steady-state distribution is well-known to be geometric

$$P(x) = \left(1 - \frac{\lambda}{\mu}\right) \cdot \left(\frac{\lambda}{\mu}\right)^x, \quad x = 0, 1, 2, \dots, \quad (21)$$

which indeed satisfies the detailed balance $P(x)\lambda = P(x+1)\mu$ for all x . Thus, the Markov process $\{X_t\}$, designating the number of customers in the queue at time t , is time-reversible.

For the sake of simplicity, from this point onward, our discussion will focus almost exclusively on discrete-time Markov processes, but the results to be stated, will hold for continuous-time Markov processes as well. We will continue to denote by $P_t(x)$ the probability of $X_t = x$, except that now t will be limited to take on integer values only. The one-step state transition probabilities will be denoted by $\{P(x'|x)\}$, as mentioned earlier.

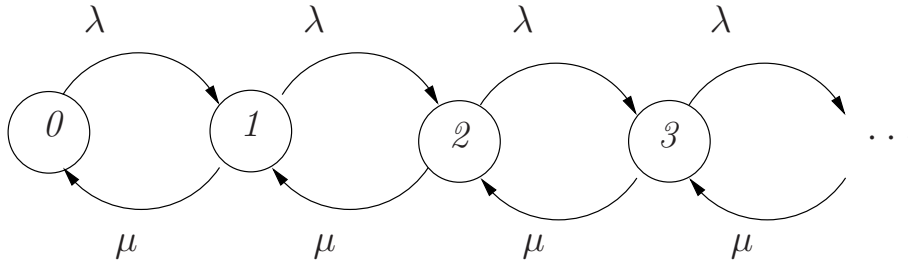


Figure 1: State transition diagram of an M/M/1 queue.

How does the H–theorem extend to situations where the stationary state distribution is not uniform? In [5, p. 82], it is shown (among other things) that the divergence,

$$D(P_t||P) = \sum_{x \in \mathcal{X}} P_t(x) \log \frac{P_t(x)}{P(x)}, \quad (22)$$

where $P = \{P(x), x \in \mathcal{X}\}$ is a stationary state distribution, is a monotonically non–increasing function of t . Does this result have a physical interpretation, like the H–theorem and the second law of thermodynamics? When it comes to non–isolated systems, where the steady state distribution is non–uniform, the extension of the second law of thermodynamics, replaces the principle of increase of entropy by the principle of decrease of free energy, or equivalently, the decrease of the difference between the free energy at time t and the free energy in equilibrium. The information–theoretic counterpart of this free energy difference is the divergence $D(P_t||P)$ (see, e.g., [2]). Thus, the monotonic decrease of $D(P_t||P)$ has a simple physical interpretation of free energy decrease, which is the natural extension of the entropy increase. Indeed, particularizing this to the case where P is the uniform distribution (as in an isolated system), then

$$D(P_t||P) = \log |\mathcal{X}| - H(X_t), \quad (23)$$

which means that the decrease of the divergence is equivalent to the increase of entropy, as before. However, here the result is more general than the H–theorem from an additional aspect: It does not require detailed balance. It only requires the existence of the stationary state distribution. Note that even in the earlier case, of an isolated system, detailed balance, which means symmetry of the state transition probability matrix ($P(x'|x) = P(x|x')$), is a stronger requirement than uniformity of the stationary state distribution, as the latter requires merely that the matrix $\{P(x'|x)\}$ would be doubly stochastic, i.e., $\sum_x P(x|x') = \sum_x P(x'|x) = 1$ for all $x' \in \mathcal{X}$, which is weaker than symmetry of the matrix itself. The results shown in [5] are, in fact, somewhat more general: Let $P_t = \{P_t(x)\}$

and $P'_t = \{P'_t(x)\}$ be two time-varying state-distributions pertaining to the same Markov chain, but induced by two different initial state distributions, $\{P_0(x)\}$ and $\{P'_0(x)\}$, respectively. Then $D(P_t||P'_t)$ is monotonically non-increasing. This is easily seen as follows:

$$\begin{aligned}
D(P_t||P'_t) &= \sum_x P_t(x) \log \frac{P_t(x)}{P'_t(x)} \\
&= \sum_{x,x'} P_t(x) P(x'|x) \log \frac{P_t(x) P(x'|x)}{P'_t(x) P(x'|x)} \\
&= \sum_{x,x'} P(X_t = x, X_{t+1} = x') \log \frac{P(X_t = x, X_{t+1} = x')}{P'(X_t = x, X_{t+1} = x')} \\
&\geq D(P_{t+1}||P'_{t+1})
\end{aligned} \tag{24}$$

where the last inequality follows from the data processing theorem of the divergence: the divergence between two joint distributions of (X_t, X_{t+1}) is never smaller than the divergence between corresponding marginal distributions of X_{t+1} . Another interesting special case of this result is obtained if we now take the first argument of the divergence to be a stationary state distribution: This will mean that $D(P||P_t)$ is also monotonically non-increasing.

