

# Chapter 8

## Inference

In Chapter 7 the process model was introduced as a way of accounting for flow of information through processes that are discrete, finite, and memoryless, and which may be nondeterministic and nontransparent. Although the model was motivated by the way many communication systems work, it is more general.

Formulas were given for input information  $I$ , loss  $L$ , mutual information  $M$ , noise  $N$ , and output information  $J$ . Each of these is measured in bits, although in a setting in which many symbols are chosen, one after another, they may be multiplied by the rate of symbol selection and then expressed in bits per second. The information flow is shown in Figure 8.1. All these quantities depend on the input probability distribution  $p(A_i)$ .

If the input probabilities are already known, and a particular output outcome is observed, it is possible to make inferences about the input event that led to that outcome. Sometimes the input event can be identified with certainty, but more often the inferences are in the form of changes in the initial input probabilities. This is typically how communication systems work—the output is observed and the “most likely” input event is inferred. Inference in this context is sometime referred to as **estimation**. It is the topic of Section 8.1.

On the other hand, if the input probabilities are not known, this approach does not work. We need a way to get the initial probability distribution. An approach that is based on the information analysis is discussed in Section 8.2 and in subsequent chapters of these notes. This is the Principle of Maximum Entropy.

### 8.1 Estimation

It is often necessary to determine the input event when only the output event has been observed. This is the case for communication systems, in which the objective is to infer the symbol emitted by the source so that it can be reproduced at the output. It is also the case for memory systems, in which the objective is to recreate the original bit pattern without error.

In principle, this estimation is straightforward if the input probability distribution  $p(A_i)$  and the conditional output probabilities, conditioned on the input events,  $p(B_j | A_i) = c_{ji}$ , are known. These “forward” conditional probabilities  $c_{ji}$  form a matrix with as many rows as there are output events, and as many columns as there are input events. They are a property of the process, and do not depend on the input probabilities  $p(A_i)$ .

The unconditional probability  $p(B_j)$  of each output event  $B_j$  is

$$p(B_j) = \sum_i c_{ji} p(A_i) \tag{8.1}$$

---

Author: [Paul Penfield, Jr.](#)

This document: <http://www.mtl.mit.edu/Courses/6.050/2007/notes/chapter8.pdf>

Version 1.4, April 2, 2007. Copyright © 2007 Massachusetts Institute of Technology

[Start of notes](#) · [back](#) · [next](#) | [6.050J/2.110J home page](#) | [Site map](#) | [Search](#) | [About this document](#) | [Comments and inquiries](#)

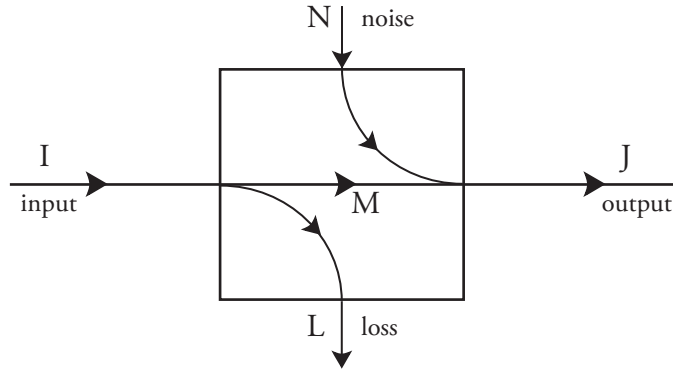


Figure 8.1: Information flow in a discrete memoryless process

and the joint probability of each input with each output  $p(A_i, B_j)$  and the backward conditional probabilities  $p(A_i | B_j)$  can be found using Bayes' Theorem:

$$\begin{aligned} p(A_i, B_j) &= p(B_j)p(A_i | B_j) \\ &= p(A_i)p(B_j | A_i) \\ &= p(A_i)c_{ji} \end{aligned} \quad (8.2)$$

