

Introduction to Information Entropy

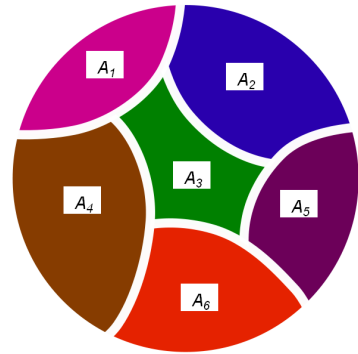
Adapted from Papoulis (1991)

Federico Lombardo

1. INTRODUCTION

Definitions

- The probability $P(A)$ of an event A_i can be interpreted as a measure of our uncertainty about the occurrence or nonoccurrence of A in a single performance of the underlying experiment S (certain event).
- We are interested in assigning a measure of uncertainty to the occurrence or nonoccurrence not of a single event of S , but of any event A_i of a partition \mathbf{A} of S , where a partition is a collection of mutually exclusive events whose union equals S .
- The measure of uncertainty about \mathbf{A} will be denoted by $H(\mathbf{A})$ and will be called the entropy of the partitioning \mathbf{A} .



- In his landmark paper, Shannon (1948) derived the functional $H(\mathbf{A})$ from a number of postulates based on our heuristic understanding of uncertainty. The following is a typical set of such postulates:
 - $H(\mathbf{A})$ is a continuous function of $p_i = P(A_i)$.
 - If $p_1 = \dots = p_N = 1/N$, then $H(\mathbf{A})$ is an increasing function of N .
 - If a new partition \mathbf{B} is formed by subdividing one of the sets of \mathbf{A} , then $H(\mathbf{B}) \geq H(\mathbf{A})$.

- It can be shown that the following sum satisfies these postulates and it is unique within a constant factor:

$$H(\mathbf{A}) = -\sum_{i=1}^N p_i \log p_i$$

- The above assertion can be proven, but here we propose to follow the Papoulis (1991) approach by introducing the above formula as the definition of entropy and developing axiomatically all its properties within the framework of probability.

Shannon, C.E., *A Mathematical Theory of Communication*, Bell System Technical Journal, vol. 27, pp. 379–423, 623–656, July, October, 1948.
 Papoulis, A., *Probability, Random Variables and Stochastic Processes*, 3rd edition, McGraw Hill, 1991. 4

The choice of a logarithmic base corresponds to the choice of a unit for measuring information.

If the base 2 is used the resulting units may be called binary digits, or more briefly *bits*, a word suggested by J. W. Tukey.

A device with two stable positions, such as a relay or a flip-flop circuit, can store one bit of information.

N such devices can store N bits, since the total number of possible states is 2^N and $\log_2 2^N = N$. If the base 10 is used the units may be called decimal digits.

- The applications of entropy can be divided into two categories:
 1. Problems involving the determination of unknown distributions:
 - the available information is in the form of known expected values or other statistical functionals, and the solution is based on the principle of maximum entropy;
 - determine the unknown distributions so as to maximize the entropy $H(\mathbf{A})$ of some partition \mathbf{A} subject to the given constraints.
 2. Coding theory:
 - in this second category, we are given $H(\mathbf{A})$ (source entropy) and we wish to construct various random variables (code lengths) so as to minimize their expected values;
 - the solution involves the construction of optimum mappings (codes) of the random variables under consideration, into the given probability space.

- In the heuristic interpretation of entropy the number $H(\mathbf{A})$ is a measure of our uncertainty about the events A_i of the partition \mathbf{A} prior to the performance of the underlying experiment.
- If the experiment is performed and the results concerning A_i become known, then the uncertainty is removed.
- We can thus say that the experiment provides information about the events A_i equal to the entropy of their partition.
- Thus uncertainty equals information and both are measured by entropy.

1. Determine the entropy of the partition $\mathbf{A} = [\text{even}, \text{odd}]$ in the fair-die experiment. Clearly, $P\{\text{even}\} = P\{\text{odd}\} = 1/2$. Hence

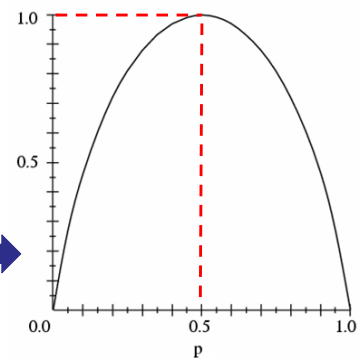
$$H(\mathbf{A}) = -1/2 \log 1/2 - 1/2 \log 1/2 = \log 2$$

2. In the same experiment, \mathbf{G} is the partition consisting of the elementary events $\{f_i\}$. In this case, $P\{f_i\} = 1/6$; hence

$$H(\mathbf{G}) = -\sum_{i=1}^6 P\{f_i\} \log P\{f_i\} = \log 6$$

3. We consider now the coin experiment where $P\{\text{heads}\} = p$. In this case, the entropy of \mathbf{G} equals

$$H(\mathbf{G}) = -p \log p - (1-p) \log(1-p) \equiv r(p) \quad \rightarrow$$



If the die is rolled and we are told which face showed, then we gain information about the partition \mathbf{G} equal to its entropy $\log 6$.

If we are told merely that "even" or "odd" showed, then we gain information about the partition \mathbf{A} equal to its entropy $\log 2$.

In this case, the information gained about the partition \mathbf{G} equals again $\log 2$.

As we shall see, the difference $\log 6 - \log 2 = \log 3$ is the uncertainty about \mathbf{G} assuming \mathbf{A} (conditional entropy).

- An important application of entropy is the determination of the probabilities p_i of the events of a partition \mathbf{A} subject to various constraints, with the maximum entropy (ME) method.
- ME principle states that the unknown p_i 's must be so chosen as to maximize the entropy of \mathbf{A} subject to the given constraints (Jaynes, 1957).
- The ME principle is equivalent to the principle of insufficient reason (Bernoulli, 1713): "*In the absence of any prior knowledge, we must assume that the events A_i have equal probabilities*". This conclusion is based on the subjective interpretation of probability as a measure of our state of knowledge about the events A_i .
- Operationally, the ME method simplifies the analysis drastically when, as is the case in most applications, the constraints are phrased in terms of probabilities in the space S^n of repeated trials (i.e., the resulting product space from the experiment S repeated n times).