In [9, Theorem 1.6], there is a further extension of all the above monotonicity results, where the ordinary divergence is actually replaced by the f-divergence (though the relation to the f-divergence is not mentioned in [9]): If $\{X_t\}$ is a Markov process with a given state transition probability matrix $\{P(x'|x)\}$, then the function

$$U(t) = D_Q(P||P_t) = \sum_{x \in \mathcal{X}} P(x) \cdot Q\left(\frac{P_t(x)}{P(x)}\right) \tag{25}$$

is monotonically non-increasing, provided that Q is convex. Moreover, $U(t)$ monotonically strictly decreasing if Q is strictly convex and $\{P_t(x)\}$ is not identical to $\{P(x)\}$. To see why this is true, define the backward transition probability matrix by

$$\tilde{P}(x|x') = \frac{P(x)P(x'|x)}{P(x')} \tag{26}$$

Obviously,

$$\sum_x \tilde{P}(x|x') = 1 \tag{27}$$

for all $x' \in \mathcal{X}$, and so,

$$\frac{P_{t+1}(x)}{P(x)} = \sum_{x'} \frac{P_t(x')P(x|x')}{P(x)} = \sum_{x'} \frac{\tilde{P}(x'|x)P_t(x')}{P(x')} \tag{28}$$

By the convexity of Q :

$$\begin{aligned}
U(t+1) &= \sum_x P(x) \cdot Q\left(\frac{P_{t+1}(x)}{P(x)}\right) \\
&= \sum_x P(x) \cdot Q\left(\sum_{x'} \tilde{P}(x'|x) \frac{P_t(x')}{P(x')}\right) \\
&\leq \sum_x \sum_{x'} P(x) \tilde{P}(x'|x) \cdot Q\left(\frac{P_t(x')}{P(x')}\right) \\
&= \sum_x \sum_{x'} P(x') P(x|x') \cdot Q\left(\frac{P_t(x')}{P(x')}\right) \\
&= \sum_{x'} P(x') \cdot Q\left(\frac{P_t(x')}{P(x')}\right) = U(t).
\end{aligned} \tag{29}$$

Now, a few interesting choices of the function Q may be considered: As proposed in [9, p. 19], for $Q(u) = u \ln u$, we have $U(t) = D(P_t \| P)$, and we are back to the aforementioned result in [5]. Another interesting choice is $Q(u) = -\ln u$, which gives $U(t) = D(P \| P_t)$. Thus, the monotonicity of $D(P \| P_t)$ is also obtained as a special case.⁵ Yet another choice is $Q(u) = -u^s$, where $s \in [0, 1]$ is a parameter. This would yield the increasing monotonicity of $\sum_x P^{1-s}(x) P_t^s(x)$, a ‘metric’ that plays a role in the theory of asymptotic exponents of error probabilities pertaining to the optimum likelihood ratio test between two probability distributions [13, Chapter 3]. In particular, the choice $s = 1/2$ yields balance between the two kinds of error and it is intimately related to the Bhattacharyya distance. In the case of detailed balance, there is another physical interpretation of the approach to equilibrium and the growth of $U(t)$ [9, p. 20]: Returning, for a moment, to the realm of continuous-time Markov processes, we can write the master equations as follows:

$$\frac{dP_t(x)}{dt} = \sum_{x'} \frac{1}{R_{xx'}} \left[\frac{P_t(x')}{P(x')} - \frac{P_t(x)}{P(x)} \right] \tag{30}$$

where $R_{xx'} = [P(x')W_{x'x}]^{-1} = [P(x)W_{xx'}]^{-1}$. Imagine now an electrical circuit where the indices $\{x\}$ designate the various nodes. Nodes x and x' are connected by a wire with resistance $R_{xx'}$ and every node x is grounded via a capacitor with capacitance $P(x)$ (see Fig. 2). If $P_t(x)$ is the charge at node x at time t , then the master equations are the Kirchoff equations of the currents at each node in the circuit. Thus, the way in which probability spreads across the states is analogous to

⁵We are not yet in a position to obtain the monotonicity of $D(P_t \| P'_t)$ as a special case of the monotonicity of $D_Q(P \| P_t)$. This will require a slight further extension of this information measure, to be carried out later on.

the way charge spreads across the circuit and probability fluxes are now analogous to electrical currents. If we now choose $Q(u) = \frac{1}{2}u^2$, then

$$U(t) = \frac{1}{2} \sum_x \frac{P_t^2(x)}{P(x)}, \quad (31)$$

which means that the energy stored in the capacitors dissipates as heat in the wires until the system reaches equilibrium, where all nodes have the same potential, $P_t(x)/P(x) = 1$, and hence detailed balance corresponds to the situation where all individual currents vanish (not only their algebraic sum).

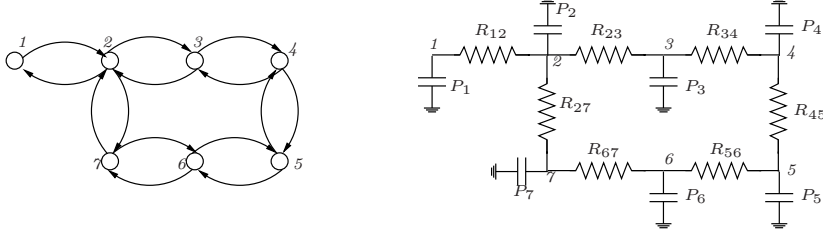


Figure 2: State transition diagram of a Markov chain (left part) and the electric circuit that emulates the dynamics of $\{P_t(x)\}$ (right part).