Now let us suppose that a particular output event  $B_j$  has been observed. The input event that “caused” this output can be estimated only to the extent of giving a probability distribution over the input events. For each input event  $A_i$  the probability that it was the input is simply the backward conditional probability  $p(A_i | B_j)$  for the particular output event  $B_j$ , which can be written using Equation 8.2 as

$$p(A_i | B_j) = \frac{p(A_i)c_{ji}}{p(B_j)} \quad (8.3)$$

If the process has no loss ( $L = 0$ ) then for each  $j$  exactly one of the input events  $A_i$  has nonzero probability, and therefore its probability  $p(A_i | B_j)$  is 1. In the more general case, with nonzero loss, estimation consists of refining a set of input probabilities so they are consistent with the known output. Note that this approach only works if the original input probability distribution is known. All it does is refine that distribution in the light of new knowledge, namely the observed output.

It might be thought that the new input probability distribution would have less uncertainty than that of the original distribution. Is this always true?

The uncertainty of a probability distribution is, of course, its entropy as defined earlier. The uncertainty (about the input event) before the output event is known is

$$U_{\text{before}} = \sum_i p(A_i) \log_2 \left( \frac{1}{p(A_i)} \right) \quad (8.4)$$

The residual uncertainty, after some particular output event is known, is

$$U_{\text{after}}(B_j) = \sum_i p(A_i | B_j) \log_2 \left( \frac{1}{p(A_i | B_j)} \right) \quad (8.5)$$

The question, then is whether  $U_{\text{after}}(B_j) \leq U_{\text{before}}$ . The answer is often, but not always, yes. However, it is not difficult to prove that the average (over all output states) of the residual uncertainty is less than the original uncertainty:

$$\sum_j p(B_j) U_{\text{after}}(B_j) \leq U_{\text{before}} \quad (8.6)$$

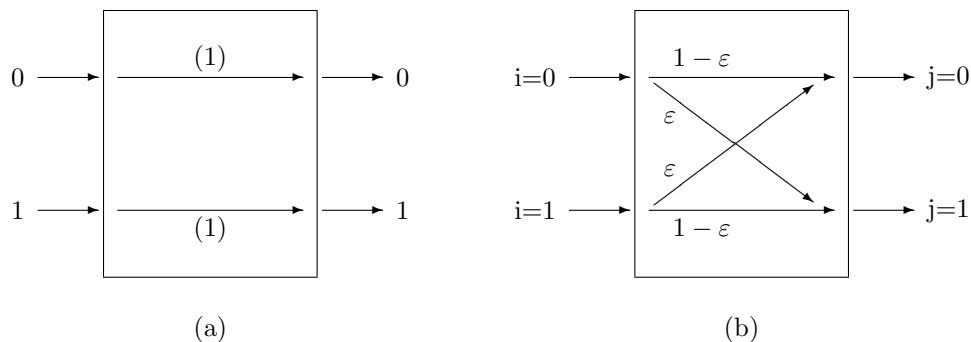


Figure 8.2: (a) Binary Channel without noise (b) Symmetric Binary Channel, with errors

In words, this statement says that on average, our uncertainty about the input state is never increased by learning something about the output state. In other words, on average, this technique of inference helps us get a better estimate of the input state.

Two of the following examples will be continued in subsequent chapters including the next chapter on the Principle of Maximum Entropy—the symmetric binary channel and Berger’s Burgers.

### 8.1.1 Symmetric Binary Channel

The noiseless, lossless binary channel shown in Figure 8.2(a) is a process with two input values which may be called 0 and 1, two output values similarly named, and a transition matrix  $c_{ji}$  which guarantees that the output equals the input:

$$\begin{bmatrix} c_{00} & c_{01} \\ c_{10} & c_{11} \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad (8.7)$$

This channel has no loss and no noise, and the mutual information, input information, and output information are all the same.