Jaynes, E. T., *Information Theory and Statistical Mechanics*, Physical Review, 106(4), 620–630, 1957.
Bernoulli, J., *Ars Conjectandi*, 1713.

8

Classical Definition

The classical definition of probability was introduced as a consequence of the principle of insufficient reason.

According to the classical definition, the probability $P(A)$ of an event A is determined a priori without actual experimentation. It is given by the ratio

$$P(A) = N_A/N$$

where N is the number of possible outcomes and N_A is the number of outcomes that are favorable to the event A , provided that all outcomes are equally likely.

1. Determine the probabilities p_i of the six faces of a die, having access to no prior information. The ME principle states that the p_i 's must be such as to maximize the sum

$$H(\mathbf{G}) = -p_1 \log p_1 - \dots - p_6 \log p_6$$

Since $p_1 + \dots + p_6 = 1$, this yields $p_1 = \dots = p_6 = 1/6$, in agreement with the classical definition.

2. A player places a bet of one euro on "odd" and he wins, on the average, 20 cents per game. We wish again to determine the p_i 's using the ME method; however, now we must satisfy the constraints

$$p_1 + p_3 + p_5 = 0.6 \qquad p_2 + p_4 + p_6 = 0.4$$

This is a consequence of the available information because an average gain of 20 cents means that $P\{\text{odd}\} - P\{\text{even}\} = 0.2$. Maximizing $H(\mathbf{G})$ subject to the above constraints, we obtain

$$p_1 = p_3 = p_5 = 0.2 \qquad p_2 = p_4 = p_6 = 0.133$$

The result of the second experiment agrees again with the classical definition if we apply the principle of insufficient reason to the outcomes of the events {odd} and {even} separately.

- The ME method is thus a valuable tool in the solution of applied problems. It can be used, in fact, even in deterministic problems involving the estimation of unknown parameters from insufficient data.
- We should emphasize, however, that as in the case of the classical definition of probability, the conclusions drawn from the ME method must be accepted with skepticism particularly when they involve elaborate constraints.
- Concerning the previous examples, we conclude that all p_i 's must be equal in the absence of prior constraints, which is not in conflict with our experience concerning dice. The second conclusion, however, is not as convincing, we would think, even though we have no basis for any other conclusion.
- One might argue that this apparent conflict between the ME method and our experience is due to the fact that we did not make total use of our prior knowledge. This might be true; however, it is not always clear how such constraints can be phrased analytically and, even if they can, how complex the required computations might be.

1. **Axiomatic:** $P(A)$ is a number assigned to the event A . This number satisfies the following three postulates but is otherwise arbitrary
 - The probability of an event A is a positive number, $P(A) \geq 0$
 - The probability of the certain event S equals 1, $P(S) = 1$
 - If the events A and B are mutually exclusive, $P(A + B) = P(A) + P(B)$
2. **Empirical:** For large n , $P(A) \approx k/n$, where k is the number of times A occurs in n repetitions of the underlying experiment S .
3. **Subjective:** $P(A)$ is a measure of our uncertainty about the occurrence of A in a single performance of S .
4. **Principle of insufficient reason:** If A_i are N events of a partition \mathbf{A} of S and nothing is known about their probabilities, then $P(A_i) = 1/N$.

1. **Axiomatic:** $H(\mathbf{A})$ is a number assigned to each partition $\mathbf{A} = [A_1, \dots, A_N]$ of S . This number equals the sum $-\sum p_i \ln p_i$, where $p_i = P(A_i)$ and $i = 1, \dots, N$
2. **Empirical:** This interpretation involves the repeated performance not of the experiment S , but of the experiment S^n of repeated trials. In this experiment, each specific typical sequence $\mathbf{t}_j = \{A_i \text{ occurs } n_i \approx np_i \text{ times in a specific order } j\}$ is an event with probability

$$P(\mathbf{t}_j) = p_1^{n_1} \cdots p_N^{n_N} \approx e^{np_1 \ln p_1 + \cdots + np_N \ln p_N} = e^{-nH(\mathbf{A})}$$

Applying the relative frequency interpretation of probability to this event, we conclude that if the experiment S^n is repeated m times and the event \mathbf{t}_j occurs m_j times, then for sufficiently large m ,

$$P(\mathbf{t}_j) = e^{-nH(\mathbf{A})} \approx \frac{m_j}{m}; \quad \text{hence} \quad H(\mathbf{A}) \approx -\frac{1}{n} \ln \frac{m_j}{m}$$

This relates the theoretical quantity $H(\mathbf{A})$ to the experimental numbers m_j and m .

3. **Subjective:** $H(\mathbf{A})$ is a measure of our uncertainty about the occurrence of the events A_i of the partition \mathbf{A} in a single performance of S .

4. **Principle of maximum entropy:** The probabilities $p_i = P(A_i)$ must be such as to maximize $H(\mathbf{A})$ subject to the given constraints. Since it can be demonstrated that the number of typical sequences is $n_t = e^{nH(\mathbf{A})}$, the ME principle is equivalent to the principle of maximizing n_t . If there are no constraints, that is, if nothing is known about the probabilities p_i , then the ME principle leads to the estimates $p_i = 1/N$, $H(\mathbf{A}) = \ln N$, and $n_t = N^n$.

2. BASIC CONCEPTS

Conditional Entropy

- The entropy $H(\mathbf{A})$ of a partition $\mathbf{A} = [A_i]$ gives us a measure of uncertainty about the occurrence of the events A_i at a given trial.
- If in the definition of entropy we replace the probabilities $P(A_i)$ by the conditional probabilities $P(A_i|M)$, we obtain the conditional entropy $H(\mathbf{A}|M)$ of \mathbf{A} assuming M

$$H(\mathbf{A}|M) = - \sum_i P(A_i|M) \log P(A_i|M)$$

- From this it follows that if at a given trial we know that M occurs, then our uncertainty about \mathbf{A} equals $H(\mathbf{A}|M)$.
- If we know that the complement M^c of M occurs, then our uncertainty equals $H(\mathbf{A}|M^c)$.
- Assuming that the binary partition $\mathbf{M} = [M, M^c]$ is observed, the uncertainty per trial about \mathbf{A} is given by the weighted sum

$$H(\mathbf{A}|\mathbf{M}) = P(M)H(\mathbf{A}|M) + P(M^c)H(\mathbf{A}|M^c)$$

The expression a partition \mathbf{B} is observed will mean that we know which of the events of \mathbf{B} has occurred.

