

Teddy Seidenfeld ¹

ENTROPY AND UNCERTAINTY

ABSTRACT

This essay is, primarily, a discussion of four results about the principle of maximizing entropy (MAXENT) and its connections with Bayesian theory. Result 1 provides a restricted equivalence between the two where the Bayesian model for MAXENT inference uses an *a priori* probability that is

(Frieden, 1972) and estimating missing proportions in contingency tables for socio-economic survey data (Denzau *et al.*, 1984). But I doubt there is a more staunch defender of the generality of entropy as a basis for quantifying

(probabilistic) uncertainty than the physicist E. T. Jaynes.

Almost thirty years ago Jaynes (1957) offered his celebrated papers

on "Information Theory and Statistical Mechanics." There he argued that

$1/6$ ($i = 1, \dots, 6$).² If, instead the constraint specifies

$$E[\text{number of spots on next roll}] = 4.5 \quad (3)$$

instead of the value 3.5 (for a fair die), the MAXENT solution (Jaynes, 1978) is (to five places):

$$\{p_1, \dots, p_6\} = \{.05435, .07877, .11416, .16545, .23977, .34749\}. \quad (4)$$

Note that in (4) the probabilities are shifted away from the uniform distribution to lie on a smooth (convex) curve, increasing (decreasing) in n .

Not all who have examined these supporting arguments find them convincing. (See especially Dias and Shimony, 1981; Frieden, 1984; Friedman and Shimony, 1971; Rowlinson, 1970; Shimony, 1973. Jaynes offers selected rebuttal in (1978).) In what follows I present concerns I have primarily

(S_1) U_S is a continuous function of the p_i 's.

(S_2) When $P = \{1/n, \dots, 1/n\}$ is the uniform distribution on n -states, U_S is monotonically increasing in n , the number of states over which one is uncertain.

(S_3) U_S is additive over decomposition of the sample space of possible outcomes. That is, let $\Omega = \{s_1, \dots, s_n\}$ be the set of (n) possible outcomes, and let Ω be partitioned into $m \leq n$ disjoint subsets $\Omega' = \{r_1, \dots, r_m\}$ with r_i a subset of Ω . If P is a probability function

Now, it is clear that $P_1(\cdot) = P_1(\cdot | e)$, since P_1 satisfies c_{k+1} . Define a probability $P'_1(\cdot) = P_1(\cdot | e)$.

dicting the assumption that P_0 is the MAXENT solution for constraints C_0 . To verify that P'_0 satisfies C_0 , note that the class of distributions satisfying

outcomes $\{1, 3, 5\}$ —which is the conditional probability $P'_0(\cdot | e_1)$ —but instead is the distribution (see Appendix).

$$P'_1(i) = \{.21624, .31752, .46624\} \quad (i = 1, 3, 5) \quad (5a)$$

where D_0 is discrete and by the same reasoning as above.

$$\int_{D_0} (D_1 - D_0) \int_{D_0} (D_1 - D_0) \dots$$

variables, some (Bayesian) conditionalizations do not agree with the revision from P^0 to P^1 by minimizing

Besides generalizing U_S with discrete distributions, I_K affords a consistent extension of entropy to continuous distributions, unlike the (natural)

of a cubical die (with spots from 1 to 6 arranged in conventional order), then

Table 1. ('Yes'/'No' identifies which arrangements are possible.)

spots showing

Source

2	Yes	No	Yes
3	Yes	Yes	Yes
4	Yes	No	Yes
5	Yes	Yes	Yes
6	Yes	Yes	Yes

priority in the application of Insufficient Reason. Of course, what is lookin

where p_i ($i = 1, \dots, 6$) is the probability of i -spots showing up.

However, since the alternative partition (Table 1) is a refinement of the six-fold partition used above, the constraint (10) applies there too. Specifically, define $f(\text{state } j)$ ($j = 1, \dots, 14$ —counting across possible states in Table 1) as follows:

$$f(\text{state}_1) = 1, f(\text{state}_2) = f(\text{state}_3) = 2,$$

$$f(\text{state}_4) = f(\text{state}_5) = f(\text{state}_6) = 3,$$

$$f(\text{state}_7) = f(\text{state}_8) = 4,$$

$$f(\text{state}_9) = f(\text{state}_{10}) = f(\text{state}_{11}) = 5,$$

and

$$f(\text{state}_{12}) = f(\text{state}_{13}) = f(\text{state}_{14}) = 6$$

Then (10) is equivalent to the constraint:

$$E[f] = 55/14. \tag{12}$$

But the distribution over the 14 states which maximizes entropy subject to

a constant and constants are (vacuously) probabilistically independent of other variables. In Section 4, where MAXENT is contrasted with Bayesian inference, the device of using a degenerate 0-1 distribution with nuisance factors is key to understanding an important objective.