The symmetric binary channel (Figure 8.2(b)) is similar, but occasionally makes errors. Thus if the input is 1 the output is not always 1, but with the “bit error probability”  $\varepsilon$  is flipped to the “wrong” value 0, and hence is “correct” only with probability  $1 - \varepsilon$ . Similarly, for the input of 0, the probability of error is  $\varepsilon$ . Then the transition matrix is

$$\begin{bmatrix} c_{00} & c_{01} \\ c_{10} & c_{11} \end{bmatrix} = \begin{bmatrix} 1 - \varepsilon & \varepsilon \\ \varepsilon & 1 - \varepsilon \end{bmatrix} \quad (8.8)$$

This channel is symmetric in the sense that the errors in both directions (from 0 to 1 and vice versa) are equally likely.

Because of the loss, the input event associated with, say, output event  $B_0$  cannot be determined with certainty. Nevertheless, the formulas above can be used. In the important case where the two input probabilities are equal (and therefore each equal to 0.5) an output of 0 implies that the input event  $A_0$  has probability  $1 - \varepsilon$  and input event  $A_1$  has probability  $\varepsilon$ . Thus if, as would be expected in a channel designed for low-error communication,  $\varepsilon$  is small, then it would be reasonable to infer that the input that produced the output event  $B_0$  was the event  $A_0$ .

### 8.1.2 Non-symmetric Binary Channel

A non-symmetric binary channel is one in which the error probabilities for inputs 0 and 1 are different, i.e.,  $c_{01} \neq c_{10}$ . We illustrate non-symmetric binary channels with an extreme case based on a medical test for Huntington’s Disease.

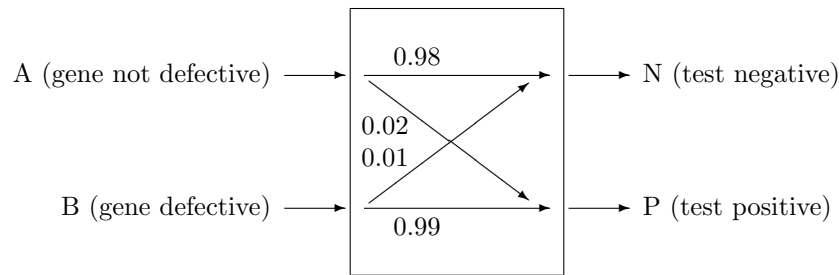


Figure 8.3: Huntington's Disease Test

Huntington's Disease is a rare, progressive, hereditary brain disorder with no known cure. It was named after Dr. George Huntington (1850–1915), a Long Island physician who published a description in 1872. It is caused by a defective gene, which was identified in 1983. Perhaps the most famous person afflicted with the disease was the songwriter Woodie Guthrie.

The biological child of a person who carries the defective gene has a 50% chance of inheriting the defective gene. For the population as a whole, the probability of carrying the defective gene is much lower. According to the Huntington's Disease Society of America, <http://www.hdsa.org>, "More than a quarter of a million Americans have HD or are 'at risk' of inheriting the disease from an affected parent." This is about 1/1000 of the population, so for a person selected at random, with an unknown family history, the probability of carrying the affected gene is 1/2000.

People carrying the defective gene all eventually develop the disease unless they die of another cause first. The symptoms most often appear in middle age, in people in their 40's or 50's, perhaps after they have already produced a family and thereby possibly transmitted the defective gene to another generation. Although the disease is not fatal, those in advanced stages generally die from its complications. Until recently, people with a family history of the disease were faced with a life of uncertainty, not knowing whether they carried the defective gene, and not knowing how to manage their personal and professional lives.

In 1993 a test was developed which can tell if you carry the defective gene. Unfortunately the test is not perfect; there is a probability of a false positive (reporting you have it when you actually do not) and of a false negative (reporting your gene is not defective when it actually is). For our purposes we will assume that the test only gives a yes/no answer, and that the probability of a false positive is 2% and the probability of a false negative is 1%. (The real test is actually better—it also estimates the severity of the defect, which is correlated with the age at which the symptoms start.)

If you take the test and learn the outcome, you would of course like to infer whether you will develop the disease eventually. The techniques developed above can help.