Mutual Information

- If at each trial we observe the partition $\mathbf{B} = [B_j]$, then we show that the uncertainty per trial about \mathbf{A} equals $H(\mathbf{A}|\mathbf{B})$
- Indeed, in a sequence of n trials, the number of times the event B_j occurs equals $n_j \approx nP(B_j)$; in this subsequence, the uncertainty about \mathbf{A} equals $H(\mathbf{A}|B_j)$ per trial. Hence, the total uncertainty about \mathbf{A} equals

$$\sum_j n_j H(\mathbf{A}|B_j) \approx \sum_j n P(B_j) H(\mathbf{A}|B_j) = n H(\mathbf{A}|\mathbf{B})$$

and the uncertainty per trial equals $H(\mathbf{A}|\mathbf{B})$

- Thus the observation of \mathbf{B} reduces the uncertainty about \mathbf{A} from $H(\mathbf{A})$ to $H(\mathbf{A}|\mathbf{B})$. The mutual information

$$I(\mathbf{A}, \mathbf{B}) = H(\mathbf{A}) - H(\mathbf{A}|\mathbf{B})$$

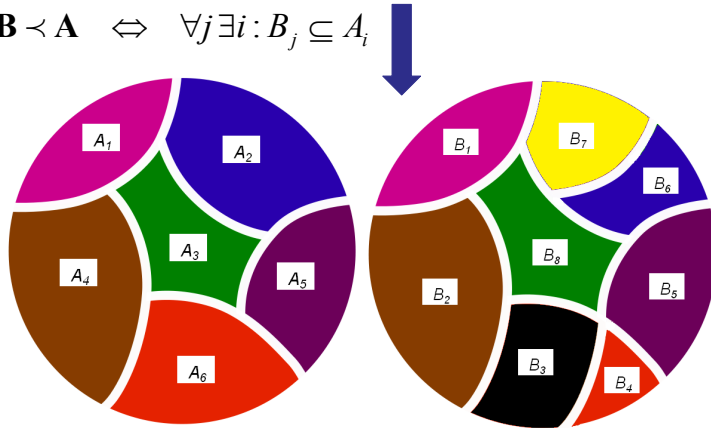
is the reduction of the uncertainty about \mathbf{A} resulting from the observation of \mathbf{B} . $I(\mathbf{A}, \mathbf{B})$ can be interpreted as the information about \mathbf{A} contained in \mathbf{B} .

Definitions

1. A partition whose elements are the elementary events $\{\zeta_j\}$ of the space S will be denoted by \mathbf{G} and will be called the element partition.
2. A refinement of a partition \mathbf{A} is a partition \mathbf{B} such that each element B_j of \mathbf{B} is a subset of some element A_i of \mathbf{A} . We shall use the following notation:

$$\mathbf{B} \prec \mathbf{A} \Leftrightarrow \forall j \exists i : B_j \subseteq A_i$$

3. The product of two partitions \mathbf{A} and \mathbf{B} is a partition whose elements are all intersections $A_i \cap B_j$ of the elements of \mathbf{A} and \mathbf{B} . This new partition is denoted by $\mathbf{A} \cdot \mathbf{B}$.



Refinement = rifinitura

Considerations

1. If **B** is a refinement of **A**, it can be shown that $H(\mathbf{A}) \leq H(\mathbf{B})$.
Then, for any **A** we have $H(\mathbf{A}) \leq H(\mathbf{G})$, where **G** is the element partition.
2. If **B** is a refinement of **A** and **B** is observed, then we know which of the events of **A** occurred. Hence $H(\mathbf{A}|\mathbf{B}) = 0$.
3. Thus, for any **A** we have $H(\mathbf{A}|\mathbf{G}) = 0$.
4. For any **A** and **B**, we have that $H(\mathbf{A} \cdot \mathbf{B}) \geq H(\mathbf{A})$ and $H(\mathbf{A} \cdot \mathbf{B}) \geq H(\mathbf{B})$, because $\mathbf{A} \cdot \mathbf{B}$ is a refinement of both **A** and **B**.
5. If the partitions **A** and **B** are independent (i.e., their events are all independent of each other) and **B** is observed, then no information about **A** is gained. Hence $H(\mathbf{A}|\mathbf{B}) = H(\mathbf{A})$.

6. If we observe **B**, our uncertainty about **A** cannot increase.
Hence $H(\mathbf{A}|\mathbf{B}) \leq H(\mathbf{A})$.
7. To observe $\mathbf{A} \cdot \mathbf{B}$, we must observe **A** and **B**. If only **B** is observed, the information gained equals $H(\mathbf{B})$. Therefore, the uncertainty about **A** assuming **B**, equals the remaining uncertainty, $H(\mathbf{A}|\mathbf{B}) = H(\mathbf{A} \cdot \mathbf{B}) - H(\mathbf{B})$.
8. Combining 6 and 7, we conclude that $H(\mathbf{A} \cdot \mathbf{B}) \leq H(\mathbf{A}) + H(\mathbf{B})$.
9. If **B** is observed, then the information that is gained about **A** equals $I(\mathbf{A}, \mathbf{B})$.
 - If **B** is a refinement of **C** and **B** is observed, then **C** is known.
 - But knowledge of **C** yields information about **A** equal to $I(\mathbf{A}, \mathbf{C})$.
 - Hence, if **B** is a refinement of **C**, then $I(\mathbf{A}, \mathbf{B}) \geq I(\mathbf{A}, \mathbf{C})$.
 - Equivalently, we have also that $H(\mathbf{A}|\mathbf{B}) \leq H(\mathbf{A}|\mathbf{C})$.