We have seen, in the above examples, that various choices of the function Q yield various f-divergences, or ‘metrics’, between $\{P(x)\}$ and $\{P_t(x)\}$, which are both marginal distributions of a single symbol x . What about joint distributions of two or more symbols? Consider, for example, the function

$$J(t) = \sum_{x,x'} P(X_0 = x, X_t = x') \cdot Q \left(\frac{P(X_0 = x)P(X_t = x')}{P(X_0 = x, X_t = x')} \right), \quad (32)$$

where Q is convex as before. Here, by the same token, $J(t)$ is the f-divergence between the joint probability distribution $\{P(X_0 = x, X_t = x')\}$ and the product of marginals $\{P(X_0 = x)P(X_t = x')\}$, namely, it is the generalized mutual information of [6],[14], and [15], as mentioned in the Introduction. Now, using a similar chain of inequalities as before, we get the non-decreasing monotonicity of $J(t)$ as follows:

$$\begin{aligned} J(t) &= \sum_{x,x',x''} P(X_0 = x, X_t = x', X_{t+1} = x'') \times \\ &Q \left(\frac{P(X_0 = x)P(X_t = x')}{P(X_0 = x, X_t = x')} \cdot \frac{P(X_{t+1} = x''|X_t = x')}{P(X_{t+1} = x''|X_t = x')} \right) \\ &= \sum_{x,x''} P(X_0 = x, X_{t+1} = x'') \sum_{x'} P(X_t = x'|X_0 = x, X_{t+1} = x'') \times \end{aligned}$$

$$\begin{aligned}
& Q \left(\frac{P(X_0 = x)P(X_t = x', X_{t+1} = x'')}{P(X_0 = x, X_t = x', X_{t+1} = x'')} \right) \\
& \leq \sum_{x, x''} P(X_0 = x, X_{t+1} = x'') \cdot Q \left(\sum_{x'} P(X_t = x' | X_0 = x, X_{t+1} = x'') \times \right. \\
& \quad \left. \frac{P(X_0 = x)P(X_t = x', X_{t+1} = x'')}{P(X_0 = x, X_t = x', X_{t+1} = x'')} \right) \\
& = \sum_{x, x''} P(X_0 = x, X_{t+1} = x'') Q \left(\sum_{x'} \frac{P(X_0 = x)P(X_t = x', X_{t+1} = x'')}{P(X_0 = x, X_{t+1} = x'')} \right) \\
& = \sum_{x, x''} P(X_0 = x, X_{t+1} = x'') \cdot Q \left(\frac{P(X_0 = x)P(X_{t+1} = x'')}{P(X_0 = x, X_{t+1} = x'')} \right) \\
& = J(t + 1). \tag{33}
\end{aligned}$$

This time, we assumed only the Markov property of (X_0, X_t, X_{t+1}) (not even homogeneity). This is, in fact, nothing but the 1973 version of the generalized data processing theorem of Ziv and Zakai [14], which was mentioned in the Introduction.

3 A Unified Framework

In spite of the general resemblance (via the notion of the f-divergence), the last monotonicity result, concerning $J(t)$, and the monotonicity of $D(P_t \| P'_t)$, do not seem, at first glance, to fall in the framework of the monotonicity of the f-divergence $D_Q(P \| P_t)$. This is because in the latter, there is an additional dependence on a stationary state distribution that appears neither in $D(P_t \| P'_t)$ nor in $J(t)$. However, two simple observations can put them both in the framework of the monotonicity of $D_Q(P \| P_t)$.

The first observation is that the monotonicity of $U(t) = D_Q(P \| P_t)$ continues to hold (with a straightforward extension of the proof) if $P_t(x)$ is extended to be a vector of time varying state distributions $(P_t^1(x), P_t^2(x), \dots, P_t^k(x))$, and Q is taken to be a convex function of k variables. Moreover, each component $P_t^i(x)$ does not have to be necessarily a probability distribution. It can be any function $\mu_t^i(x)$ that satisfies the recursion

$$\mu_{t+1}^i(x) = \sum_{x'} \mu_t^i(x') P(x|x'), \quad 1 \leq i \leq k. \tag{34}$$

Let us then denote $\boldsymbol{\mu}_t(x) = (\mu_t^1(x), \mu_t^2(x), \dots, \mu_t^k(x))$ and assume that Q is jointly convex in all its

k arguments. Then the redefined function

$$\begin{aligned} U(t) &= \sum_{x \in \mathcal{X}} P(x) \cdot Q\left(\frac{\boldsymbol{\mu}_t(x)}{P(x)}\right) \\ &= \sum_{x \in \mathcal{X}} P(x) \cdot Q\left(\frac{\mu_t^1(x)}{P(x)}, \dots, \frac{\mu_t^k(x)}{P(x)}\right) \end{aligned} \quad (35)$$

is monotonically non-increasing with t .

The second observation is rooted in convex analysis, and it is related to the notion of the perspective of a convex function and its convexity property [4]. Here, a few words of background are in order. Let $Q(\mathbf{u})$ be a convex function of the vector $\mathbf{u} = (u_1, \dots, u_k)$ and let $v > 0$ be an additional variable. Then, the function

$$\tilde{Q}(v, u_1, u_2, \dots, u_k) \triangleq v \cdot Q\left(\frac{u_1}{v}, \frac{u_2}{v}, \dots, \frac{u_k}{v}\right) \quad (36)$$

is called the *perspective function* of Q . A well-known property of the perspective operation is conservation of convexity, in other words, if Q is convex in \mathbf{u} , then \tilde{Q} is convex in (v, \mathbf{u}) . The proof of this fact, which is straightforward, can be found, for example, in [4, p. 89, Subsection 3.2.6] (see also [7]) and it is brought here for the sake of completeness: Letting λ_1 and λ_2 be two non-negative numbers summing to unity and letting (v_1, \mathbf{u}_1) and (v_2, \mathbf{u}_2) be given, then