Let us model the test as a discrete memoryless process, with input  $A$  (no defective gene) and  $B$  (defective gene), and outputs  $P$  (positive) and  $N$  (negative). The process is shown in Figure 8.3.

First, consider the application of this test to someone with a family history, for which  $p(A) = p(B) = 0.5$ . Then, if the test is negative, the probability, for that person, of having the defect is  $1/99 = 0.0101$  and the probability of not having it is  $98/99 = 0.9899$ . On the other hand, if the test is positive, the probability, for that person, of carrying the defective gene is  $99/101 = 0.9802$  and the probability of not doing so is  $2/101 = 0.0198$ . The test is very effective, in that the two outputs imply, to high probability, different inputs.

An interesting question that is raised by the existence of this test but is not addressed by our mathematical model is whether a person with a family history would elect to take the test, or whether he or she would prefer to live not knowing what the future holds in store. The development of the test was funded by a group including Guthrie's widow and headed by a man named Milton Wexler (1908–2007), who was concerned about his daughters because his wife and her brothers all had the disease. Wexler's daughters, whose situation inspired the development of the test, decided not to take it.

Next, consider the application of this test to someone with an unknown family history, so that  $p(A) =$

0.9995 and  $p(B) = 0.0005$ . Then, if the test is negative, the probability of that person carrying the defective gene  $p(B | N)$  is

$$\frac{0.0005 \times 0.01}{0.0005 \times 0.01 + 0.9995 \times 0.98} = 0.000005105 \quad (8.9)$$

and the probability of that person carrying the normal gene  $p(A | N)$  is

$$\frac{0.9995 \times 0.98}{0.0005 \times 0.01 + 0.9995 \times 0.98} = 0.999994895 \quad (8.10)$$

On the other hand, if the test is positive, the probability of that person carrying the defective gene  $p(B | P)$  is

$$\frac{0.0005 \times 0.99}{0.0005 \times 0.99 + 0.9995 \times 0.02} = 0.02416 \quad (8.11)$$

and the probability of not having the defect  $p(A | P)$  is

$$\frac{0.9995 \times 0.02}{0.0005 \times 0.99 + 0.9995 \times 0.02} = 0.97584 \quad (8.12)$$

The test does not seem to distinguish the two possible inputs, since the overwhelming probability is that the person has a normal gene, regardless of the test results. In other words, if you get a positive test result, it is more likely to have been caused by a testing error than a defective gene. There seems to be no useful purpose served by testing people without a family history. (Of course repeated tests could be done to reduce the false positive rate.)

An information analysis shows clearly the difference between these two cases. First, recall that probabilities are subjective, or observer-dependent. The lab technician performing the test presumably does not know whether there is a family history, and so would not be able to infer anything from the results. Only someone who knows the family history could make a useful inference. Second, it is instructive to calculate the information flow in the two cases. Recall that all five information measures ( $I$ ,  $L$ ,  $M$ ,  $N$ , and  $J$ ) depend on the input probabilities. A straightforward calculation of the two cases leads to the information quantities (in bits) in Table 8.1 (note how much larger  $N$  is than  $M$  if there is no known family history).

	$p(A)$	$p(B)$	$I$	$L$	$M$	$N$	$J$
Family history	0.5	0.5	1.00000	0.11119	0.88881	0.11112	0.99993
Unknown Family history	0.9995	0.0005	0.00620	0.00346	0.00274	0.14141	0.14416

Table 8.1: Huntington's Disease Test Process Model Characteristics

Clearly the test conveys information about the patients' status by reducing the uncertainty in the case where there is a family history of the disease. On the other hand, without a family history there is very little information that could possibly be conveyed because there is so little initial uncertainty.

### 8.1.3 Berger's Burgers

A former 6.050J/2.110J student opened a fast-food restaurant, and named it in honor of the very excellent Undergraduate Assistant of the course. At Berger's Burgers, meals are prepared with state-of-the-art high-tech equipment using reversible computation for control. To reduce the creation of entropy there are no warming tables, but instead the entropy created by discarding information is used to keep the food warm. Because the rate at which information is discarded in a computation is unpredictable, the food does not always stay warm. There is a certain probability, different for the different menu items, of a meal being "COD" (cold on delivery).

