

# Il valore dell'informazione

L. P.

SMRI (Italy)

## Sommario

Un ben noto risultato di Kelly mostra che l'informazione mutua tra un processo e un segnale relativo ad esso misura il vantaggio che può esserne ottenuto nell'ambito del gioco d'azzardo (p.es., scommesse su corse di cavalli). Questa misura del valore dell'informazione può essere ritrovata in altre forme, come in termodinamica (estrazione di lavoro da un serbatoio di calore) o in dinamica delle popolazioni (adattamento fenotipico all'ambiente). In questo articolo, percorro a grandi linee il ragionamento che porta a queste conclusioni.

## 1 La misura dell'incertezza

La teoria matematica dell'informazione nasce con il lavoro di Claude Shannon "A Mathematical Theory of Communication" [1, 2], che considera un modello ideale di trasmissione di messaggi tramite un canale di comunicazione soggetto a rumore. Il problema fondamentale della comunicazione è quello di ridurre l'incertezza che il recipiente del messaggio ha rispetto al contenuto del messaggio stesso. Prima dell'emissione del messaggio, il recipiente ha una certa informazione sui *possibili* messaggi che possono essere trasmessi. Avendo ricevuto il messaggio, che è parzialmente corrotto dal rumore, la sua incertezza è diminuita. Al limite, se è possibile ricostruire il messaggio senza errori, la sua incertezza è pari a zero. Questo suggerisce di quantificare il contenuto d'informazione del messaggio come la differenza fra l'incertezza a priori e l'incertezza a posteriori. Dobbiamo considerare il messaggio  $X$  come una variabile aleatoria, che ha una certa probabilità a priori  $p_x$  di assumere il valore  $x$ . Allora Shannon dimostra che l'unica misura dell'incertezza su  $X$  che soddisfa degli assiomi semplici ed intuitivi è data dall'espressione

$$H(X) = -k \sum_{x=1}^N p_x \ln p_x, \quad (1)$$

dove  $k$  è una costante positiva, per altro arbitraria. Shannon scelse per  $k$  il valore  $k = 1/\ln 2$ . In questo modo, l'incertezza unitaria corrisponde a una semplice alternativa tra due possibili valori con uguale probabilità  $p_x = \frac{1}{2}$ . Seguendo un suggerimento di John W. Tukey, Shannon chiamò "bit" questa unità, come abbreviazione

di *binary digit* (cifra binaria), ma corrispondente anche all'espressione inglese che traduce "un po' ". In questo articolo useremo  $k = 1$ , a meno di avviso contrario. L'espressione in eq. (1) è nota come **entropia di Shannon**.<sup>1</sup> Shannon (citato in [3]) rivelò che il nome "entropia" gli era stato suggerito da John von Neumann:

Dovresti chiamarla "entropia" per due ragioni: la prima perché la tua espressione è identica a quella usata da Gibbs per l'entropia termodinamica, e la seconda e più importante, perché nessuno sa sul serio che cosa sia l'entropia e quindi avresti sempre un vantaggio nelle discussioni.

In effetti, Gibbs nel suo trattato [4] aveva mostrato che l'entropia termodinamica di un sistema all'equilibrio, in cui la probabilità di trovarsi nello stato microscopico  $x$  è pari a  $p_x$ , è data dall'espressione

$$S = -k_B \sum_x p_x \ln p_x, \quad (2)$$

che corrisponde alla (2) con  $k$  uguale alla costante di Boltzmann  $k_B$ . Shannon fece un brillante uso della sua misura dell'incertezza per impostare e risolvere il problema della codifica ottimale di messaggi lungo un canale di comunicazione soggetto a rumore. In particolare dimostrò che esisteva un valore limite quantificabile alla capacità di trasmissione di un canale. Trasmettendo messaggi a una cadenza inferiore a questa capacità è possibile codificarli in modo da recuperarli con un tasso d'errore arbitrariamente piccolo, mentre questo non è possibile se la cadenza supera questo limite.

Tuttavia la misura dell'incertezza ha trovato applicazione in campi molto diversi. In particolare vorrei descrivere in questo articolo alcune applicazioni, partendo dal classico esempio del gioco d'azzardo, alla termodinamica, e perfino a problemi di interesse biologico, come la percezione, la diversificazione fenotipica, e la possibilità di interpretare l'evoluzione darwiniana come un processo di acquisizione dell'informazione.

## 2 Corse di cavalli

Nel 1956, Kelly propose un'interpretazione della capacità di trasmissione che andava al di là del problema della codifica ottimale [5]. Il problema non è tanto quello di trasmettere dei messaggi senza errori, quanto quello di assegnare un valore a una comunicazione ricevuta con errore. L'idea fondamentale è che il destinatario dovrebbe poter approfittare della conoscenza (anche se incerta) del messaggio trasmesso, e che quindi una misura del valore dell'informazione può essere basata sul vantaggio ottenuto.

La situazione considerata da Kelly ricorda il classico film del 1973 "La Stangata", con Paul Newman e Robert Redford. Consideriamo un giocatore che può scommettere su una corsa di cavalli che si svolge in un'altra città. Un complice è in grado di comunicare al giocatore il risultato della corsa prima che essa sia nota all'allibratore, mettendolo in grado di vincere a colpo (quasi) sicuro. Il "quasi" è dovuto al fatto che il canale di comunicazione è soggetto a errori, e talvolta

---

<sup>1</sup>La dimostrazione dell'espressione (1) della misura dell'incertezza è presentata in appendice A.

trasmette il risultato sbagliato. Supponiamo per semplicità che ci siano solo due cavalli, con uguale probabilità di vittoria a priori. Il segnale ricevuto è esatto con probabilità  $p$  e sbagliato con probabilità  $q = 1 - p$ . Se il giocatore punta sul cavallo corretto, vince 2 volte la posta. Qual è la migliore strategia per il giocatore?

Se il giocatore può giocare una sola volta, e  $p > \frac{1}{2}$ , la situazione è chiara. Il giocatore tenderà ad ottimizzare il valor medio della vincita, che è pari a due volte la posta se sceglie il cavallo giusto (con probabilità  $p$ ) e a 0 se sceglie il cavallo sbagliato. Quindi il giocatore punterà tutto il suo capitale sul cavallo indicato dal complice, e ritornerà a tasche vuote con probabilità  $q$ .

Le cose cambiano se le corse si ripetono indefinitamente (sempre con le stesse probabilità) e il giocatore può reinvestire il proprio capitale. In questo caso, a ogni corsa, con probabilità  $p$  il giocatore raddoppierebbe il suo capitale, e con probabilità  $q$  lo perderebbe tutto. È vero che il valor medio della vincita viene ad essere uguale a  $(2p)^N$  dopo  $N$  corse, ma il giocatore avrebbe perso tutto con probabilità  $1 - q^N$ , che tende rapidamente a 1 al crescere di  $N$ .