3. RANDOM VARIABLES AND STOCHASTIC PROCESSES

- We are given an experiment specified by the space S , the field of subsets of S called events, and the probability assigned to these events.
- To every outcome ζ of this experiment, we assign a number $\mathbf{x}(\zeta)$. We have thus created a function \mathbf{x} with domain the set S and range a set of numbers. This function is called random variable (RV) if it satisfies the following conditions but is otherwise arbitrary:
 - The set of experimental outcomes $\{\mathbf{x} \leq x\}$ is an event for every x .
 - The probabilities of the events $\{\mathbf{x} = \infty\}$ and $\{\mathbf{x} = -\infty\}$ equal 0.
- The elements of the set S that are contained in the event $\{\mathbf{x} \leq x\}$ change as the number x takes various values. The probability $P\{\mathbf{x} \leq x\}$ is, therefore, a number that depends on x .
- This number is denoted by $F(x)$ and is called the cumulative distribution function (CDF) of the RV \mathbf{x} : $F(x) = P\{\mathbf{x} \leq x\}$

The set of experimental outcomes $\{\mathbf{x} \leq x\}$ can be interpreted as follows. Given an arbitrary number x , we find all numbers $\mathbf{x}(\zeta_i)$ that do not exceed x . Then, the corresponding experimental outcomes ζ_i form the set $\{\mathbf{x} \leq x\}$.

- The RV \mathbf{x} is of continuous type if its CDF $F(x)$ is continuous. In this case, we have: $P\{\mathbf{x} = x\} = 0$.
- The RV \mathbf{x} is of discrete type if its CDF $F(x)$ is a staircase function. Denoting by x_i the discontinuity points of $F(x)$, we have: $P\{\mathbf{x} = x_i\} = p_i$.
- The derivative $f(x)$ of $F(x)$ is called the probability density function (PDF) of the RV \mathbf{x}

$$f(x) = \frac{dF(x)}{dx}$$

- If the RV \mathbf{x} is of discrete type taking the values x_i with probabilities p_i , then

$$f(x) = \sum_i p_i \delta(x - x_i) \quad p_i = P\{\mathbf{x} = x_i\}$$

where $\delta(x)$ is the impulse function. The term $p_i \delta(x - x_i)$ can be shown as a vertical arrow at $x = x_i$ with length equal to p_i .

- Entropy is a number assigned to a partition. To define the entropy of an RV we must, therefore, form a suitable partition.
- This is simple if the RV is of discrete type. However, for continuous-type RVs we can do so only indirectly.
- Suppose that the RV \mathbf{x} is of discrete type taking the values x_i with probabilities $P\{\mathbf{x} = x_i\} = p_i$.
 - The events $\{\mathbf{x} = x_i\}$ are mutually exclusive and their union is the certain event; hence they form a partition.
 - This partition will be denoted by $\mathbf{A}_{\mathbf{x}}$ and will be called the partition of \mathbf{x} .
- *Definition:* The entropy $H(\mathbf{x})$ of a discrete-type RV \mathbf{x} is the entropy $H(\mathbf{A}_{\mathbf{x}})$ of its partition $\mathbf{A}_{\mathbf{x}}$:

$$H(\mathbf{x}) = H(\mathbf{A}_{\mathbf{x}}) = -\sum_i p_i \ln p_i$$

- The entropy of a continuous-type RV cannot be so defined because the events $\{\mathbf{x} = x_i\}$ do not form a partition (they are not countable).
- To define $H(\mathbf{x})$, we form, first, the discrete-type RV \mathbf{x}_δ obtained by rounding off \mathbf{x} , so as to make it a staircase function: $\mathbf{x}_\delta = n\delta$ if $n\delta - \delta < \mathbf{x} \leq n\delta$, hence

$$P(\mathbf{x}_\delta = n\delta) = P(n\delta - \delta < \mathbf{x} \leq n\delta) = \int_{n\delta - \delta}^{n\delta} f(x) dx = \delta \bar{f}(n\delta)$$

where $\bar{f}(n\delta)$ is a number between the maximum and the minimum of $f(x)$ in the interval $(n\delta - \delta, n\delta)$.

- Applying the definition of the entropy of a discrete-type RV to \mathbf{x}_δ we obtain

$$H(\mathbf{x}_\delta) = - \sum_{n=-\infty}^{\infty} \delta \bar{f}(n\delta) \ln[\delta \bar{f}(n\delta)]$$

where:

$$\sum_{n=-\infty}^{\infty} \delta \bar{f}(n\delta) = \int_{-\infty}^{\infty} f(x) dx = 1$$

Continuous type

- After algebraic manipulations, we conclude that

$$H(\mathbf{x}_\delta) = -\ln \delta - \sum_{n=-\infty}^{\infty} \delta \bar{f}(n\delta) \ln \bar{f}(n\delta)$$

- As $\delta \rightarrow 0$, the RV $\mathbf{x}_\delta \rightarrow \mathbf{x}$, but its entropy $H(\mathbf{x}_\delta) \rightarrow \infty$ because: $-\ln \delta \rightarrow \infty$.
- For this reason, we define the entropy $H(\mathbf{x})$ of \mathbf{x} not as the limit of $H(\mathbf{x}_\delta)$ but as the limit of the sum: $H(\mathbf{x}_\delta) + \ln \delta$, as $\delta \rightarrow 0$. This yields

$$H(\mathbf{x}_\delta) + \ln \delta \xrightarrow{\delta \rightarrow 0} -\int_{-\infty}^{\infty} f(x) \ln f(x) dx$$

- Definition: The entropy of a continuous-type RV \mathbf{x} is by definition the integral

$$H(\mathbf{x}) = -\int_{-\infty}^{\infty} f(x) \ln f(x) dx$$

- Example: If \mathbf{x} is uniform in the interval $(0, a)$, where $f(x) = 1/a$, then

$$H(\mathbf{x}) = -\frac{1}{a} \int_0^a \ln \frac{1}{a} dx = \ln a$$

Considerations

- The entropy $H(\mathbf{x}_\delta)$ of \mathbf{x}_δ is a measure of our uncertainty about the RV \mathbf{x} rounded off to the nearest $n\delta$. If δ is small, the resulting uncertainty is large and it tends to ∞ as $\delta \rightarrow 0$.
- This conclusion is based on the assumption that \mathbf{x} can be observed perfectly; that is, its various values can be recognized as distinct no matter how close they are.
- In a physical experiment, however, this assumption is not realistic. Values of \mathbf{x} that differ slightly cannot always be treated as distinct (noise considerations or round-off errors, for example).
- Accounting for the term $\ln \delta$ in the definition of entropy of a continuous-type RV \mathbf{x} is, in a sense, a recognition of this ambiguity.