The three original menu items are Value Meals 1, 2, and 3. Value Meal 1 (burger) costs \$1, contains 1000 Calories, and has a probability 0.5 of arriving cold. Value Meal 2 (chicken) costs \$2, has 600 Calories, and

a probability 0.2 of arriving cold. Value Meal 3 (fish) costs \$3, has 400 Calories, and has a 0.1 probability of being cold.

Item	Entree	Cost	Calories	Probability of arriving hot	Probability of arriving cold
Value Meal 1	Burger	\$1.00	1000	0.5	0.5
Value Meal 2	Chicken	\$2.00	600	0.8	0.2
Value Meal 3	Fish	\$3.00	400	0.9	0.1

Table 8.2: Berger's Burgers

There are several inference questions that can be asked about Berger's Burgers. All require an initial assumption about the buying habits of the public, i.e., about the probability of each of the three meals being ordered  $p(B)$ ,  $p(C)$ , and  $p(F)$ . Then, upon learning another fact, such as a particular customer's meal arriving cold, these probabilities can be refined to lead to a better estimate of the meal that was ordered.

Suppose you arrive at Berger's Burgers with your friends and place your orders. Assume that money is in plentiful supply so you and your friends are equally likely to order any of the three meals. Also assume that you do not happen to hear what your friends order or see how much they pay. Also assume that you do not know your friends' taste preferences and that the meals come in identical packages so you cannot tell what anyone else received by looking.

Before the meals are delivered, you have no knowledge of what your friends ordered and might assume equal probability of  $1/3$  for  $p(B)$ ,  $p(C)$ , and  $p(F)$ . You can estimate the average amount paid per meal (\$2.00), the average Calorie count (667 Calories), and the probability that any given order would be COD (0.267).

Now suppose your friend Alice remarks that her meal is cold. Knowing this, what is the probability she ordered a burger? (0.625) Chicken? (0.25) Fish? (0.125). And what is the expected value of the amount she paid for her meal? (\$1.50) And what is her expected Calorie count? (825 Calories)

Next suppose your friend Bob says he feels sorry for her and offers her some of his meal, which is hot. Straightforward application of the formulas above can determine the refined probabilities of what he ordered, along with the expected calorie count and cost.

### 8.1.4 Inference Strategy

Often, it is not sufficient to calculate the probabilities of the various possible input events. The correct operation of a system may require that a definite choice be made of exactly one input event. For processes without loss, this can be done accurately. However, for processes with loss, some strategy must be used to convert probabilities to a single choice.

One simple strategy, "maximum likelihood," is to decide on whichever input event has the highest probability after the output event is known. For many applications, particularly communication with small error, this is a good strategy. It works for the symmetric binary channel when the two input probabilities are equal. However, sometimes it does not work at all. For example, if used for the Huntington's Disease test on people without a family history, this strategy would never say that the person has a defective gene, regardless of the test results.

Inference is important in many fields of interest, such as machine learning, natural language processing and other areas of artificial intelligence. An open question, of current research interest, is which inference strategies are best suited for particular purposes.

## 8.2 Principle of Maximum Entropy: Simple Form

In the last section, we discussed one technique of estimating the input probabilities of a process given that the output event is known. This technique, which relies on the use of Bayes' Theorem, only works if the

process is lossless (in which case the input can be identified with certainty) or an a priori input probability distribution is assumed (in which case the technique refines the initial probability distribution to take account of the known output).

The Principle of Maximum Entropy is a technique that can be used to estimate input probabilities more generally. The result is a probability distribution that is consistent with known constraints expressed in terms of averages, or expected values, of one or more quantities, but is otherwise as unbiased as possible (the word “bias” is used here not in the technical sense of statistics, but the everyday sense of a preference that inhibits impartial judgment). This principle is described first for the simple case of one constraint and three input events, in which case the technique can be carried out analytically. Then it is described more generally in Chapter 9.