$$\begin{aligned} \tilde{Q}(\lambda_1(v_1, \mathbf{u}_1) + \lambda_2(v_2, \mathbf{u}_2)) &= (\lambda_1 v_1 + \lambda_2 v_2) \cdot Q\left(\frac{\lambda_1 \mathbf{u}_1 + \lambda_2 \mathbf{u}_2}{\lambda_1 v_1 + \lambda_2 v_2}\right) \\ &= (\lambda_1 v_1 + \lambda_2 v_2) \cdot Q\left(\frac{\lambda_1 v_1}{\lambda_1 v_1 + \lambda_2 v_2} \cdot \frac{\mathbf{u}_1}{v_1} + \frac{\lambda_2 v_2}{\lambda_1 v_1 + \lambda_2 v_2} \cdot \frac{\mathbf{u}_2}{v_2}\right) \\ &\leq \lambda_1 v_1 Q\left(\frac{\mathbf{u}_1}{v_1}\right) + \lambda_2 v_2 Q\left(\frac{\mathbf{u}_2}{v_2}\right) \\ &= \lambda_1 \tilde{Q}(v_1, \mathbf{u}_1) + \lambda_2 \tilde{Q}(v_2, \mathbf{u}_2). \end{aligned} \quad (37)$$

Putting these two observations together, we can now state the following result:

Theorem 1 *Let*

$$V(t) = \sum_x \mu_t^0(x) Q\left(\frac{\mu_t^1(x)}{\mu_t^0(x)}, \frac{\mu_t^2(x)}{\mu_t^0(x)}, \dots, \frac{\mu_t^k(x)}{\mu_t^0(x)}\right), \quad (38)$$

where Q is a convex function of k variables and $\{\mu_t^i(x)\}_{i=0}^k$ are arbitrary functions that satisfy the recursion

$$\mu_{t+1}^i(x) = \sum_{x'} \mu_t^i(x') P(x|x'), \quad i = 0, 1, 2, \dots, k, \quad (39)$$

and where $\mu_t^0(x)$ is moreover strictly positive. Then, $V(t)$ is a monotonically non-increasing function of t .

Using the above mentioned observations, the proof of Theorem 1 is straightforward: Letting P be a stationary state distribution of $\{X_t\}$, we have:

$$\begin{aligned} V(t) &= \sum_x \mu_t^0(x) Q \left(\frac{\mu_t^1(x)}{\mu_t^0(x)}, \frac{\mu_t^2(x)}{\mu_t^0(x)}, \dots, \frac{\mu_t^k(x)}{\mu_t^0(x)} \right) \\ &= \sum_x P(x) \cdot \frac{\mu_t^0(x)}{P(x)} Q \left(\frac{\mu_t^1(x)/P(x)}{\mu_t^0(x)/P(x)}, \dots, \frac{\mu_t^k(x)/P(x)}{\mu_t^0(x)/P(x)} \right) \\ &= \sum_x P(x) \tilde{Q} \left(\frac{\mu_t^0(x)}{P(x)}, \frac{\mu_t^1(x)}{P(x)}, \dots, \frac{\mu_t^k(x)}{P(x)} \right). \end{aligned} \quad (40)$$

Since \tilde{Q} is the perspective of the convex function Q , then it is convex as well, and so, the monotonicity of $V(t)$ follows from the first observation above. It is now readily seen that both $D(P_t \| P'_t)$ and $J(t)$ are special cases of $V(t)$ and hence we have covered all special cases seen thus far under the umbrella of the more general information functional $V(t)$.

It is important to observe that the same idea exactly can be applied, first of all, to the 1973 version of the Ziv–Zakai data processing theorem (regardless of the above described monotonicity results concerning Markov processes): Consider the generalized mutual information functional

$$J^Q(X; Y) \triangleq \sum_{x,y} \mu_0(x, y) Q \left(\frac{\mu_1(x, y)}{\mu_0(x, y)} \right), \quad (41)$$

where $\mu_0(x, y) > 0$ and $\mu_1(x, y)$ are arbitrary functions that are consistent with the Markov conditions, i.e., for any Markov chain $X \rightarrow Y \rightarrow Z$, these functions satisfy

$$\mu_i(x, z) = \sum_y \mu_i(x, y) P(z|y) = \sum_y \mu_i(y, z) P(x|y), \quad i = 0, 1. \quad (42)$$

Then, $J^Q(X; Y)$ satisfies a data processing inequality, because, again

$$\begin{aligned} J^Q(X; Y) &= \sum_{x,y} P(x, y) \cdot \frac{\mu_0(x, y)}{P(x, y)} Q \left(\frac{\mu_1(x, y)/P(x, y)}{\mu_0(x, y)/P(x, y)} \right) \\ &= \sum_{x,y} P(x, y) \tilde{Q} \left(\frac{\mu_0(x, y)}{P(x, y)}, \frac{\mu_1(x, y)}{P(x, y)} \right), \end{aligned} \quad (43)$$

which is a Zakai–Ziv information functional of the 1975 version [15] and hence it satisfies a data processing inequality.

What functions, $\mu_0(x, y)$ and $\mu_1(x, y)$, can be consistent with the Markov conditions? Two such functions are, of course, $\mu_0(x, y) = P(x, y)$ and $\mu_1(x, y) = P(x)P(y)$, which bring us back to the 1973 Ziv–Zakai information measure. We can, of course, swap their roles and obtain a generalized version of the lautum information [11], which is also known to satisfy a data processing inequality. For additional options, let us consider a communication system, operating on single symbols (block length 1), where the source symbol u is mapped into a channel input $x = f(u)$, by a deterministic encoder f , which is then fed into the channel $P(y|x)$, and the channel output y is in turn mapped into the reconstruction symbol $v = g(y)$. As is argued in [15], the function $\mu(u, y) = P(u)P(y|u_0)$ is consistent with the Markov conditions for any given source symbol u_0 . Indeed, since the encoder is assumed deterministic, $P(y|u_0) = P(y|f(u_0)) = P(y|x_0)$, and it is easily seen that

$$\mu(u, v) = P(u)P(v|u_0) = \sum_y P(u)P(y|u_0)P(v|y) = \sum_y \mu(u, y)P(v|y) \quad (44)$$

and

$$\begin{aligned} \mu(u, y) &= P(u)P(y|u_0) \\ &= \sum_x P(u|x)P(x)P(y|u_0) \\ &= \sum_x P(u|x)P(x)P(y|x_0) = \sum_x P(u|x)\mu(x, y). \end{aligned} \quad (45)$$