Di conseguenza conviene che il giocatore non punti tutto il suo capitale, in modo da garantirsi nel caso che le informazioni ricevute siano sbagliate. Indichiamo con  $b$  la frazione del capitale puntata sul cavallo suggerito. Quale criterio adottare per scegliere  $b$ ? Consideriamo il capitale del giocatore dopo  $N$  giocate, in cui il cavallo suggerito ha vinto  $N_+$  volte, e ha perso  $N_-$  volte. Otteniamo

$$V_N = (1 + b)^{N_+} (1 - b)^{N_-} V_0, \quad (3)$$

dove  $V_0$  è il capitale iniziale. Da questa espressione è chiaro che il capitale vinto tende a crescere (o decrescere) esponenzialmente. Questo suggerisce di ottimizzare non il valor medio di  $V_N$ , ma il suo logaritmo, o meglio, il logaritmo del rapporto fra  $V_N$  e  $V_0$ . Definiamo

$$\Lambda(b) = \lim_{N \rightarrow \infty} \frac{1}{N} \ln \frac{V_N}{V_0}. \quad (4)$$

La quantità  $\Lambda(b)$  è detta **tasso di crescita del capitale**. Si ha

$$\ln \frac{V_N}{V_0} = \sum_{i=1}^N \ln [\epsilon_i (1 + b) + (1 - \epsilon_i) (1 - b)], \quad (5)$$

dove  $\epsilon_i = 1$  se nella corsa  $i$  vince il cavallo suggerito, e  $\epsilon_i = 0$  se esso invece perde. Questa espressione è la somma di variabili indipendenti, identicamente distribuite, e per la legge dei grandi numeri la sua media empirica, che appare in eq. (4), tende con probabilità 1 alla sua media valutata sulla distribuzione di probabilità.<sup>2</sup> Dato che  $\langle \epsilon_i \rangle = p$ ,  $\forall i$ , otteniamo

$$\Lambda(b) = \langle \ln [\epsilon(1 + b) + (1 - \epsilon)(1 - b)] \rangle = p \ln(1 + b) + q \ln(1 - b). \quad (6)$$

La strategia ottimale  $b^*$  si ottiene imponendo la condizione

$$\Lambda'(b^*) = \frac{p}{1 + b^*} - \frac{q}{1 - b^*} = 0, \quad (7)$$

<sup>2</sup>In questo articolo, utilizzo le parentesi angolari per denotare la media. Data una funzione  $f_x$  della variabile aleatoria discreta  $X$ , si ha  $\langle f \rangle = \sum_x p_x f_x$ .

che dà

$$b^* = 2p - 1. \quad (8)$$

Il tasso di crescita ottimale  $\Lambda(b^*)$  è dato da

$$\Lambda(b^*) = \ln 2 + p \ln p + q \ln q, \quad (9)$$

dove abbiamo usato la relazione  $q = 1 - p$ . Ma questa espressione può essere scritta

$$\Lambda(b^*) = \ln 2 - H(X), \quad (10)$$

dove

$$H(X) = -(p \ln p + q \ln q) \quad (11)$$

è la misura dell'incertezza contenuta nel messaggio dato dal complice (con la costante  $k = 1$ ). Possiamo interpretare questo risultato nel modo seguente: se il risultato della corsa fosse certo, il giocatore raddoppierebbe il suo capitale a ogni corsa, e si avrebbe  $\Lambda = \ln 2$ . A causa dell'incertezza della comunicazione, la vincita viene ridotta di una quantità esattamente uguale all'entropia del messaggio. Alla peggio, se l'incertezza del messaggio è totale, si ha  $H(X) = \ln 2$  e la vincita è nulla (in media). O detto in altro modo, la riduzione dell'incertezza sul vincitore della corsa, da  $\ln 2$  in assenza d'informazione a  $H(X)$  in presenza del messaggio, è uguale al tasso di crescita del capitale, e costituisce una misura del valore dell'informazione contenuta nel messaggio.

Il fatto che il valore dell'informazione ricevuta stia nella riduzione dell'entropia viene confermato dall'osservazione che il tasso ottimale di crescita del capitale non dipende dai dettagli della strategia del giocatore. Supponiamo, p.es., che il giocatore, invece di trattenere una frazione  $1 - b$  del capitale, decida di utilizzarla per scommettere sull'altro cavallo, quello *non* suggerito dal messaggio. In questo caso,  $\Lambda(b)$  viene espresso da

$$\Lambda(b) = p \ln(2b) + q \ln(2(1 - b)) = \ln 2 + p \ln b + q \ln(1 - b). \quad (12)$$

La strategia ottimale è data da

$$b^* = p. \quad (13)$$

Sostituendo, otteniamo

$$\Lambda(b^*) = \ln 2 + p \ln p + q \ln q = \ln 2 - H(X), \quad (14)$$

esattamente come nel caso precedente.

### 3 Un caso più generale

Possiamo analizzare il problema in una situazione più generale, in cui ci sono  $N$  cavalli, ciascuno con una probabilità  $p_x$  di vincere ( $x \in \{1, \dots, N\}$ ), e in cui l'allibratore restituisce  $o_x$  volte la posta se vince il cavallo  $x$ . Questo paragrafo è necessariamente più tecnico, perché richiede di introdurre un certo numero di

quantità analoghe all'entropia di Shannon, di cui faremo uso nel seguito. Supponiamo di avere a che fare con un allibratore "onesto", che non trattiene per sé (in media) parte delle scommesse. In questo caso si ha

$$\sum_{x=1}^N \frac{1}{o_x} = 1. \quad (15)$$

In questa situazione  $1/o_x$  può essere interpretata come la stima  $q_x$  che l'allibratore fa della probabilità che vinca il cavallo  $x$ . Indichiamo con  $b = (b_x)$  la frazione del capitale che il giocatore punta sul cavallo  $x$ , con  $\sum_x b_x = 1$ . Si ha allora

$$\Lambda(b) = \sum_x p_x \ln \frac{b_x}{q_x}. \quad (16)$$

Ottimizzando rispetto a  $b$ ,<sup>3</sup> e tenendo conto della condizione  $\sum_x b_x = 1$  si ottiene

$$b_x^* = p_x, \quad \forall x. \quad (17)$$

Sorprendentemente, le quote  $o = (o_x)$  non influenzano la strategia ottimale  $b^*$ , ma solo il tasso di crescita del capitale. Si ha infatti

$$\Lambda(b^*) = \sum_x p_x \ln(o_x b_x^*) = \sum_x p_x \ln \frac{p_x}{q_x}. \quad (18)$$

L'espressione a secondo membro è definita non negativa, e si annulla solo se  $p_x = q_x, \forall x$ .<sup>4</sup> Essa costituisce una misura di quanto le distribuzioni di probabilità  $p = (p_x)$  e  $q = (q_x)$  sono differenti, e viene chiamata **divergenza di Kullback-Leibler** e denotata con  $D_{\text{KL}}(p||q)$ . Non può essere considerata una distanza: in particolare non è simmetrica, perché si può vedere che se, p.es.,  $q_x = 0$  per un qualche  $x$ , mentre  $p_x > 0, \forall x$ , si ha  $D_{\text{KL}}(p||q) = \infty$ , ma  $D_{\text{KL}}(q||p)$  è finito.