This principle has applications in many domains, but was originally motivated by statistical physics, which attempts to relate macroscopic, measurable properties of physical systems to a description at the atomic or molecular level. It can be used to approach physical systems from the point of view of information theory, because the probability distributions can be derived by avoiding the assumption that the observer has more information than is actually available. Information theory, particularly the definition of information in terms of probability distributions, provides a quantitative measure of ignorance (or uncertainty, or entropy) that can be maximized mathematically to find the probability distribution that best avoids unnecessary assumptions.

This approach to statistical physics was pioneered by Edwin T. Jaynes (1922–1998), a professor at Washington University in St. Louis, and previously Stanford University. The seminal publication was

- E. T. Jaynes, “Information Theory and Statistical Mechanics,” *Physical Review*, vol. 106, no. 4, pp. 620-630; May 15, 1957.  
(<http://bayes.wustl.edu/etj/articles/theory.1.pdf>)

Other references of interest by Jaynes include:

- a continuation of this paper, E. T. Jaynes, “Information Theory and Statistical Mechanics. II,” *Physical Review*, vol. 108, no. 2, pp. 171-190; October 15, 1957.  
(<http://bayes.wustl.edu/etj/articles/theory.1.pdf>)
- a review paper, including an example of estimating probabilities of an unfair die, E. T. Jaynes, “Information Theory and Statistical Mechanics,” pp. 181-218 in “Statistical Physics,” Brandeis Summer Institute 1962, W. A. Benjamin, Inc., New York, NY; 1963.  
(<http://bayes.wustl.edu/etj/articles/brandeis.pdf>)
- personal history of the approach, Edwin T. Jaynes, “Where Do We Stand on Maximum Entropy?,” pp. 15-118, in “The Maximum Entropy Formalism,” Raphael D. Levine and Myron Tribus, editors, The MIT Press, Cambridge, MA; 1979.  
(<http://bayes.wustl.edu/etj/articles/stand.on.entropy.pdf>)

The philosophy of assuming maximum uncertainty as an approach to thermodynamics is discussed in

- Chapter 3 of M. Tribus, “Thermostatistics and Thermodynamics,” D. Van Nostrand Co, Inc., Princeton, NJ; 1961.

Before the Principle of Maximum Entropy can be used the problem domain needs to be set up. In cases involving physical systems, this means that the various states in which the system can exist need to be identified, and all the parameters involved in the constraints known. For example, the energy, electric charge, and other quantities associated with each of the states is assumed known. Often quantum mechanics is needed for this task. It is not assumed in this step which particular state the system is in (or, as often expressed, which state is actually “occupied”); indeed it is assumed that we do not know and cannot know this with certainty, and so we deal instead with the probability of each of the states being occupied. Thus we use probability as a means of coping with our lack of complete knowledge. Naturally we want to avoid

inadvertently assuming more knowledge than we actually have, and the Principle of Maximum Entropy is the technique for doing this. In the application to nonphysical systems, again the various possible events have to be enumerated and the properties of each type associated with each of the possibilities known. In these notes we will derive a simple form of the Principle of Maximum Entropy and apply it to the restaurant example set up in Section 8.1.3.

### 8.2.1 Berger's Burgers

The Principle of Maximum Entropy will be introduced by means of an example. This example was described in Section 8.1.3. A fast-food restaurant, Berger's Burgers, offers three meals: burger, chicken, and fish. The price, Calorie count, and probability of each meal being delivered cold are as listed in Table 8.2.

### 8.2.2 Probabilities

This example has been defined so that the choice of one of the three meals constitutes an outcome. If we do not know this outcome we may still have some knowledge, and we use probabilities to express this knowledge. The question is how to assign probabilities that are consistent with whatever information we may have.