Of course, every linear combination of all these functions is also consistent with the Markov conditions. Thus, we can take

$$\mu_0(x, y) = s_0P(x, y) + \sum_{x_i \in \mathcal{X}} s_i P(x)P(y|x_i) \quad (46)$$

and

$$\mu_1(x, y) = t_0P(x, y) + \sum_{x_i \in \mathcal{X}} t_i P(x)P(y|x_i), \quad (47)$$

where $\{s_i\}$ and $\{t_i\}$ are the (arbitrary) coefficients of these linear combinations (with the limitation that $s_i \geq 0$ for all i , with at least one $s_i > 0$). Thus, we may define

$$J^Q(X; Y) = \sum_{x, y} \left[s_0P(x, y) + \sum_{x_i \in \mathcal{X}} s_i P(x)P(y|x_i) \right] \cdot Q \left(\frac{t_0P(x, y) + \sum_{x_i \in \mathcal{X}} t_i P(x)P(y|x_i)}{s_0P(x, y) + \sum_{x_i \in \mathcal{X}} s_i P(x)P(y|x_i)} \right), \quad (48)$$

or, equivalently,

$$J^Q(X; Y) = \sum_{x, y} P(x) \left[s_0P(y|x) + \sum_{x_i \in \mathcal{X}} s_i P(y|x_i) \right] \cdot Q \left(\frac{t_0P(y|x) + \sum_{x_i \in \mathcal{X}} t_i P(y|x_i)}{s_0P(y|x) + \sum_{x_i \in \mathcal{X}} s_i P(y|x_i)} \right). \quad (49)$$

Moreover, to eliminate the dependence on the specific encoder, we can think of $\{x_i\}$ as independent random variables, take the expectation w.r.t. their randomness (in the same spirit as in [15]), and obtain the following information measure

$$\mathbf{E} \left\{ \sum_{x,y} P(x) \left[s_0 P(y|x) + \sum_i s_i P(y|X_i) \right] \cdot Q \left(\frac{t_0 P(y|x) + \sum_i t_i P(y|X_i)}{s_0 P(y|x) + \sum_i s_i P(y|X_i)} \right) \right\}, \quad (50)$$

where the expectation is w.r.t. the product measure of $\{X_i\}$, $P(x_1, x_2, \dots) = \prod_i P(x_i)$. These are the most general information measures, that obey a data processing inequality, that we can get with a univariate convex function Q . For example, returning to eq. (49) and taking $s_0 = 1$, $t_0 = 0$, $s_i = sP(x_i)$ ($s \geq 0$, a parameter), and $t_i = P(x_i)$, $x_i \in \mathcal{X}$, we have $\mu_0(x, y) = P(x, y) + sP(x)P(y)$, and $\mu_1(x, y) = P(x)P(y)$, and the resulting generalized mutual information reads

$$J^Q(X; Y) = \sum_{x,y} P(x)[P(y|x) + sP(y)] \cdot Q \left(\frac{P(y)}{P(y|x) + sP(y)} \right). \quad (51)$$

The interesting point concerning these generalized mutual information measures is that even if we remain in the framework of the 1973 version of the Ziv–Zakai data processing theorem (as opposed to the 1975 version), we have added an extra degrees of freedom (in the above example, the parameter s), which may be used in order to improve the obtained bounds. If the inequality $R^Q(d) \leq C^Q$ can be transformed into an inequality on the distortion d , where the lower bound depends on s , then this bound can be maximized w.r.t. the parameter s . If the optimum $s > 0$ yields a distortion bound which is larger than that of $s = 0$, then we have improved on [14] for the given choice of the convex function Q . Sometimes this optimization may not be a trivial task, but even if we can just identify one positive value of s (including the limit $s \rightarrow \infty$) that is better than $s = 0$, then we have improved on the generalized data processing bound of [14], which corresponds to $s = 0$. This additional degree of freedom may be important, because, as mentioned in the Introduction, the variety of convex functions $\{Q\}$ which are convenient to work with, is somewhat limited (most notably, the functions $Q(z) = z^2$, $Q(z) = 1/z$, $Q(z) = -\sqrt{z}$ and some piecewise linear functions [14],[15]). The next example demonstrates this point.