Anche in questa situazione possiamo considerare l'effetto di un canale di comunicazione che produce un segnale  $Y$ . Indichiamo con  $p = (p_x)$  le probabilità a priori che vinca il cavallo  $x$ , e con  $p_{x|y}$  la probabilità condizionata che il cavallo vincente sia  $x$  dato che il giocatore abbia ricevuto il segnale  $y$ . Indichiamo inoltre con  $p_y$  la probabilità che il segnale ricevuto sia  $y$ . Queste probabilità sono legate alla probabilità congiunta  $p_{x,y}$  che vinca il cavallo  $x$  e che il segnale ricevuto sia  $y$ :

$$p_{x,y} = p_{x|y}p_y, \quad p_x = \sum_y p_{x,y}, \quad p_y = \sum_x p_{x,y}. \quad (19)$$

Ci si riferisce a  $p_x$  e  $p_y$ , in questo contesto, come alle **marginali** della distribuzione congiunta  $p_{x,y}$ . Dobbiamo anche introdurre  $b_{x|y}$ , che è la frazione del capitale scommessa sul cavallo  $x$  quando sia stato ricevuto il segnale  $y$ . Essa soddisfa le relazioni

$$\sum_x b_{x|y} = 1, \quad \forall y. \quad (20)$$

Si ha allora

$$\Lambda_Y(b) = \sum_{x,y} p_{x,y} \ln(o_x b_{x|y}) = \sum_{x,y} p_{x|y} p_y \ln \frac{b_{x|y}}{q_x}, \quad (21)$$

<sup>3</sup>La dimostrazione è riportata nell'appendice B.

<sup>4</sup>Queste proprietà sono dimostrate nell'appendice C.

dove  $\Lambda_Y(b)$  ricorda che stiamo considerando il processo in presenza del segnale  $Y$ . Otteniamo così la strategia ottimale

$$b_{x|y}^* = p_{x|y}, \quad \forall x, y, \quad (22)$$

e il corrispondente tasso di crescita:

$$\Lambda_Y(b^*) = \sum_{x,y} p_{x,y} \ln \frac{p_{x|y}}{q_x} = \sum_{x,y} p_{x,y} \ln \frac{p_{x,y}}{q_x p_y}. \quad (23)$$

Possiamo riscrivere questa espressione come segue:

$$\Lambda_Y(b^*) = \sum_{x,y} p_{x,y} \ln \frac{p_{x,y}}{p_x p_y} + \sum_x p_x \ln \frac{p_x}{q_x}. \quad (24)$$

Il secondo termine è  $D_{\text{KL}}(p||q)$ , che corrisponde al tasso ottimale di vincita in assenza di comunicazione parallela. Il primo termine è la divergenza di Kullback-Leibler fra la probabilità congiunta  $p_{x,y}$  e il prodotto delle marginali  $p_x$  e  $p_y$ . Questa quantità è non negativa e si annulla solo se  $X$  e  $Y$  sono indipendenti. Indichiamola con  $I(X : Y)$ . Per interpretarla, date due variabili  $X$  e  $Y$ , con distribuzione congiunta  $p_{x,y}$ , definiamo l'**entropia condizionata**  $H(X|Y)$  mediante l'espressione

$$H(X|Y) = - \sum_{x,y} p_{x,y} \ln p_{x|y} = - \sum_y p_y \sum_x p_{x|y} \ln p_{x|y}. \quad (25)$$

L'entropia condizionata è il valor medio dell'entropia della distribuzione condizionata di  $x$  dato  $y$ , valutato sulla marginale di  $y$ . Si ha allora

$$I(X : Y) = H(X) - H(X|Y) = H(Y) - H(Y|X). \quad (26)$$

Vediamo così che  $I(X : Y)$  dà una misura di quanto la conoscenza di una variabile diminuisce l'incertezza sull'altra variabile. Essa è nota come **informazione mutua**. Dalla sua espressione si vede che è simmetrica:

$$I(X : Y) = I(Y : X). \quad (27)$$

Concludendo, abbiamo ottenuto

$$\Lambda_Y(b^*) = I(X : Y) + D_{\text{KL}}(p||q), \quad (28)$$

che mostra che l'informazione mutua è pari all'incremento del tasso di crescita del capitale dovuto alla presenza del segnale  $Y$ . In particolare, nel caso di un allibratore non solo onesto ( $\sum_x 1/o_x = 1$ ), ma anche competente ( $q_x = p_x, \forall x$ ), il tasso di crescita ottimale del capitale è esattamente uguale all'informazione mutua fra segnale e risultato.

## 4 La macchina di Szilard

Il secondo principio della termodinamica, nella formulazione di Kelvin, stipula che non è possibile disegnare una macchina che, operando ciclicamente, estragga