In the case of Berger's Burgers, there are three probabilities which for simplicity we denote  $p(B)$ ,  $p(C)$ , and  $p(F)$  for the three meals. A probability distribution  $p(A_i)$  has the property that each of the probabilities is between or equal to 0 and 1, and, since the input events are mutually exclusive and exhaustive, the sum of all the probabilities is 1:

$$\begin{aligned} 1 &= \sum_i p(A_i) \\ &= p(B) + p(C) + p(F) \end{aligned} \quad (8.13)$$

If any of the probabilities is equal to 1 then all the other probabilities are 0, and we know exactly which state the system is in; in other words, we have no uncertainty and there is no need to resort to probabilities.

### 8.2.3 Entropy

More generally our uncertainty is expressed quantitatively by the information which we do not have about the meal chosen, or the state occupied. This is the entropy

$$\begin{aligned} S &= \sum_i p(A_i) \log_2 \left( \frac{1}{p(A_i)} \right) \\ &= p(B) \log_2 \left( \frac{1}{p(B)} \right) + p(C) \log_2 \left( \frac{1}{p(C)} \right) + p(F) \log_2 \left( \frac{1}{p(F)} \right) \end{aligned} \quad (8.14)$$

Here, information is in bits because we are using logarithms to base 2.

In the context of physical systems this uncertainty is known as the entropy. In communication systems the uncertainty regarding which actual message is to be transmitted is also known as the entropy of the source. Note that in general the entropy, because it is expressed in terms of probabilities, depends on the observer. One person may have different knowledge of the system from another, and therefore would calculate a different numerical value for entropy. The Principle of Maximum Entropy is used to discover the probability distribution which leads to the highest value for this uncertainty, thereby assuring that no information is inadvertently assumed.



### 8.2.4 Constraints

It is a property of the entropy formula above that it has its maximum value when all probabilities are equal (we assume the number of possible states is finite). This property is easily proved using the Gibbs inequality. If we have no additional information about the system, then such a result seems reasonable. However, if we have additional information then we ought to be able to find a probability distribution that is better in the sense that it has less uncertainty.

For simplicity we consider only one such constraint, namely that we know the expected value of some quantity (the Principle of Maximum Entropy can handle multiple constraints but the mathematical procedures and formulas become more complicated). The quantity in question is one for which each of the states of the system has its own amount, and the expected value is found by averaging the values corresponding to each of the states, taking into account the probabilities of those states. Thus if there is a quantity  $G$  for which each of the states has a value  $g(A_i)$  then we want to consider only those probability distributions for which the expected value is  $G$

$$G = \sum_i p(A_i)g(A_i) \quad (8.15)$$

Note that this constraint cannot be achieved if  $G$  is less than the smallest  $g(A_i)$  or larger than the largest  $g(A_i)$ .

For our Berger's Burgers example, suppose we are told that the average price of a meal is \$1.75, and we want to estimate the separate probabilities of the various meals without making any other assumptions. Then our constraint would be

$$\$1.75 = \$1.00p(B) + \$2.00p(C) + \$3.00p(F) \quad (8.16)$$

Note that the probabilities are dimensionless and so both the expected value of the constraint and the individual values must be expressed in the same units, in this case dollars.

### 8.2.5 Maximum Entropy, Analytic Form

Here we demonstrate the Principle of Maximum Entropy for the simple case in which there is one constraint and three variables. It will be possible to go through all the steps analytically.

Suppose you have been hired by Carnivore Corporation, the parent company of Berger's Burgers, to analyze their worldwide sales. You visit Berger's Burgers restaurants all over the world, and determine that, on average, people are paying \$1.75 for their meals. (As part of Carnivore's commitment to global homogeneity, the price of each meal is exactly the same in every restaurant, after local currencies are converted to U.S. dollars.)

After you return, your supervisors ask about the probabilities of a customer ordering each of the three value meals. In other words, they want to know  $p(B)$ ,  $p(C)$ , and  $p(F)$ . You are horrified to realize that you did not keep the original data, and there is no time to repeat your trip. You have to make the best estimate of the probabilities  $p(B)$ ,  $p(C)$ , and  $p(F)$  consistent with the two things you do know:

$$1 = p(B) + p(C) + p(F) \quad (8.17)$$

$$\$1.75 = \$1.00p(B) + \$2.00p(C) + \$3.00p(F) \quad (8.18)$$

Since you have three unknowns and only two equations, there is not enough information to solve for the unknowns.