Example. Consider the information functional (51) with the convex function $Q(z) = -\sqrt{z}$. Then, the corresponding generalized mutual information is

$$J^Q(U; V) = - \sum_{u,v} P(u)[P(v|u) + sP(v)] \cdot \sqrt{\frac{P(v)}{P(v|u) + sP(v)}}$$

$$\begin{aligned}
&= - \sum_{u,v} P(u) \sqrt{P(v)[P(v|u) + sP(v)]} \\
&= - \sum_{u,v} P(u)P(v) \sqrt{s + \frac{P(v|u)}{P(v)}}.
\end{aligned} \tag{52}$$

Consider now the above-described problem of joint source–channel coding, for the following source and channel: The source is designated by a random variable U , which is uniformly distributed over the alphabet $\mathcal{U} = \{0, 1, \dots, K-1\}$. The reproduction variable, V , takes on values in the same alphabet, i.e., $\mathcal{V} = \mathcal{U} = \{0, 1, \dots, K-1\}$ and the distortion function is

$$d(u, v) = \begin{cases} 0 & v = u \\ 1 & v = (u+1) \bmod K \\ \infty & \text{elsewhere} \end{cases} \tag{53}$$

which means that errors other than $v = (u+1) \bmod K$ are strictly forbidden. Therefore the channel from U to V must be of the form

$$P(v|u) = \begin{cases} 1 - \epsilon_u & v = u \\ \epsilon_u & v = (u+1) \bmod K \\ 0 & \text{elsewhere} \end{cases} \tag{54}$$

where $\{\epsilon_u\}$ are parameters taking values in $[0, 1]$ and complying with the distortion constraint

$$\mathbf{E}\{d(U, V)\} = \frac{1}{K} \sum_{u=0}^{K-1} \epsilon_u \leq d. \tag{55}$$

The channel is a noise-free L -ary channel, i.e., its input and output alphabets are $\mathcal{X} = \mathcal{Y} = \{0, 1, \dots, L-1\}$ with $P(y|x) = 1$ for $y = x$, and $P(y|x) = 0$ otherwise.

Obviously, the case $K \leq L$ is not interesting because the data can be conveyed error-free by trivially connecting the source to the channel. In the other extreme, where $K > 2L$, there must be some channel input symbol to which at least three source symbols are mapped. In such a case, it is impossible to avoid at least one of the forbidden errors in the reconstruction. Thus, the interesting cases are those for which $L < K \leq 2L$, or equivalently, $\theta \in (1, 2]$, where $\theta \triangleq K/L$.

We next derive a distortion bound based on the generalized data processing theorem, in the spirit of [14] and [15], where we now have the parameter s as a degree of freedom.

As for the source, let us suppose that in addition to the distortion constraint, we impose the constraint that the distribution of the reproduction variable V , just like U , must be uniform over

its alphabet, namely, $P(v) = 1/K$ for all $v \in \mathcal{V}$. In this case,

$$\begin{aligned}
-J^Q(U; V) &= \sum_{u,v} P(u)P(v) \sqrt{s + \frac{P(v|u)}{P(v)}} \\
&= \frac{1}{K^2} \sum_{u=0}^{K-1} \left[\sqrt{s + K\epsilon_u} + \sqrt{s + K(1 - \epsilon_u)} + (K - 2)\sqrt{s} \right] \\
&= \frac{1}{K^2} \sum_{u=0}^{K-1} \left[\sqrt{s + K\epsilon_u} + \sqrt{s + K(1 - \epsilon_u)} \right] + \left(1 - \frac{2}{K}\right) \sqrt{s} \\
&\leq \frac{1}{K^2} \cdot K \left[\sqrt{s + Kd} + \sqrt{s + K(1 - d)} \right] + \left(1 - \frac{2}{K}\right) \sqrt{s} \\
&= \frac{1}{K} \left[\sqrt{s + Kd} + \sqrt{s + K(1 - d)} \right] + \left(1 - \frac{2}{K}\right) \sqrt{s}, \tag{56}
\end{aligned}$$

where the inequality follows from the fact that the maximum of the concave function

$$\sum_u [\sqrt{s + K\epsilon_u} + \sqrt{s + K(1 - \epsilon_u)}],$$

subject to the distortion constraint (55), is achieved when $\epsilon_u = d$ for all $u \in \mathcal{U}$. Thus,

$$R^Q(d) = -\frac{1}{K} \left[\sqrt{s + Kd} - \sqrt{s + K(1 - d)} \right] - \left(1 - \frac{2}{K}\right) \sqrt{s}. \tag{57}$$

As for the channel, we have:

$$\begin{aligned}
-J^Q(X; Y) &= \sum_{x,y} P(x)P(y) \sqrt{s + \frac{P(y|x)}{P(y)}} \\
&= \sum_{x' \neq x} P(x)P(x') \sqrt{s} + \sum_x P^2(x) \sqrt{s + \frac{1}{P(x)}} \\
&= \sqrt{s} \left[1 - \sum_x P^2(x) \right] + \sum_x P^2(x) \sqrt{s + \frac{1}{P(x)}} \\
&= \sqrt{s} + \sum_x P^2(x) \left(\sqrt{s + \frac{1}{P(x)}} - \sqrt{s} \right) \\
&= \sqrt{s} + \sum_x P^2(x) \cdot \frac{1/P(x)}{\sqrt{s + 1/P(x)} + \sqrt{s}} \\
&= \sqrt{s} + \sum_x \frac{P(x)}{\sqrt{s + 1/P(x)} + \sqrt{s}}. \tag{58}
\end{aligned}$$

The function $f(t) = t/[\sqrt{s + 1/t} + \sqrt{s}]$ is convex in t (for fixed s) since $f''(t) \geq 0$ for all $t \geq 0$, as can readily be verified. Thus, $-J^Q(X; Y)$ is minimized by the uniform distribution $P(x) = 1/L$,

$\forall x$, which leads to the ‘capacity’ expression:

$$C^Q = -\sqrt{s} - \frac{1}{\sqrt{s} + \sqrt{s+L}}. \quad (59)$$

Applying now the data processing theorem,

$$R^Q(d) \leq C^Q, \quad (60)$$

we obtain, after rearranging terms

$$\sqrt{s+Kd} + \sqrt{s+K(1-d)} \geq \frac{K}{\sqrt{s} + \sqrt{s+L}} + 2\sqrt{s}. \quad (61)$$