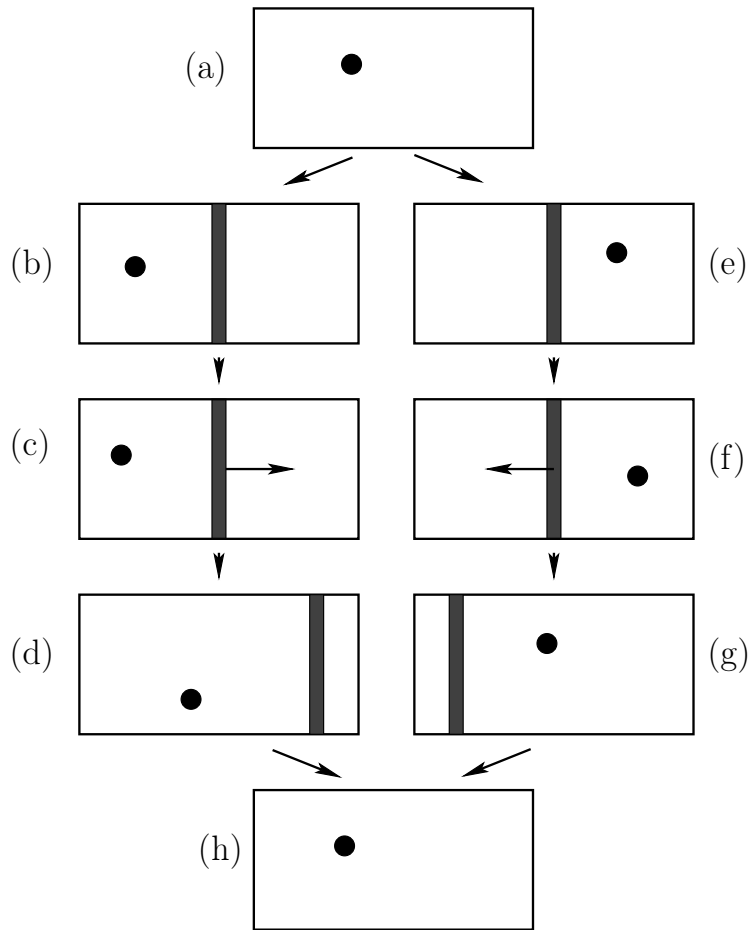


Figura 1: Il ciclo della macchina di Szilard. Si pone una partizione in un cilindro, dividendolo in due sezioni di volume  $p_S \mathcal{V}$  e  $p_D \mathcal{V}$  rispettivamente. Quindi si effettua una misura della posizione (S o D) della molecola. Dato il risultato  $y \in \{S, D\}$  della misura, la partizione viene spostata reversibilmente dalla posizione iniziale alla posizione finale. La partizione viene rimossa e il ciclo ricomincia.

energia sotto forma di calore da un serbatoio a temperatura uniforme e la restituisca sotto forma di lavoro. Nel 1929, Leo Szilard (noto fra l'altro per avere scritto con Einstein la lettera a Roosevelt che lanciò il Progetto Manhattan) propose un esperimento concettuale che sembrava violare questo principio [6]. La macchina è costituita da un cilindro di volume  $\mathcal{V}$  contenente una sola molecola di gas, e in contatto con un serbatoio di calore alla temperatura  $T$  (cf. fig. 1). Un "demone" dispone nel cilindro una partizione mobile, dividendolo in due sezioni, non necessariamente uguali. Indichiamo con  $p_S$  la frazione del volume contenuta nella sezione di sinistra e con  $p_D$  quella della sezione destra. A questo punto il demone identifica in quale sezione si trova la molecola. Se essa, per esempio, si trova nella sezione sinistra, esso attacca una barra alla parte destra della sezione mobile. (Se essa si trova a destra, l'attacca alla parte sinistra.) Quindi lascia espandere reversibilmente la sezione contenente la molecola, finché la partizione tocca la pa-

rete del cilindro. In questa espansione viene compiuto lavoro sull'esterno, pari a  $k_B T \ln(\mathcal{V}/\mathcal{V}_0)$ , dove  $\mathcal{V}_0$  è il volume iniziale della partizione, pari a  $p_S \mathcal{V}$  se la molecola si trovava a sinistra, o a  $p_D \mathcal{V}$  se si trovava a destra. A questo punto il demone rimuove la partizione e il sistema ritorna apparentemente alla condizione iniziale.

Valutiamo il valor medio  $W$  del lavoro estratto dal demone in un ciclo. La molecola si trova a sinistra con probabilità  $p_S$ . Supponendo che il demone non si sbaglia mai, esso estrarrà in questo caso una quantità di lavoro pari a  $k_B T \ln(\mathcal{V}/\mathcal{V}_0) = -k_B T \ln p_S$ . Ripetendo l'argomento per il caso in cui la molecola si trova a destra, otteniamo il valor medio

$$W = -k_B T (p_S \ln p_S + p_D \ln p_D) = k_B T H(X), \quad (29)$$

dove  $X$  è una variabile aleatoria che assume i valori  $(S, D)$  rispettivamente con probabilità  $(p_S, p_D)$ . Vediamo così che la riduzione dell'incertezza sulla posizione della molecola, dovuta alla misura, può essere utilizzata per convertire calore in lavoro.

Se la misura è affetta da errore, dobbiamo tener conto del fatto che, se la misura è sbagliata, comprimere il gas contro la parete (riducendo il suo volume a 0) costerebbe un lavoro infinito. Quindi la partizione non può toccare la parete del cilindro alla fine della manipolazione. Indichiamo con  $b_{x|y}$  la frazione del volume  $\mathcal{V}$  in cui si trova la particella alla fine della manipolazione, data la sua posizione iniziale  $x \in \{S, D\}$  e il risultato della misura  $y \in \{S, D\}$ . La probabilità congiunta di  $x$  e  $y$  è denotata con  $p_{x,y}$ , mentre  $p_{x|y}$  è, come nel paragrafo precedente, la probabilità condizionata che la molecola sia in  $x$  se il risultato della misura è  $y$ . Il lavoro ottenuto dalla manipolazione è quindi dato da  $k_B T \ln(b_{x|y}/p_x)$ , e può essere positivo (se la misura è corretta) o negativo (se è sbagliata). Il valor medio  $W$  del lavoro estratto è dato da

$$W(b) = k_B T \sum_{x,y} p_{x,y} \ln \frac{b_{x|y}}{p_x}. \quad (30)$$

Vediamo che questa espressione coincide con il tasso di crescita del capitale dato in eq. (21), per il caso dell'allibratore competente (che sa cioè valutare correttamente le probabilità) moltiplicato per l'energia "termica"  $k_B T$ . In analogia con quel caso, vediamo che la scelta ottimale di  $b$  è data da

$$b_{x|y}^* = p_{x|y}, \quad \forall y, \quad (31)$$

e il lavoro medio  $W$  corrispondentemente estratto è dato da

$$W(b^*) = k_B T \sum_{x,y} p_{x,y} \ln \frac{p_{x|y}}{p_x} = k_B T \sum_{x,y} p_{x,y} \ln \frac{p_{x,y}}{p_x p_y} = k_B T I(X : Y). \quad (32)$$

Otteniamo così il risultato che il lavoro medio  $W$  che può essere estratto da una macchina di Szilard soddisfa la disuguaglianza

$$W \leq k_B T I(X : Y), \quad (33)$$

dove  $I(X : Y)$  è l'informazione mutua tra la misura  $Y$  e la posizione  $X$  della molecola. Questo risultato è stato esteso a situazioni più generali da Sagawa e Ueda [7].