- As in the case of arbitrary partitions, the entropy of a discrete-type RV \mathbf{x} is positive and it is used as a measure of uncertainty about \mathbf{x} .
- This is not so, however, for continuous-type RVs. Their entropy can take any value from $-\infty$ to ∞ and it is used to measure only changes in uncertainty.
- The various properties of partitions also apply to continuous-type RVs if, as is generally the case, they involve only differences of entropies.

- The entropy of a continuous-type RV \mathbf{x} can be expressed as the expected value of the RV $\mathbf{y} = -\ln f(\mathbf{x})$:

$$H(\mathbf{x}) = E\{-\ln f(\mathbf{x})\} = -\int_{-\infty}^{\infty} f(x) \ln f(x) dx$$

- Similarly, the entropy of a discrete-type RV \mathbf{x} can be written as the expected value of the RV $-\ln p(\mathbf{x})$:

$$H(\mathbf{x}) = E\{-\ln p(\mathbf{x})\} = -\sum_i p_i \ln p_i$$

where now $p(x)$ is a function defined only for $x = x_i$ and such that $p(x_i) = p_i$.

- If the RV \mathbf{x} is *exponentially* distributed, then $f(x) = \lambda e^{-\lambda x} U(x)$, where $U(x)$ is the Heaviside step function. Hence:

$$H(\mathbf{x}) = E\{-\ln f(\mathbf{x})\} = 1 - \ln \lambda = \ln \frac{e}{\lambda}$$

- If the RV \mathbf{x} is *normally* distributed, then

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad H(\mathbf{x}) = E\{-\ln f(\mathbf{x})\} = \ln(\sigma\sqrt{2\pi}e)$$

Functions of one RV

Suppose that \mathbf{x} is an RV and $g(x)$ is a function of the real variable x .

The expression $\mathbf{y} = g(\mathbf{x})$ is a new RV defined as follows: For a given ζ , $\mathbf{x}(\zeta)$ is a number and $g[\mathbf{x}(\zeta)]$ is another number specified in terms of $\mathbf{x}(\zeta)$ and $g(x)$.

This number is the value $\mathbf{y}(\zeta) = g[\mathbf{x}(\zeta)]$ with the domain set S of experimental outcomes.

The distribution function $F(y)$ of the RV so formed is the probability of the event $\{\mathbf{y} \leq y\}$ consisting of all outcomes ζ such that $\mathbf{y}(\zeta) = g[\mathbf{x}(\zeta)] \leq y$. Thus

$$F(y) = P\{\mathbf{y} \leq y\} = P\{g(\mathbf{x}) \leq y\}$$

For a specific y , the values of x such that $g(x) \leq y$ form a set on the x axis denoted by R_y . Clearly, $g(\mathbf{x}) \leq y$ if $\mathbf{x}(\zeta)$ is a number in the set R_y .

The above leads to the conclusion that for $g(\mathbf{x})$ to be an RV, the function $g(x)$ must have the following properties:

- Its domain must include the range of the RV \mathbf{x} .
- It must be a *Baire* function, that is, for every y , the set R_y such that $g(x) \leq y$ must consist of the union and intersection of a countable number of intervals. Only then $\{\mathbf{y} \leq y\}$ is an event.
- The events $\{g(\mathbf{x}) = \pm \infty\}$ must have zero probability.

- Suppose that \mathbf{x} and \mathbf{y} are two discrete-type RVs taking the values x_i and y_j respectively with $P\{\mathbf{x} = x_i, \mathbf{y} = y_j\} = p_{ij}$.
- Their joint entropy, denoted by $H(\mathbf{x}, \mathbf{y})$, is by definition the entropy of the product of their respective partitions. Clearly, the elements of $\mathbf{A}_x \cdot \mathbf{A}_y$ are the events $\{\mathbf{x} = x_i, \mathbf{y} = y_j\}$. Hence

$$H(\mathbf{x}, \mathbf{y}) = H(\mathbf{A}_x \cdot \mathbf{A}_y) = - \sum_{i,j} p_{ij} \ln p_{ij}$$

- The above can be written as an expected value:

$$H(\mathbf{x}, \mathbf{y}) = E\{-\ln p(\mathbf{x}, \mathbf{y})\}$$

where $p(x, y)$ is a function defined only for $x = x_i$ and $y = y_j$ and it is such that $p(x_i, y_j) = p_{ij}$.

- The joint entropy $H(\mathbf{x}, \mathbf{y})$ of two continuous-type RVs \mathbf{x} and \mathbf{y} is defined as the limit of the sum: $H(\mathbf{x}_\delta, \mathbf{y}_\delta) + 2 \ln \delta$, as $\delta \rightarrow 0$, where \mathbf{x}_δ and \mathbf{y}_δ are their staircase approximation. Thus we have:

$$H(\mathbf{x}, \mathbf{y}) = - \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) \ln f(x, y) dx dy = E\{-\ln f(\mathbf{x}, \mathbf{y})\}$$