Summary: The question addressed in this section is

conditional probability $P_{BK}(\cdot | \cdot)$ —coherence.

(B_2) $P_{BK}(\cdot | \cdot)$ is relativized to background evidence BK (consistent and closed under entailment), where BK depicts the agent's *total background evidence*.

(B_3) As regulated by Bayes' theorem for conditional probability $P_{BK}(\cdot | \cdot \& e)$ is the agent's hypothetical belief state for the hypothesis that he accepts only the new (consistent) evidence e , i.e., under the hypothesis that BK is enlarged by addition of e (and its consequences).

given BK)—conditionalization.¹⁰

We have Result 1 (p. 263 from Shannon's property (S_1)) establishing

a restricted equivalence between revising probabilities through MAXENT

-the MAXENT principle is coherent (from a Bayesian point of view) returns us to the question of the previous section. Under which conditions can we extend (refine) the field of events, while preserving MAXENT solutions for a given set of constraints?

It is from this perspective I propose we consider the interesting case

tion Jaynes wants for the "constraint" in his Brandeis Dice problem.¹⁶ So,
instead, let us examine Jaynes' constraint.

In (1978) he writes,

exchangeable P is a mixture of i.i.d. multinomial distributions (each on a sample space of n -outcomes) for some "mixing" prior π on the multinomial parameter. Recall, when $r = (n + 1)/2$, that is when the "sample average" equals the average of the number of spots showing on the n faces of the

5. COMMENTS ON THE CONCENTRATION THEOREM
(JAYNES, 1979 AND SEE (1963, PP. 51-52))

Theorem (Jaynes): Consider M repetitions of an experiment with n possi-

ble outcomes on a given trial. Let f_i ($1 \leq i \leq n$) be the observed relative frequency of the i th outcome in these M trials. Then the class of sequences

the MAXENT probability is determined by the (asymptotic) proportion of these states with frequencies close to the MAXENT distribution. Why is this a problem? It is because, if the concentration about the MAXENT solution

Support for this research came from the Department of Preventive
Medicine, Washington University (St. Louis) and NSF grant SES 8607300

paradoxes" (due to Dawid *et al.*, 1973). As Dawid *et al.* use their anomalies to question this view, Dawid *et al.*

(1971) formulation, it is supposed there is one state whose magnitude a_m equals the average of the n magnitudes $a_m = (1/n) \cdot \sum_{i=1}^n a_i$. This condi-

(Recall, $\lambda = \infty$ corresponds to the point probability distribution.)

I thank Prof. E. Greenberg for alerting me to Frieden's recent work.

APPENDIX A: ON THE MAXENT FORMALISM

Here we review some of the mathematics for solving the MAXENT

where $c_1 = -\lambda_1$. (The value of (A5) is $\lambda_1 = 55/14$.)

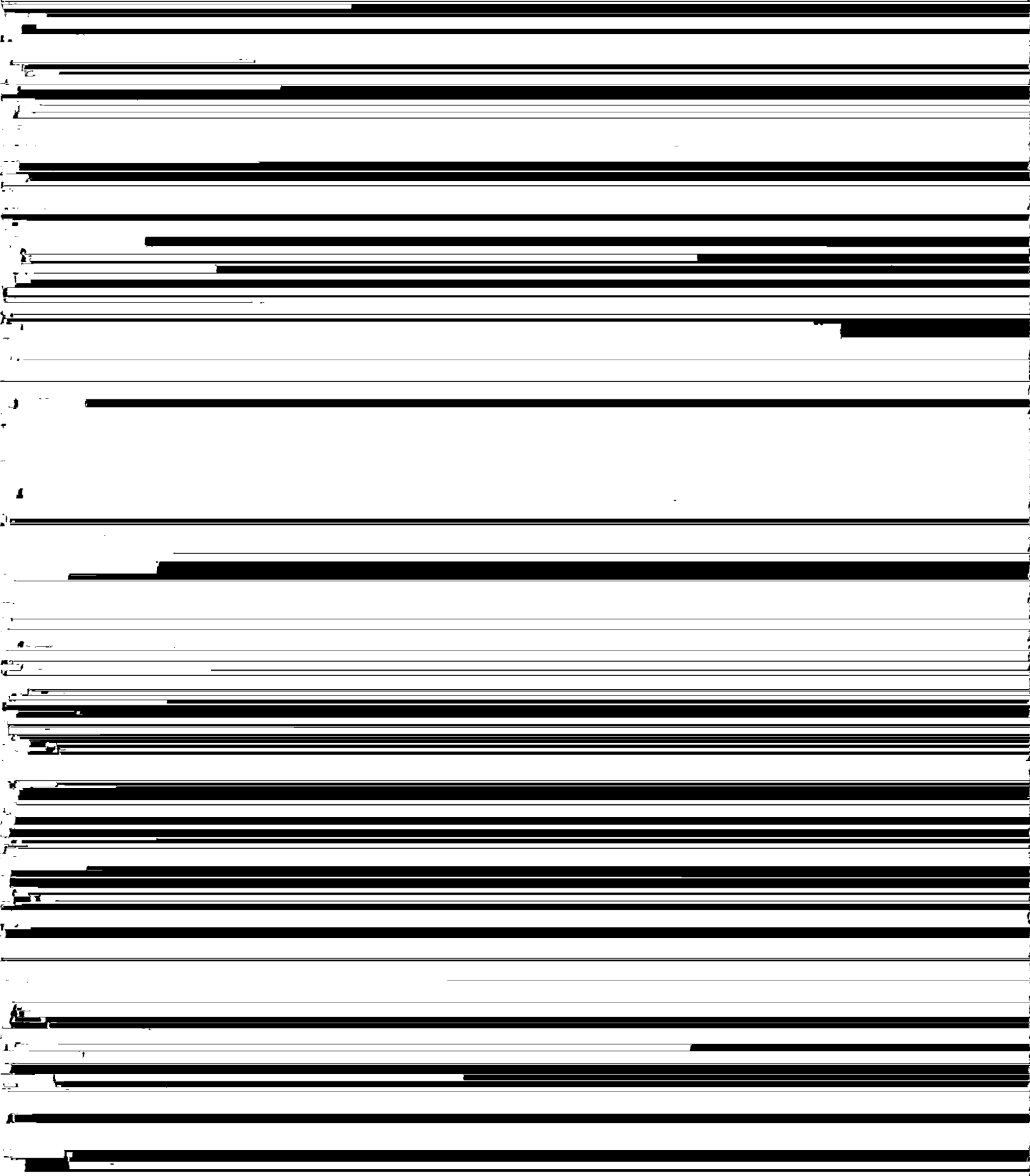
Then, by (A3) and (A4)