Come si concilia l'operazione della macchina di Szilard con il secondo principio della termodinamica? Szilard [6] e altri ricercatori, fra cui Léon Brillouin [8] hanno suggerito che il processo di misura è intrinsecamente dissipativo, e che la dissipazione implicata non può essere inferiore a  $k_B T \Delta H(X)$ , dove  $\Delta H(X)$  è la differenza fra l'incertezza su  $X$  prima e dopo la misura. Tuttavia Bennett [9] ha fatto vedere che è possibile, in linea di principio, eseguire delle misure senza dissipazione, purché l'apparato di misura sia in uno stato ben definito prima della misura stessa. Conseguentemente, il ciclo della macchina di Szilard non riporta il sistema esattamente nello stato iniziale, perché l'apparato di misura conserva il risultato della misura stessa. Per ottenere una trasformazione perfettamente ciclica è necessario riportarlo nello stato definito iniziale. In questa operazione l'incertezza sullo stato dell'apparato di misura viene cancellata, in altri termini, l'entropia  $H(Y)$  del sistema diminuisce. Questo violerebbe il secondo principio, a meno che non venga dissipata una quantità di lavoro almeno pari a  $k_B T \Delta H(Y)$ , dove  $\Delta H(Y)$  è la differenza dell'incertezza su  $Y$  prima e dopo il ripristino dello stato iniziale dell'apparato. Un'analisi basata sulla termodinamica stocastica mostra che in condizioni molto generali questo è esattamente ciò che avviene. Questo risultato venne suggerito da Rolf Landauer nel 1961 [10], ed è stato sperimentalmente dimostrato in diversi sistemi mesoscopici in anni recenti [11, 12]. Esso è noto come il **principio di Landauer** ed è stato enunciato da Bennett [9] come segue:

Ogni manipolazione logicamente irreversibile dell'informazione, come la cancellazione di un bit o la confluenza di due percorsi computazionali, deve essere accompagnata da un corrispondente incremento d'entropia nei gradi di libertà non portanti informazione dell'apparato di manipolazione dell'informazione o del suo ambiente.

Una volta tenuto conto della dissipazione conseguente al ripristino dell'apparato di misura, si vede che il valor medio  $W$  del lavoro netto ottenuto dalla macchina di Szilard non può essere positivo.

## 5 Analogia fra gioco d'azzardo ed estrazione di lavoro

L'analogia fra il funzionamento della macchina di Szilard e le corse di Kelly è stata notata da Vinkler e collaboratori [13]. Questa analogia è dettagliata nella tabella 1. Notiamo che, a differenza dell'allibratore delle corse di Kelly, nella macchina di Szilard si ha a che fare con un allibratore sempre competente, in cui l'analogo della quota  $o_x$  è sempre uguale all'inverso della probabilità  $p_x$  della posizione iniziale. Tuttavia, a differenza delle corse di cavalli, in questo caso la posizione iniziale della partizione può essere scelta dal demone, il che corrisponde a scegliere arbitrariamente la quota  $o_x$  (sempre in maniera competente). Quindi nella macchina di Szilard è possibile scegliere la posizione della partizione iniziale in modo da ottimizzare la mutua informazione  $I(X : Y)$ , e disporre conseguentemente le posizioni finali  $b_{x|y}$  in modo da ottimizzare  $W$ . Otteniamo così il valore massimo  $W_{\max}$  del lavoro estraibile:

$$W_{\max} = \max_{p_x} I(X : Y). \quad (34)$$

Riassumendo: l'analogia non è perfetta, perché nel caso di Kelly l'allibratore ha la possibilità di scegliere le quote  $o_x$  arbitrariamente (ma la scelta ottimale è sempre

$o_x = 1/p_x$ ), mentre nel caso di Szilard il giocatore (cioè il demone) ha la possibilità di scegliere  $p_x$  in modo da ottimizzare il lavoro estratto.

Corse di Kelly	Macchina di Szilard
$x_i$ – risultato della corsa $i$	$x_i$ – posizione della molecola nel ciclo $i$
Informazione parallela	Risultato della misura
$y_i$ – informazione sulla corsa $i$	$y_i$ – misura imprecisa nel ciclo $i$
$(p_x)$ – distr. di prob. del risultato	$(p_x)$ – distr. di prob. della posizione
$p_{x y}$ – prob. cond. di $x$ , dato $y$	$p_{x y}$ – prob. cond. di $x$ data la misura
$o_x$ – quota della scommessa su $x$	$1/p_x$ – inverso del volume iniziale
Piazzare scommesse	Muovere la partizione al punto finale
Logaritmo del capitale	Lavoro estratto
Tasso di crescita del capitale	Lavoro medio estratto per ciclo

Tabella 1: Analogia fra le corse di Kelly e la macchina di Szilard. Adattato dalla ref. [13].

L’analogia può essere estesa al caso in cui una parte del lavoro estratto viene dissipata irreversibilmente (per esempio a causa di attriti). Questo corrisponde al caso in cui l’allibratore detiene per sé una parte delle scommesse, e si ha così  $\sum_x 1/o_x > 1$ . Per questo caso, Kelly ha elaborato una strategia sofisticata, che contempla in generale la possibilità di trattenere una parte del capitale. Nel caso della macchina di Szilard, questo corrisponde alla possibilità di non compiere la manipolazione se il lavoro dissipato viene ad essere in eccesso di quello estratto.

## 6 Popolazioni e fitness

Nel caso delle popolazioni, il capitale è costituito dal numero di individui che la compongono. Si può dire che un individuo piazza una scommessa producendo prole, e che questa scommessa viene a maturazione quando la prole è a sua volta in grado di riprodursi. La **fitness** di un individuo è il numero medio dei suoi “figli” che giunge a potersi riprodurre. Quindi la fitness è analoga alla quota nelle scommesse, e possiamo in linea di principio applicare a questo caso i ragionamenti fatti nei paragrafi precedenti.

Seguendo Donaldson-Matasci e collaboratori [14], consideriamo una forma di vita esposta ad un ambiente fluttuante  $X$ , che può assumere con probabilità  $p_x$  uno stato  $x$  fra un certo numero di stati possibili. Il fenotipo  $\Phi$  dell’organismo può anch’esso assumere un certo numero di stati  $\phi$ . Ciascun fenotipo prospera se si trova in un ambiente adatto, e altrimenti viene penalizzato. Questo significa che la fitness  $f_{x|\phi}$  del fenotipo  $\phi$  nell’ambiente  $x$  è grande se  $\phi$  è adatto ad  $x$  ed è piccola altrimenti. Adottiamo la convenzione che il fenotipo adatto all’ambiente  $x$  ha il valore  $x$ . Consideriamo il caso limite in cui il fenotipo  $\phi$ , con  $\phi \neq x$ , non è in grado di riprodursi.<sup>5</sup> Quindi la taglia media della prole di un individuo con fenotipo  $\phi$  in un ambiente  $x$  è pari a  $f_{x|x}$  se  $\phi = x$ , ed a 0 altrimenti. Di conseguenza, conviene

<sup>5</sup>Questa condizione può essere resa meno stringente (cf. Haccou e Iwasa [15]), ma non è il caso di discutere qui come procedere nel caso più generale.