Conditional Entropy

- Consider two discrete-type RVs \mathbf{x} and \mathbf{y} taking the values x_i and y_j respectively with

$$P(\mathbf{x} = x_i | \mathbf{y} = y_j) = \pi_{ji} = p_{ji} / p_j$$

- The conditional entropy $H(\mathbf{x}|y_j)$ of \mathbf{x} assuming $\mathbf{y} = y_j$ is by definition the conditional entropy of the partition \mathbf{A}_x of \mathbf{x} assuming $\{\mathbf{y} = y_j\}$. From the above it follows that:

$$H(\mathbf{x}|y_j) = -\sum_i \pi_{ji} \ln \pi_{ji}$$

- The conditional entropy $H(\mathbf{x}|\mathbf{y})$ of \mathbf{x} assuming \mathbf{y} is the conditional entropy of \mathbf{A}_x assuming \mathbf{A}_y . Thus

$$H(\mathbf{x}|\mathbf{y}) = -\sum_j p_j H(\mathbf{x}|y_j) = -\sum_{i,j} p_{ji} \ln \pi_{ji}$$

- For continuous-type RVs the corresponding concepts are defined similarly

$$H(\mathbf{x}|\mathbf{y}) = -\int_{-\infty}^{\infty} f(x|\mathbf{y}) \ln f(x|\mathbf{y}) dx = E\{-\ln f(\mathbf{x}|\mathbf{y}) | \mathbf{y} = y\}$$

$$H(\mathbf{x}|\mathbf{y}) = -\int_{-\infty}^{\infty} f(y) H(\mathbf{x}|y) dy = -\int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) \ln f(x|\mathbf{y}) dx dy = E\{-\ln f(\mathbf{x}|\mathbf{y})\}$$

The last equation can be easily demonstrated by accounting for the following:

$$f(x, y) = f(x|\mathbf{y})f(y)$$

- We shall define the mutual information of the RVs \mathbf{x} and \mathbf{y} as follows

$$I(\mathbf{x}, \mathbf{y}) = H(\mathbf{x}) + H(\mathbf{y}) - H(\mathbf{x}, \mathbf{y})$$

- $I(\mathbf{x}, \mathbf{y})$ can be written as an expected value

$$I(\mathbf{x}, \mathbf{y}) = E \left\{ \ln \frac{f(\mathbf{x}, \mathbf{y})}{f(\mathbf{x})f(\mathbf{y})} \right\}$$

- Since $f(x, y) = f(x|y)f(y)$ it follows from the above that

$$I(\mathbf{x}, \mathbf{y}) = H(\mathbf{x}) - H(\mathbf{x}|\mathbf{y}) = H(\mathbf{y}) - H(\mathbf{y}|\mathbf{x})$$

- The properties of entropy, developed before for arbitrary partitions, are obviously true for the entropy of discrete-type RVs and can be established as appropriate limits for continuous-type RVs.

- We shall compare the entropy of the RVs \mathbf{x} and $\mathbf{y} = g(\mathbf{x})$.
- If the RV \mathbf{x} is of discrete type, then $H(\mathbf{y}) \leq H(\mathbf{x})$
with equality if and only if the transformation $y = g(x)$ has a unique inverse $x = g^{(-1)}(y)$.
- If the transformation $y = g(x)$ has not a unique inverse (it is not one-to-one), then $\mathbf{y} = y_i$ for more than one value of \mathbf{x} . This results in a reduction of $H(\mathbf{x})$.
- If the RV \mathbf{x} is of continuous type, then $H(\mathbf{y}) \leq H(\mathbf{x}) + E\{\ln|g'(\mathbf{x})|\}$
where $g'(x)$ is the derivative of $g(x)$. The equality holds if and only if the transformation $y = g(x)$ has a unique inverse.
- Similarly, if $\mathbf{y}_i = g_i(\mathbf{x}_1, \dots, \mathbf{x}_n)$, $i = 1, \dots, n$, are n functions of the RVs \mathbf{x}_i , then

$$H(\mathbf{y}_1, \dots, \mathbf{y}_n) \leq H(\mathbf{x}_1, \dots, \mathbf{x}_n) + E\{\ln|J(\mathbf{x}_1, \dots, \mathbf{x}_n)|\}$$

In the last equation $J(x_1, \dots, x_n)$ is the jacobian of the transformation $y_i = g_i(x_1, \dots, x_n)$.

Linear transformations

Suppose that $\mathbf{y}_i = a_{i1}\mathbf{x}_1 + \dots + a_{in}\mathbf{x}_n$; denoting by Δ the determinant of the coefficients, we conclude that if $\Delta \neq 0$ then

$$H(\mathbf{y}_1, \dots, \mathbf{y}_n) = H(\mathbf{x}_1, \dots, \mathbf{x}_n) + \ln|\Delta|$$

because the transformation has a unique inverse and Δ does not depend on \mathbf{x}_i

- The statistics of most stochastic processes are determined in terms of the joint density $f(x_1, \dots, x_m)$ of the RVs $\mathbf{x}(t_1), \dots, \mathbf{x}(t_m)$.
- The joint entropy of these RVs is the m th-order entropy of the process $\mathbf{x}(t)$

$$H(\mathbf{x}_1, \dots, \mathbf{x}_m) = E\{-\ln f(\mathbf{x}_1, \dots, \mathbf{x}_m)\}$$
- This function equals the uncertainty about the above RVs and it equals the information gained when they are observed.
- In general, the uncertainty about the values of $\mathbf{x}(t)$ on the entire t axis or even on a finite interval, no matter how small, is infinite.
- However, we assume $\mathbf{x}(t)$ expressed in terms of its values on a countable set of points, then a rate of uncertainty can be introduced. It suffices, therefore, to consider only discrete-time processes \mathbf{x}_n .

Stochastic Processes

As we recall, an RV \mathbf{x} is a rule for assigning to every outcome ζ of an experiment S a number $\mathbf{x}(\zeta)$.

A stochastic process $\mathbf{x}(t)$ is a rule for assigning to every ζ a function $\mathbf{x}(t, \zeta)$.

Thus a stochastic process is a family of time functions depending on the parameter ζ or, equivalently, a function of t and ζ .

The domain of ζ is the set of all experimental outcomes and the domain of t is a set R of real numbers.

If R is the real axis, then $\mathbf{x}(t)$ is a continuous-time process. If R is the set of integers, then $\mathbf{x}(t)$ is a discrete-time process.

A discrete-time process is, thus, a sequence of random variables. Such a sequence will be denoted by \mathbf{x}_n .