$$-\frac{\partial}{\partial \lambda_1} \log Z = (1 - 7x^6 + 6x^7) / [(1-x)(1-x^6)] = c_1. \quad (\text{A6})$$

In the problem discussed on p. 269, (10) sets the constraint: $c_1 = 55/14$. Solving (A6) for this value yields:

(as obtained on my TI 58C). This results in the MAXENT distribution (11) in accord with (A1). The MAXENT distributions (5a) and (5b) are

In other words, Result 4 establishes that the MAXENT ~~prob. 1.1.1. 2~~



and use the argument which follows to establish the desired property for each P_{iX}^1 . By continuity of cross-entropy shifts, the desired property obtains for their limit, P_X^1 .) Refine X to Y so that P_Y^0 is uniform on Y . (This is possible

c in the measure space (Y, \mathcal{Y}) . By the lemma (above), the minimum cross-entropy shift from P_Y^0 to P_Y^1 agrees with the minimum cross-entropy shift from P_X^0 to P_X^1 on X . But, with P_Y^0 uniform on Y , the minimum cross-entropy shift is just the MAXENT distribution P , in the measure space (Y, \mathcal{Y}) , subject to the constraint c . Then apply Result 4 to show that P_Y^1 has the desired property on E . To wit: $D^1(E) \geq D^0(E)$ unless $D^1 = D^0$.

The inequality (B6) is demonstrated as follows. Let

$$k = rm, \quad (B7)$$

so $P^U(E_1)/P^U(E_2) = r$. Substituting (B7) into (B5), we obtain

$$P(E_1 \cup E_2) \cdot Z = mr^{1-\alpha} (1/[\alpha^\alpha + (1-\alpha)^{1-\alpha}]). \quad (B8)$$

The inequality (B6) obtains just in case

$$1/[\alpha^\alpha + (1-\alpha)^{1-\alpha}] < (1-\alpha)/\alpha$$

Taking the derivative (with respect to r) of the r.h.s. of (B9) and setting it...
equal to 0...

$$r = (1-\alpha)/\alpha \{= P(E_1)/P(E_2)\} \quad (B10)$$

University Press.

4th printing. New York: Interscience Publishers.

Dawid, A. P., M. Stone, and J. V. Zidek (1973), "Marginalization paradoxes in Bayesian and structural inference." *Journal of the Royal Statistical Society, Series B* 35, 189-233 (with discussion).

Denzau, A. T., P. C. Gibbons, and E. Greenberg (1984). "Bayesian estimation of

- Jaynes, E. T. (1983), *Papers on Probability, Statistics and Statistical Physics*, ed. R. Rosenkrantz. Dordrecht: D. Reidel.
- Jaynes, E. T. (1983), "Highly informative priors." In *Bayesian Statistics 2*, ed. J.M. Bernardo, M. H. DeGroot, D. V. Lindley, and A. F. M. Smith, pp. 329-