agli individui produrre prole con fenotipi diversi, allo stesso modo che conviene ai giocatori fare puntate su diversi cavalli. Questo fenomeno è stato osservato, p.es., in certe popolazioni batteriche, in cui una parte degli individui hanno un metabolismo più lento (persisters) a parità di genotipo, il che permette loro di resistere agli antibiotici (vedi, p.es. [16]). Denotiamo  $f_{x|x}$  con  $o_x$ , e con  $b_x$  la frazione di individui della prole che posseggono il fenotipo  $x$ . Allora il tasso medio di incremento della popolazione, definito da

$$\Lambda(b) = \lim_{N \rightarrow \infty} \frac{1}{N} \ln \frac{\mathcal{N}_N}{\mathcal{N}_0}, \quad (35)$$

dove  $\mathcal{N}_N$  è il numero di individui alla  $N$ -esima generazione, e  $\mathcal{N}_0$  è il numero iniziale di individui, è dato da

$$\Lambda(b) = \sum_x p_x \ln(o_x b_x). \quad (36)$$

Questo corrisponde all'equazione (16), con la sola differenza che  $\sum_x 1/o_x$  non è necessariamente uguale a 1. Anzi ci aspettiamo che in condizioni favorevoli si abbia  $\sum_x 1/o_x < 1$ , cioè che la popolazione tenda ad aumentare. Questo corrisponde a un allibratore "generoso" che distribuisce più di quanto ricava dalle scommesse. Possiamo tuttavia definire  $Z = (\sum_x 1/o_x)^{-1}$  e  $q_x = Z/o_x$ , così che  $q_x$  è normalizzato. In questo caso il tasso di crescita per la strategia  $b = (b_x)$  è dato da

$$\Lambda(b) = \ln Z + \sum_x p_x \ln \frac{b_x}{q_x}. \quad (37)$$

Come al solito, il valore ottimale  $b^*$  di  $b$  è dato da

$$b_x^* = p_x, \quad \forall x. \quad (38)$$

Il corrispondente tasso medio di crescita è dato da

$$\Lambda(b^*) = \sum_x p_x \ln(o_x p_x) = \ln Z + \sum_x p_x \ln \frac{p_x}{q_x} = \ln Z + D_{\text{KL}}(p||q). \quad (39)$$

Vediamo che  $\ln Z \geq 0$  è il *minimo* tasso ottimale di crescita, che si ottiene quando  $p_x = q_x, \forall x$ . Questo corrisponde a un processo ambientale "perverso", che minimizza il benessere degli organismi, mentre in generale si possono avere tassi più elevati se le probabilità  $p_x$  dell'ambiente non coincidono con le probabilità  $q_x$  definite dalla distribuzione della fitness.

Consideriamo adesso le probabilità ambientali  $p_x$  fissate, e valutiamo il massimo tasso di crescita possibile, che si potrebbe ottenere se l'organismo fosse in grado di individuare con certezza l'ambiente in cui si trova, e produrre prole tutta con il fenotipo adatto. Otteniamo

$$\begin{aligned} \Lambda_{\text{max}} &= \ln Z - \sum_x p_x \ln q_x = \ln Z + \sum_x p_x \ln \frac{p_x}{q_x} - \sum_x p_x \ln p_x \\ &= \ln Z + D_{\text{KL}}(p||q) + H(X). \end{aligned} \quad (40)$$

Otteniamo così

$$\Lambda_{\text{max}} - \Lambda(b^*) = H(X), \quad (41)$$

che mostra che il tasso di crescita viene ridotto dall'entropia del processo ambientale.

Tuttavia può non essere possibile per l'organismo di aggiustare esattamente la frequenza  $b_x$  del fenotipo  $x$  alla frequenza  $p_x$  dell'ambiente  $x$  nel processo ambientale. Ci aspettiamo che in questo caso il tasso di crescita della popolazione sia più piccolo che nel caso ottimale. Un semplice calcolo dà infatti

$$\Lambda(b^*) - \Lambda(b) = D_{\text{KL}}(p||b). \quad (42)$$

Questo mostra che una valutazione accurata del processo ambientale permette una crescita più vigorosa.

## 7 Il valore di fitness della percezione

Da quanto detto nel paragrafo precedente, risulta evidente che gli organismi possono trarre vantaggio da informazioni tempestive sullo stato dell'ambiente. Supponiamo che sia disponibile un segnale  $Y$ , accessibile a tutti gli individui, tale che la probabilità che l'ambiente si trovi ad avere il valore  $x$  quando il segnale vale  $y$  è data da  $p_{x|y}$ . Il segnale gioca il ruolo del telegrafo nel caso delle corse di cavalli. Il tasso di crescita della popolazione è dato dalla stessa espressione (21) che per le corse. Come in quel caso, la strategia ottimale è una strategia condizionata  $b_{x|y}^*$ , data da

$$b_{x|y}^* = p_{x|y}, \quad \forall x, y. \quad (43)$$

Il tasso corrispondente di crescita della popolazione è dato da

$$\Lambda_Y(b^*) = \sum_{x,y} p_{x,y} \ln(o_x p_{x|y}), \quad (44)$$

dove  $\Lambda_Y(b)$  ricorda che siamo in presenza del segnale  $Y$  correlato all'ambiente. Si ha allora

$$\Lambda_Y(b^*) - \Lambda(b^*) = \sum_{x,y} p_{x,y} \ln \frac{p_{x|y}}{p_x} = \sum_{x,y} p_{x,y} \ln \frac{p_{x,y}}{p_x p_y} = I(X : Y), \quad (45)$$

come nel caso delle scommesse.

A questa analisi si può obiettare che gli individui che l'informazione sull'ambiente non è comune a tutti gli individui, ma è distribuita. In altre parole, che il segnale ricevuto non è necessariamente lo stesso per tutti gli individui, ma differisce (in principio) da individuo a individuo. Indichiamo con  $p_{y|x}$  la probabilità condizionata che l'individuo riceva il segnale  $y$  se l'ambiente è nello stato  $x$ . Allora la probabilità che il fenotipo  $\phi$  sia uguale a  $x$  è data da

$$b_x = \sum_y b_{x|y} p_{y|x}. \quad (46)$$

Il tasso di crescita della popolazione in questo caso (informazione distribuita) è dato da

$$\Lambda_Y(b) = \ln Z + \sum_x p_x \ln \left( \sum_y \frac{b_{x|y} p_{y|x}}{q_x} \right) = \Lambda_{\text{max}} + \sum_x p_x \ln \left( \sum_y b_{x|y} p_{y|x} \right), \quad (47)$$

dove  $\Lambda_{\max} = \ln Z - \sum_x p_x \ln q_x = \langle \ln o_x \rangle$ . Il secondo termine è dato dalla media su  $x$  del logaritmo del valor medio su  $y$  di  $b_{x|y}$ , valutata con la distribuzione condizionata  $p_{y|x}$ . D'altra parte, quando il segnale  $y$  è lo stesso per tutti gli individui (informazione comune), si ottiene

$$\Lambda_Y(b) = \ln Z + \sum_{x,y} p_{x,y} \ln \frac{b_{x|y}}{q_x} = \Lambda_{\max} + \sum_x p_x \sum_y p_{y|x} \ln b_{x|y}. \quad (48)$$

Per ogni valore di  $x$  abbiamo nel caso dell'informazione distribuita il logaritmo della media di  $b_{x|y}$  su  $p_{y|x}$ , e nel caso di informazione comune la media del logaritmo. Ora, data una variabile aleatoria  $W$  a valori positivi, si ha in generale

$$\ln \langle w \rangle \geq \langle \ln w \rangle. \quad (49)$$

Questa relazione è dimostrata nell'appendice C. Abbiamo ottenuto così che ci si aspetta una crescita della popolazione più grande nel caso di informazione distribuita che in quello di informazione comune. Questo risultato è dovuto a Rivoire e Leibler [17], e può essere considerato come un riflesso della "saggezza delle folle" osservata da Galton nel 1907 [18].