We shall say that $\mathbf{x}(t)$ is a discrete-state process if its values are countable. Otherwise, it is a continuous-state process.

We shall use the notation $\mathbf{x}(t)$ to represent a stochastic process omitting, as in the case of random variables, its dependence on ζ .

Thus $\mathbf{x}(t)$ has the following interpretations:

1. It is a family (or an ensemble) of functions $\mathbf{x}(t, \zeta)$. In this interpretation, t and ζ are variables.
2. It is a single time function (or a sample of the given process). In this case, t is a variable and ζ is fixed.
3. If t is fixed and ζ is variable, then $\mathbf{x}(t)$ is a random variable equal to the state of the given process at time t .
4. If t and ζ are fixed, then $\mathbf{x}(t)$ is a number.

- The m th-order entropy of a discrete-time process \mathbf{x}_n is the joint entropy $H(\mathbf{x}_1, \dots, \mathbf{x}_m)$ of the m RVs: $\mathbf{x}_n, \mathbf{x}_{n-1}, \dots, \mathbf{x}_{n-m+1}$
- We shall assume throughout that the process \mathbf{x}_n is strict-sense stationary (SSS). In this case, $H(\mathbf{x}_1, \dots, \mathbf{x}_m)$ is the uncertainty about any m consecutive values of the process \mathbf{x}_n .
- The first-order entropy will be denoted by $H(\mathbf{x})$ and equals the uncertainty about \mathbf{x}_n for a specific n .

- Recalling the properties of entropy, we have:

$$H(\mathbf{x}_1, \dots, \mathbf{x}_m) \leq H(\mathbf{x}_1) + \dots + H(\mathbf{x}_m) = m H(\mathbf{x})$$

- Example: If the process \mathbf{x}_n is *strictly white*, that is, if the RVs $\mathbf{x}_n, \mathbf{x}_{n-1}, \dots$ are independent, then $H(\mathbf{x}_1, \dots, \mathbf{x}_m) = m H(\mathbf{x})$

Stationary Stochastic Processes

A stochastic process $\mathbf{x}(t)$ is called strict-sense stationary (abbreviated SSS) if its statistical properties are invariant to a shift of the origin.

This means that the processes $\mathbf{x}(t)$ and $\mathbf{x}(t + c)$ have the same statistics for any c .

A stochastic process $\mathbf{x}(t)$ is called wide-sense stationary (abbreviated WSS) if:

- Its mean is constant: $E\{\mathbf{x}(t)\} = \eta$
- Its autocorrelation depends only on $\tau = t_1 - t_2$: $E\{\mathbf{x}(t + \tau) \mathbf{x}(t)\} = R(\tau)$

- The conditional entropy of order m , $H(\mathbf{x}_n | \mathbf{x}_{n-1}, \dots, \mathbf{x}_{n-m})$, of a process \mathbf{x}_n is the uncertainty about its present under the assumption that its m most recent values have been observed.
- Recalling that $H(\mathbf{x}|y) \leq H(\mathbf{x})$, we can readily show that:

$$H(\mathbf{x}_n | \mathbf{x}_{n-1}, \dots, \mathbf{x}_{n-m}) \leq H(\mathbf{x}_n | \mathbf{x}_{n-1}, \dots, \mathbf{x}_{n-m-1})$$

- Thus the above conditional entropy is a decreasing function of m . If, therefore, it is bounded from below, it tends to a limit. This is certainly the case if the RVs \mathbf{x}_n are of discrete type because then all entropies are positive.
- The limit will be denoted by $H_c(\mathbf{x})$ and will be called the conditional entropy of the process \mathbf{x}_n :

$$H_c(\mathbf{x}) = \lim_{m \rightarrow \infty} H(\mathbf{x}_n | \mathbf{x}_{n-1}, \dots, \mathbf{x}_{n-m})$$
- The function $H_c(\mathbf{x})$ is a measure of our uncertainty about the present of \mathbf{x}_n under the assumption that its entire past is observed.

Example

If \mathbf{x}_n is *strictly white*, then: $H_c(\mathbf{x}) = H(\mathbf{x})$

Entropy Rate

The ratio $H(\mathbf{x}_1, \dots, \mathbf{x}_m)/m$ is the average uncertainty per sample in a block of m consecutive samples.

The limit of this average as $m \rightarrow \infty$ will be denoted by $\underline{H}(\mathbf{x})$ and will be called the entropy rate of the process \mathbf{x}_n .

It can be shown that the entropy rate of a process \mathbf{x}_n equals its conditional entropy:

$$\underline{H}(\mathbf{x}) = H_c(\mathbf{x})$$

4. MAXIMUM ENTROPY METHOD

- The ME method is used to determine various parameters of a probability space subject to given constraints.
- The resulting problem can be solved, in general, only numerically and it involves the evaluation of the maximum of a function of several variables.
- In a number of important cases, however, the solution can be found analytically or it can be reduced to a system of algebraic equations.
- We consider herein certain special cases, concentrating on constraints in the form of expected values.
- For most problems under consideration, the following inequality is used. If $f(x)$ and $\varphi(x)$ are two arbitrary densities, then it can be proven that:

$$-\int_{-\infty}^{\infty} \varphi(x) \ln \varphi(x) dx \leq -\int_{-\infty}^{\infty} \varphi(x) \ln f(x) dx$$

- In the coin experiment, the probability of heads is often viewed as an RV \mathbf{p} (bayesian estimation).
- We shall show that if no prior information about \mathbf{p} is available, then, according to the ME principle, its density $f(p)$ is uniform in the interval $(0,1)$.
- In this problem we must maximize $H(\mathbf{p})$ subject to the constraint (dictated by the meaning of \mathbf{p}) that $f(p) = 0$ outside the interval $(0, 1)$.
- The corresponding entropy is, therefore, given by $H(\mathbf{p}) = -\int_0^1 f(p) \ln f(p) dp$ and our problem is to find $f(p)$ such as to maximize the above integral.
- We maintain that $H(\mathbf{p})$ is maximum if $f(p) = 1$, hence $H(\mathbf{p}) = 0$.
- Indeed, if $\varphi(p)$ is any other density such that $\varphi(p) = 0$ outside $(0, 1)$, then

$$-\int_0^1 \varphi(p) \ln \varphi(p) dp \leq -\int_0^1 \varphi(p) \ln f(p) dp = 0 = H(\mathbf{p})$$

Bayesian Estimation

We investigate the problem of estimating the parameter θ of a distribution $F(x, \theta)$.