Squaring both sides, we have:

$$2s + K + 2\sqrt{(s+Kd)[s+K(1-d)]} \geq \left[\frac{K}{\sqrt{s} + \sqrt{s+L}} + 2\sqrt{s} \right]^2 \quad (62)$$

or

$$2\sqrt{(s+Kd)[s+K(1-d)]} \geq \left[\frac{K}{\sqrt{s} + \sqrt{s+L}} + 2\sqrt{s} \right]^2 - 2s - K, \quad (63)$$

which after squaring again and applying some further straightforward algebraic manipulations, gives eventually the following inequality on the distortion d :

$$4d(1-d) \geq \psi(s), \quad (64)$$

where

$$\psi(s) \triangleq \frac{1}{K^2} \left[\left(\frac{K}{\sqrt{s} + \sqrt{s+L}} + 2\sqrt{s} \right)^2 - 2s - K \right]^2 - \frac{4s(s+K)}{K^2}. \quad (65)$$

The resulting lower bound on the distortion is the smaller of the two solutions of the equation $4d(1-d) = \psi(s)$, which is

$$d_s \triangleq \frac{1}{2} - \frac{1}{2}\sqrt{1 - \psi(s)}. \quad (66)$$

Thus, the larger is $\psi(s)$, the better is the bound. The choice $s = 0$, which corresponds to the usual Ziv–Zakai bound for $Q(z) = -\sqrt{z}$, yields

$$\psi(0) = \frac{1}{K^2} \left[\left(\frac{K}{\sqrt{L}} \right)^2 - K \right]^2 = \left(\frac{K}{L} - 1 \right)^2 = (\theta - 1)^2. \quad (67)$$

However, it turns out that $s = 0$ is not the best choice of s . We next examine the limit $s \rightarrow \infty$. To this end, we derive a lower bound to $\psi(s)$ which is more convenient to analyze in this limit.

Note that for $s \geq L/8$, it is guaranteed that the expression in the square brackets of the expression defining $\psi(s)$, is positive, which means that an upper bound on $\sqrt{s+L}$ would yield a lower bound to $\psi(s)$. Thus, upper bounding $\sqrt{s+L}$ by

$$\sqrt{s+L} = \sqrt{s} \cdot \sqrt{1+L/s} \leq \sqrt{s} \left(1 + \frac{L}{2s}\right),$$

we get

$$\begin{aligned} K^2\psi(s) &= \left[\left(\frac{K}{\sqrt{s} + \sqrt{s+L}} + 2\sqrt{s} \right)^2 - 2s - K \right]^2 - 4s^2 - 4Ks \\ &\geq \left[\left(\frac{K}{\sqrt{s}(2+L/2s)} + 2\sqrt{s} \right)^2 - 2s - K \right]^2 - 4s^2 - 4Ks \\ &= K^2 \left(\frac{4s-L}{4s+L} \right)^2 + \frac{16K^4s^2}{(4s+L)^4} - \frac{8K^2Ls}{4s+L} + \frac{16K^2s^2}{(4s+L)^2} + \frac{8K^3s(4s-L)}{(4s+L)^3} \\ &\triangleq K^2\psi_0(s), \end{aligned} \tag{68}$$

where between the second and the third lines, we have skipped some standard algebraic operations. Taking now the limit $s \rightarrow \infty$, we obtain

$$\psi_\infty = \lim_{s \rightarrow \infty} \psi_0(s) = \frac{1}{K^2}(K^2 + 0 - 2KL + K^2 + 0) = 2 \left(1 - \frac{L}{K}\right) = 2 \left(1 - \frac{1}{\theta}\right), \tag{69}$$

which yields a better bound than the bound of $s = 0$ since

$$2 \left(1 - \frac{1}{\theta}\right) > (\theta - 1)^2 \tag{70}$$

for all $\theta \in (1, 2)$.

It is interesting to compare this also to the classical data processing theorem: Since

$$R(d) = \log K - h_2(d) \tag{71}$$

and

$$C = \log L, \tag{72}$$

then the ordinary data processing theorem yields the bound

$$h_2(d) \geq \log \theta. \tag{73}$$

Since

$$h_2(d) \geq 4d(1-d) \tag{74}$$

and

$$2 \left(1 - \frac{1}{\theta} \right) \geq \log_2 \theta \quad (75)$$

within the relevant range of θ , the bound pertaining to $s \rightarrow \infty$ is also better than the classical bound for this case. This completes the description of the example. \square

Finally, we should comment that the monotonicity result concerning $V(t)$ contains as special cases, not only the H–theorem, as well as all other earlier mentioned monotonicity results, but also the 1975 Zakai–Ziv generalized data processing [15]. Consider a Markov chain $U \rightarrow V \rightarrow W$, where U , V and W are random variables that take on values in (finite) alphabets, \mathcal{U} , \mathcal{V} , and \mathcal{W} , respectively. Let us now map between the Markov chain (U, V, W) and the Markov process $\{X_t\}$ in the following manner: $(u, v) \in \mathcal{U} \times \mathcal{V}$ is assigned to the state x' of the process at time t , whereas $(u, w) \in \mathcal{U} \times \mathcal{W}$ corresponds⁶ to x at time $t + 1$. Now, defining accordingly,

$$\mu_t^0(x') = P(u, v), \quad (76)$$

$$\mu_t^1(x') = P(u)P(v), \quad (77)$$

$$\mu_{t+1}^0(x) = P(u, w), \quad (78)$$

and

$$\mu_{t+1}^1(x) = P(u)P(w), \quad (79)$$

then due to the Markov property of (U, V, W) , both measures satisfy the recursion with $P(w|v)$ playing the role⁷ of $P(x|x')$. I.e.,