## 8 Discussione

Abbiamo visto che l'informazione su un processo aleatorio  $X$  contenuta in un segnale  $Y$ , e misurata dall'informazione mutua  $I(X : Y)$ , è una misura del *valore* dell'informazione stessa in ambienti tanto diversi come le scommesse, l'estrazione di lavoro, e la crescita delle popolazioni. Un ingrediente essenziale di tutte le situazioni è di avere a che fare con un processo  $X$  stazionario, dove  $Y$  fornisce informazioni relative alle successive istanze del processo, e dove i guadagni si accumulano senza interferire (nel caso delle scommesse e della popolazione, il guadagno deve essere misurato su scala logaritmica). Il discorso può essere generalizzato al caso in cui il processo  $X$  ha memoria (e può quindi essere rappresentato da un processo di Markov). In questo caso si pone il problema della memoria del segnale, e del costo e del valore della *predizione* del processo. Alcuni lavori sono stati dedicati a questi aspetti, ma discuterne qui ci porterebbe troppo lontano. Per ora mi accontento di avere richiamato l'attenzione dei lettori sull'inaspettata generalità delle applicazioni di alcuni concetti base della teoria dell'informazione.

## Bibliografia

- [1] Shannon, C. E., A mathematical theory of communication, *Bell System Technical Journal* **27** 379–423 & 623–656 (1948).
- [2] Shannon, C. E., & Weaver, W., *The Mathematical Theory of Communication* (Urbana and Chicago: University of Illinois Press, 1949). Tradotto in: Shannon, C. E. & Weaver, W., *La teoria matematica delle comunicazioni* (Milano: ETAS Compass, 1971).
- [3] Tribus, M., & McIrvine, E. C., Energy and information, *Scientific American* **225**:(3) 179-190 (September 1971).

- [4] Gibbs, J. W., *Elementary Principles in Statistical Mechanics* (New Haven: Yale University Press, 1902).
- [5] Kelly, J. L., A new interpretation of information rate, *Bell System Technical Journal* **35** 917–926 (1956).
- [6] Szilard, L., Über die Entropieverminderung in einem thermodynamischen System bei Eingriffen intelligenter Wesen, *Zeitschrift für Physik* **53** 840–856 (1929). Tradotto in inglese in: Leff, H. S., & Rex, A. F. (eds.), *Maxwell's Demon 2: Entropy, Classical and Quantum Information, Computing* (Bristol and Philadelphia: IoP Press, 2003), p. 110-119.
- [7] Sagawa, T., & Ueda, M., Nonequilibrium thermodynamics of feedback control, *Physical Review E* **85** 021104 (2012).
- [8] Brillouin, L., Maxwell's demon cannot operate: Information and Entropy. I, *Journal of Applied Physics* **22** 334-337 (1951).
- [9] Bennett, C. H., The thermodynamics of computation—A review, *International Journal of Theoretical Physics* **21** 905-940 (1982).
- [10] Landauer, R., Irreversibility and heat generation in the computing process. *IBM Journal of Research and Development* **3** 183-191 (1961).
- [11] Bérut, A., Arakelyan, A., Petrosyan, A., Ciliberto, S., Dillenschneider, R., & Lutz, E., Experimental verification of Landauer's principle linking information and thermodynamics, *Nature* **483** 187-190 (2012).
- [12] Jun, Y., Gavrilov, M., & Bechhoefer, J. High precision test of Landauer's principle in a feedback trap, *Physical Review Letters* **113** 190601 (2014).
- [13] Vinkler, D. A., Permuter, H. H., & Merhav, N., Analogy between gambling and measurement-based work extraction, *Journal of Statistical Mechanics: Theory and Experiment* 043403 (2016).
- [14] Donaldson-Matasci, M. C., Bergstrom, C. T., & Lachmann, M., The fitness value of information, *Oikos* **119** 219–230 (2010).
- [15] Haccou, P., & Iwasa, Y., Optimal mixed strategies in stochastic environments, *Theoretical Population Biology* **47** 212-243 (1995).
- [16] Gollan, B., Grabe, G., Michaux, C., & Helaine, S. Bacterial persists and infection: Past, present, and progressing, *Annual review of microbiology* **73**, 359-385 (2019).
- [17] Rivoire, O., & Leibler, S., The value of information for populations in varying environments, *Journal of Statistical Physics* **142** 1124-1166 (2011).
- [18] Galton, F., Vox Populi, *Nature* **75** 450-451 (1907).

## A Entropia di Shannon

Indichiamo con  $X$  una variabile aleatoria, con valori da 1 a  $N$ , e tale che la probabilità che essa assuma il valore  $x$  è pari a  $p_x$ , dove  $0 \leq p_x$  e  $\sum_x p_x = 1$ . Vogliamo definire una misura  $H(X)$  dell'incertezza associata a  $X$ . Questa misura non deve dipendere dalla natura di  $X$ , ma solo dalla sua distribuzione di probabilità  $p = (p_1, \dots, p_N) = (p_x)$ . Inoltre, chiaramente, non deve dipendere dalla scelta delle etichette  $1, 2, \dots, N$  che assegniamo ai valori di  $X$ . Possiamo richiedere ancora che abbia le seguenti proprietà:

1. Se un solo valore  $x_0$  di  $X$  è possibile, cioè se  $p_x = \delta_{xx_0}$ , dove  $\delta_{xy}$  è la delta di Kronecker, definita da

$$\delta_{xy} = \begin{cases} 1, & \text{se } x = y; \\ 0, & \text{altrimenti;} \end{cases} \quad (50)$$

l'incertezza si deve annullare.