In the classical approach, we view θ as an unknown constant and the estimate was based solely on the observed values x_i of the RV \mathbf{x} .

In certain applications, θ is not totally unknown.

If, for example, θ is the probability of six in the die experiment, we expect that its possible values are close to $1/6$ because most dice are reasonably fair.

In bayesian statistics, the available prior information about θ is used in the estimation problem.

In this approach, the unknown parameter θ is viewed as the value of an RV Θ and the distribution of \mathbf{x} is interpreted as the conditional distribution $F_{\mathbf{x}}(x|\theta)$ of \mathbf{x} assuming $\Theta = \theta$.

The prior information is used to assign somehow a density $f_{\Theta}(\theta)$ to the RV Θ , and the problem is to estimate the value θ of Θ in terms of the observed values x_i of \mathbf{x} and the density of Θ .

The problem of estimating the unknown parameter θ is thus changed to the problem of estimating the value θ of the RV Θ .

Thus, in bayesian statistics, estimation is changed to prediction.

- We shall consider now a class of problems involving constraints in the form of expected values. Such problems are common in hydrology.
- We wish to determine the density $f(x)$ of an RV \mathbf{x} subject to the condition that the expected values η_i of n known functions $g_i(x)$ of \mathbf{x} are given

$$E\{g_i(\mathbf{x})\} = \int_{-\infty}^{\infty} g_i(x) f(x) dx = \eta_i \quad i = 1, \dots, n$$

- We shall show that the ME method leads to the conclusion that $f(x)$ must be an exponential

$$f(x) = A \exp\{-\lambda_1 g_1(x) - \dots - \lambda_n g_n(x)\}$$

- Where λ_i are n constants determined from the above equations $E\{g_i(\mathbf{x})\}$ and A is such as to satisfy the density condition

$$A \int_{-\infty}^{\infty} \exp\{\lambda_1 g_1(x) - \dots - \lambda_n g_n(x)\} dx = 1$$

- Suppose that $f(x) = A \exp\{-\lambda_1 g_1(x) - \dots - \lambda_n g_n(x)\}$

In this case:

$$\int_{-\infty}^{\infty} f(x) \ln f(x) dx = \int_{-\infty}^{\infty} f(x) [\ln A - \lambda_1 g_1(x) - \dots - \lambda_n g_n(x)] dx$$

- Hence: $H(\mathbf{x}) = \lambda_1 \eta_1 + \dots + \lambda_n \eta_n - \ln A$
- Now it suffices, therefore, to show that, if $\varphi(x)$ is any other density satisfying the constraints $E\{g_i(\mathbf{x})\}$, then its entropy cannot exceed the right side of the above equation

$$\begin{aligned} -\int_{-\infty}^{\infty} \varphi(x) \ln \varphi(x) dx &\leq -\int_{-\infty}^{\infty} \varphi(x) \ln f(x) dx \\ &= \int_{-\infty}^{\infty} \varphi(x) [\lambda_1 g_1(x) + \dots + \lambda_n g_n(x) - \ln A] dx \\ &= \lambda_1 \eta_1 + \dots + \lambda_n \eta_n - \ln A \end{aligned}$$

Example 1:

We shall determine $f(x)$ assuming that \mathbf{x} is a positive RV with known mean η .

With $g(x) = x$, it follows from the ME density that:

$$f(x) = Ae^{-\lambda x}, \text{ if } x > 0,$$

$$f(x) = 0, \text{ if } x \leq 0.$$

We have thus shown that if an RV is positive with specified mean, then its density obtained with the MEM, is an exponential.

Example 2:

We shall find such $f(x)$ that $E\{\mathbf{x}^2\} = m_2$. With $g_1(x) = x^2$, it follows from the ME density that:

$$f(x) = Ae^{-\lambda x^2}$$

Thus, if the second moment m_2 of an RV \mathbf{x} is specified, then \mathbf{x} is $N(0, m_2)$.

We can show similarly that if the variance σ^2 of \mathbf{x} is specified, then \mathbf{x} is $N(\eta, \sigma^2)$, where η is an arbitrary constant.

- The ME method can be used to determine the statistics of a stochastic process subject to given constraints.
- Suppose that \mathbf{x}_n is a wide-sense stationary (WSS) process with autocorrelation $R[m] = E\{\mathbf{x}_{n+m} \mathbf{x}_n\}$.
- We wish to find its various densities assuming that $R[m]$ is specified either for some or for all values of m .
- The ME principle leads to the conclusion that, in both cases, \mathbf{x}_n must be a normal process with zero mean. This completes the statistical description of \mathbf{x}_n if $R[m]$ is known for all m .
- If, however, we know $R[m]$ only partially, then we must find its unspecified values. For finite-order densities, this involves the maximization of the corresponding entropy with respect to the unknown values of $R[m]$ and it is equivalent to the maximization of the correlation determinant Δ .

5. CONCLUSIONS

Conclusions

- Entropy is a valuable tool to provide a quantitative measure of uncertainty of stochastic modelling of natural processes.
- An important application of entropy is the determination of the statistics of a stochastic process subject to various constraints, with the maximum entropy (ME) method.
- We should emphasize, however, that as in the case of the classical definition of probability, the conclusions drawn from the ME method must be accepted with skepticism particularly when they involve elaborate constraints.
- Extremal entropy considerations may provide an important connection with statistical mechanics. Thus, the ME principle may provide a physical background in the stochastic representation of natural processes.

Thank you for your attention

“Von Neumann told me, ‘You should call it entropy, for two reasons. In the first place your uncertainty function has been used in statistical mechanics under that name. In the second place, and more important, no one knows what entropy really is, so in a debate you will always have the advantage.’”

Claude Elwood Shannon