$$\begin{aligned} P(u, w) &\stackrel{\Delta}{=} \mu_{t+1}^0(x) \\ &= \sum_{x'} \mu_t^0(x') P(x|x') \\ &= \sum_v P(u, v) P(w|v) \end{aligned} \quad (80)$$

and

$$P(u)P(w) \stackrel{\Delta}{=} \mu_{t+1}^1(x)$$

⁶While \mathcal{V} and \mathcal{W} may be different (finite) alphabets, x and x' , of the original Markov process, must taken on values in the same alphabet. Assuming, without loss of generality, that $\mathcal{V} = \{1, 2, \dots, |\mathcal{V}|\}$ and $\mathcal{W} = \{1, 2, \dots, |\mathcal{W}|\}$, then for the purpose of this mapping, we can unify these alphabets to be both $\{1, 2, \dots, \max\{|\mathcal{V}|, |\mathcal{W}|\}\}$ and complete the missing elements of the extended transition matrix $P(w|v)$ in a consistent manner, according to the actual support of each distribution. We omit further technical details herein.

⁷Consider the component u of $x' = (u, v)$ and $x = (u, w)$ simply as an index.

$$\begin{aligned}
&= \sum_{x'} \mu_t^1(x') P(x|x') \\
&= \sum_v P(u) P(v) P(w|v)
\end{aligned} \tag{81}$$

Thus, for $Q(z) = -\ln z$, the monotonicity of $V(t)$ is nothing but the data processing of the classical mutual information. For a general function Q of one variable ($k = 1$), this gives the generalized data processing theorem of [14]. Furthermore, letting Q be a general convex function of k variables, and $\mu_t^0(x') = P(u, v)$ as before, we get the more general form of the data processing inequality of [15].

The above extension of the H-theorem gives rise to a seemingly more general data processing theorem than in [15], as it is not necessary to let $\mu_t^0(x)$ be the actual joint probability distribution. However, when looking at the entire class of convex functions with an arbitrary number of arguments, this is not really more general, as the corresponding generalized mutual information can readily be transformed back to the form of the 1975 Zakai–Ziv information functional using again the perspective operation. Indeed, as mentioned in the Introduction and shown in [15, Theorem 7.1], the class of generalized mutual information measures studied therein cannot be improved upon in the sense that there always exist choices of Q and $\{\mu_i\}$ that provide tight bounds on the distortion of the optimum system.

4 Summary and Conclusion

The main contributions of this work can be summarized as follows: First, we have established a unified framework and a relationship between (a generalized version of) the second law of thermodynamics and the generalized data processing theorems of Zakai and Ziv. This unified framework turns out to strengthen and expand both of these pieces of theory: Concerning the second law of thermodynamics, we have identified a significantly more general information measure, which is a monotonic function of time, when it operates on a Markov process. As for the generalized Ziv–Zakai data processing theorem, we have proposed a wider class of information measures obeying the data processing theorem, which includes free parameters that may be optimized so as to tighten the distortion bounds.

Acknowledgment

Interesting discussions with J. Ziv and M. Zakai are acknowledged with thanks.

References

- [1] D. Andelman, “Bounds according to a generalized data processing theorem,” M.Sc. dissertation, Department of Electrical Engineering, Technion, Israel Institute of Technology, Haifa, Israel, October 1974.
- [2] G. B. Bağcı, “The physical meaning of Rényi relative entropies,” arXiv:cond-mat/0703008v1, March 1, 2007.
- [3] A. H. W. Beck, *Statistical Mechanics, Fluctuations and Noise*, Edward Arnold Publishers, 1976.
- [4] S. Boyd and L. Vandenberghe, *Convex Optimization*, Cambridge University Press, 2004.
- [5] T. M. Cover and J. A. Thomas, *Elements of Information Theory*, second edition, John Wiley & Sons, 2006.
- [6] I. Csiszár, “A class of measures of informativity of observation channels,” *Periodica Mathematica Hungarica*, vol. 22, no. 1–4, pp. 191–213, 1972.
- [7] B. Dacorogna and P. Maréchal, “The role of perspective functions in convexity, polyconvexity, rank-one convexity and separate convexity,”
http://caa.epfl.ch/publications/2008-The_role_of_perspective_functions_in_convexity.pdf
- [8] M. Kardar, *Statistical Physics of Particles*, Cambridge University Press, 2007.
- [9] F. P. Kelly, *Reversibility and Stochastic Networks*, J. Wiley & Sons, 1979.
- [10] C. Kittel, *Elementary Statistical Physics*, John Wiley & Sons, 1958.
- [11] D. P. Palomar and S. Verdú, “Lautum information,” *IEEE Trans. Inform. Theory*, vol. 54, no. 3, pp. 964–975, March 2008.
- [12] F. Reif, *Fundamentals of Statistical and Thermal Physics*, McGraw–Hill, 1965.
- [13] A. J. Viterbi and J. K. Omura, *Principles of Digital Communication and Coding*, McGraw–Hill, 1979.

- [14] J. Ziv and M. Zakai, "On functionals satisfying a data-processing theorem," *IEEE Trans. Inform. Theory*, vol. IT-19, no. 3, pp. 275–283, May 1973.
- [15] M. Zakai and J. Ziv, "A generalization of the rate-distortion theory and applications," in: *Information Theory New Trends and Open Problems*, edited by G. Longo, Springer-Verlag, pp. 87–123, 1975.