2. Se  $X$  può assumere uno fra  $N$  valori, tutti equiprobabili, l'incertezza  $H(X)$  deve crescere al crescere di  $N$ . Indichiamo con  $f(N)$  l'incertezza associata alla distribuzione  $p_x = 1/N$  per  $x \in \{1, \dots, N\}$ . Si ha allora

$$f(N) > f(M), \quad \text{se } N > M. \quad (51)$$

3. Valori di  $X$  che hanno probabilità nulla non contribuiscono all'incertezza.
4. Se è possibile suddividere gli  $N$  valori di  $X$  in  $K$  gruppi, l'incertezza relativa a  $X$  diventa pari all'incertezza relativa alla variabile  $Y$  che individua il gruppo, più l'incertezza di  $X$  dato  $Y$ , mediata sulla distribuzione di probabilità di  $X$ . Indichiamo con  $H_y(X)$  l'incertezza residua su  $X$ , una volta che sia noto il valore  $y$  di  $Y$ . Si ha allora

$$H(X) = H(Y) + \sum_y p_y H_y(X). \quad (52)$$

5. L'espressione  $H(X)$  intesa come funzione del vettore  $p = (p_x)$  delle probabilità di  $X$ , deve essere continua.

Mostriamo adesso che l'espressione (1) è l'unica espressione che soddisfa queste condizioni, a meno dell'arbitrarietà della costante  $k$ .

Consideriamo una variabile  $X$  che assuma  $N = K \cdot M$  valori, tutti equiprobabili. Possiamo raggruppare i valori di  $X$  in  $K$  gruppi, ciascuno con  $M$  valori. Indicando con  $Y$  la variabile che identifica i gruppi, abbiamo  $H_y(X) = f(M)$ ,  $\forall M$ . Per la (52) si ha allora

$$f(N) = f(K) + K \cdot \frac{1}{K} f(M) = f(K) + f(M). \quad (53)$$

Questo suggerisce di porre

$$f(N) = k \ln N, \quad (54)$$

con  $k > 0$  per la proprietà 2. Mostrare che questa è l'unica soluzione monotona crescente dell'equazione funzionale (53) è un facile esercizio che lascio al lettore. Questa espressione soddisfa chiaramente la proprietà 1.

Consideriamo adesso una variabile  $X$  che assume  $N$  valori, e tale che la probabilità che assuma il valore  $x$  sia pari a  $p_x$ . Data la continuità della  $H(X)$  si commette un errore arbitrariamente piccolo se supponiamo che i valori di  $p_x$  siano tutti razionali. Riportiamoli tutti al comun denominatore  $\mathcal{N}$ , ottenendo

$$p_x = \frac{n_x}{\mathcal{N}}, \quad (55)$$

con  $n_x$  intero,  $\forall x$ . Possiamo ora introdurre una variabile aleatoria  $Y$  che assume  $\mathcal{N}$  valori equiprobabili, divisi in  $N$  gruppi identificati da  $X$ , e tali che il gruppo  $x$  ha  $n_x$  valori equiprobabili. Si ha allora, per la (52),

$$H(Y) = H(X) + \sum_x p_x H_x(Y) = H(X) + \sum_x p_x f(n_x) = H(X) + k \sum_x p_x \ln n_x. \quad (56)$$

Dato che  $H(Y) = k \ln \mathcal{N}$ , otteniamo

$$H(X) = -k \sum_x p_x \ln n_x + k \ln \mathcal{N} = -k \sum_x p_x \ln p_x. \quad (57)$$

Per la proprietà 3., l'espressione  $p_x \ln p_x$  deve essere posta pari a zero se  $p_x = 0$ .

## B Strategie ottimali

Desideriamo trovare l'estremale di  $\Lambda(b)$  dato dalla eq. 16 rispetto a  $b$ , tenendo conto della condizione di normalizzazione  $\sum_x b_x = 1$ . Introduciamo il moltiplicatore di Lagrange  $\lambda$  e ottimizziamo

$$\Psi(b) = \Lambda(b) - \lambda \sum_x b_x. \quad (58)$$

Otteniamo le equazioni

$$\frac{p_x}{b_x} - \lambda = 0, \quad \forall x, \quad (59)$$

che ammettono la soluzione

$$b_x = \frac{p_x}{\lambda}, \quad x. \quad (60)$$

Dato che  $p$  è normalizzato, si deve porre  $\lambda = 1$ .

Allo stesso modo, nel caso in cui è presente un'informazione parallela  $Y$ , si introduce un moltiplicatore di Lagrange  $\lambda_y$  per ogni valore di  $y$ , che imponga il vincolo  $\sum_x b_{x|y} = 1, \forall y$ . In maniera esattamente analoga si ottengono le equazioni

$$\frac{p_{x|y}}{b_{x|y}} = \lambda_y, \quad \forall x, y, \quad (61)$$

che ammettono le soluzioni

$$b_{x|y} = p_{x|y}, \quad \forall x, y, \quad (62)$$

correttamente normalizzate.



## C Disuguaglianze

Sia  $W$  una variabile aleatoria che assume valori  $w$  positivi. Vogliamo dimostrare la disuguaglianza

$$\ln \langle w \rangle \geq \langle \ln w \rangle. \quad (63)$$

Poniamo  $w = e^x$  ed esponenziamo ambo i membri. Otteniamo la disuguaglianza

$$\langle e^x \rangle \geq e^{\langle x \rangle}, \quad (64)$$

Osserviamo che la funzione  $f(x) = e^x$  è tale che  $f''(x) > 0, \forall x$ . Allora, per lo sviluppo di Taylor si ha, dati due punti qualunque  $x_0$  e  $x$ ,

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \int_{x_0}^x dx' (x - x') f''(x'). \quad (65)$$

Per il teorema della media pesata, si ha

$$\int_{x_0}^x dx' (x - x') f''(x') = \frac{1}{2} f''(\bar{x})(x - x_0)^2, \quad (66)$$

dove  $\bar{x}$  è un punto compreso fra  $x_0$  e  $x$ . Otteniamo così

$$f(x) \geq f(x_0) + f'(x_0)(x - x_0). \quad (67)$$

Scegliamo  $x_0 = \langle x \rangle$  e prendiamo la media di ambo i membri. Otteniamo così la disuguaglianza (64), da cui segue l'eq. (63).

Osserviamo adesso che, date due distribuzioni  $p = (p_x)$  e  $q = (q_x)$  della stessa variabile  $X$ , si ha

$$D_{\text{KL}}(p||q) = - \sum_x p_x \ln \frac{q_x}{p_x} = - \left\langle \ln \frac{q}{p} \right\rangle. \quad (68)$$

Per la disuguaglianza appena dimostrata, otteniamo

$$D_{\text{KL}}(p||q) = - \left\langle \ln \frac{q}{p} \right\rangle \geq - \ln \left\langle \frac{q}{p} \right\rangle = - \ln \left( \sum_x p_x \cdot \frac{q_x}{p_x} \right) = - \ln \sum_x q_x = 0, \quad (69)$$

data la normalizzazione di  $q$ . Quindi la divergenza di Kullback-Leibler non può essere negativa.