What should you do? There are a range of values of the probabilities that are consistent with what you know. However, these leave you with different amounts of uncertainty  $S$

$$S = p(B) \log_2 \left( \frac{1}{p(B)} \right) + p(C) \log_2 \left( \frac{1}{p(C)} \right) + p(F) \log_2 \left( \frac{1}{p(F)} \right) \quad (8.19)$$

If you choose one for which  $S$  is small, you are assuming something you do not know. For example, if your average had been \$2.00 rather than \$1.75, you could have met both of your constraints by assuming that everybody bought the chicken meal. Then your uncertainty would have been 0 bits. Or you could have assumed that half the orders were for burgers and half for fish, and the uncertainty would have been 1 bit. Neither of these assumptions seems particularly appropriate, because each goes beyond what you know. How can you find that probability distribution that uses no further assumptions beyond what you already know?

The **Principle of Maximum Entropy** is based on the reasonable assumption that you should select that probability distribution which leaves you the largest remaining uncertainty (i.e., the maximum entropy) consistent with your constraints. That way you have not introduced any additional assumptions into your calculations.

For the simple case of three probabilities and two constraints, this is easy to do analytically. Working with the two constraints, two of the unknown probabilities can be expressed in terms of the third. For our case we can multiply Equation 8.17 above by \$1.00 and subtract it from Equation 8.18, to eliminate  $p(B)$ . Then we can multiply the first by \$2.00 and subtract it from the second, thereby eliminating  $p(C)$ :

$$p(C) = 0.75 - 2p(F) \quad (8.20)$$

$$p(B) = 0.25 + p(F) \quad (8.21)$$

Next, the possible range of values of the probabilities can be determined. Since each of the three lies between 0 and 1, it is easy to conclude from these results that

$$0 \leq p(F) \leq 0.375 \quad (8.22)$$

$$0 \leq p(C) \leq 0.75 \quad (8.23)$$

$$0.25 \leq p(B) \leq 0.625 \quad (8.24)$$

Next, these expressions can be substituted into the formula for entropy so that it is expressed in terms of a single probability. Thus

$$S = (0.25 + p(F)) \log_2 \left( \frac{1}{0.25 + p(F)} \right) + (0.75 - 2p(F)) \log_2 \left( \frac{1}{0.75 - 2p(F)} \right) + p(F) \log_2 \left( \frac{1}{p(F)} \right) \quad (8.25)$$

Any of several techniques can now be used to find the value of  $p(F)$  for which  $S$  is the largest. In this case the maximum occurs for  $p(F) = 0.216$  and hence  $p(B) = 0.466$ ,  $p(C) = 0.318$ , and  $S = 1.517$  bits.

After estimating the input probability distribution, any averages over that distribution can be estimated. For example, in this case the average Calorie count can be calculated (it is 743.2 Calories), or the probability of a meal being served cold (31.8%).

### 8.2.6 Summary

Let's remind ourselves what we have done. We have expressed our constraints in terms of the unknown probability distributions. One of these constraints is that the sum of the probabilities is 1. The other involves the average value of some quantity, in this case cost. We used these constraints to eliminate two of the variables. We then expressed the entropy in terms of the remaining variable. Finally, we found the value of the remaining variable for which the entropy is the largest. The result is a probability distribution that is consistent with the constraints but which has the largest possible uncertainty. Thus we have not inadvertently introduced any unwanted assumptions into the probability estimation.

This technique requires that the model for the system be known at the outset; the only thing not known is the probability distribution. As carried out in this section, with a small number of unknowns and one more unknown than constraint, the derivation can be done analytically. For more complex situations a more general approach is necessary. That is the topic of Chapter 9.