This is page i Printer: Opaque this

Contents

4	Maz	ximum entropy distributions	1	
	4.1	Bayes' theorem and consistent reasoning	1	
	4.2	Maximum entropy distributions	3	
	4.3 Entropy			
	4.4	Discrete maximum entropy examples ¹ $\ldots \ldots \ldots \ldots$	5	
		4.4.1 Discrete uniform	5	
		4.4.2 Discrete nonuniform	5	
	4.5	Normalization and partition functions	6	
		4.5.1 Special case: Logistic shape, Bernoulli distribution .	7	
	4.6	Combinatoric maximum entropy examples	8	
		4.6.1 Binomial	8	
		4.6.2 Multinomial \ldots 1	10	
		4.6.3 Hypergeometric	12	
	4.7	Continuous maximum entropy examples	13	
		4.7.1 Continuous uniform	14	
		4.7.2 Known mean	14	
		4.7.3 Gaussian (normal) distribution and known variance	14	
		4.7.4 Multivariate normal distribution	15	
		4.7.5 Lognormal distribution	19	
		4.7.6 Logistic distribution	20	

 $^1\mathrm{A}$ table summarizing some maximum entropy probability assignments is found at the end of the chapter (see Park and Bera [2009] for additional maxent assignments).

ii Contents

	4.7.7	Logistic regression as maximum entropy assignment	22
4.8	Maxim	um entropy posterior distributions	27
	4.8.1	Sequential processing for three states	28
	4.8.2	Another frame: sequential maximum relative entropy	31
	4.8.3	Simultaneous processing of moment and data condi-	
		tions	33
4.9	Conver	gence or divergence in beliefs	34
	4.9.1	diverse beliefs	34
	4.9.2	complete ignorance	35
	4.9.3	convergence in beliefs	35
4.10	Appen	dix: summary of maximum entropy probability as-	
	signme	ents	37

This is page 1 Printer: Opaque this

4 Maximum entropy distributions

Bayesian analysis illustrates scientific reasoning where consistent reasoning (in the sense that two individuals with the same background knowledge, evidence, and perceived veracity of the evidence reach the same conclusion) is fundamental. In the next few pages we survey foundational ideas: Bayes' theorem (and its product and sum rules), maximum entropy probability assignment, and consistent updating with maximum entropy probability assignments (posterior probability assignment).

4.1 Bayes' theorem and consistent reasoning

Consistency is the hallmark of scientific reasoning. When we consider probability assignment to events, whether they are marginal, conditional, or joint events their assignments should be mutually consistent (match up with common sense). This is what Bayes' product and sum rules express formally.

The *product rule* says the product of a conditional likelihood (or distribution) and the marginal likelihood (or distribution) of the conditioning variable equals the joint likelihood (or distribution).

$$p(x,y) = p(x|y) p(y)$$
$$= p(y|x) p(x)$$

The sum rule says if we sum over all events related to one (set of) variable(s) (integrate out one variable or set of variables), we are left with the

likelihood (or distribution) of the remaining (set of) variables(s).

$$p(x) = \sum_{i=1}^{n} p(x, y_i)$$
$$p(y) = \sum_{j=1}^{n} p(x_j, y)$$

Bayes' theorem combines these ideas to describe consistent evaluation of evidence. That is, the posterior likelihood associated with a proposition, θ , given the evidence, y, is equal to the product of the conditional likelihood of the evidence given the proposition and marginal or prior likelihood of the conditioning variable (the proposition) scaled by the likelihood of the evidence. Notice, we've simply rewritten the product rule where both sides are divided by p(y) and p(y) is simply the sum rule where θ is integrated out of the joint distribution, $p(\theta, y)$.

$$p(\theta \mid y) = \frac{p(y \mid \theta) p(\theta)}{p(y)}$$

For Bayesian analyses, we often find it convenient to suppress the normalizing factor, p(y), and write the posterior distribution is proportional to the product of the sampling distribution or likelihood function and prior distribution.

$$p(\theta \mid y) \propto p(y \mid \theta) p(\theta)$$

or for a particular draw $y = y_0$

$$p(\theta \mid y = y_0) \propto \ell(\theta \mid y = y_0) p(\theta)$$

where $p(y \mid \theta)$ is the sampling distribution, $\ell(\theta \mid y = y_0)$ is the likelihood function evaluated at $y = y_0$, and $p(\theta)$ is the prior distribution for θ . Bayes' theorem is the glue that holds consistent probability assignment together.

Example 1 (Bayes sum and product rules) Consider the following joint distribution:

$p\left(y=y_1,\theta=\theta_1\right)$	$p(y = y_2, \theta = \theta_1)$	$p\left(y=y_1,\theta=\theta_2\right)$	$p\left(y=y_2,\theta=\theta_2\right)$
0.1	0.4	0.2	0.3

The sum rule yields the following marginal distributions:

$p\left(y=y_1\right)$	$p\left(y=y_2\right)$
0.3	0.7

and

$p\left(\theta=\theta_{1}\right)$	$p\left(\theta=\theta_2\right)$
0.5	0.5

	$p\left(y \mid \theta = \theta_1\right)$	$p\left(y \mid \theta = \theta_2\right)$
y_1	0.2	0.4
y_2	0.8	0.6

The product rule gives the conditional distributions:

and

	$p\left(\theta \mid y = y_1\right)$	$p\left(\theta \mid y = y_2\right)$
θ_1	$\frac{1}{3}$	$\frac{4}{7}$
θ_2	$\frac{2}{3}$	$\frac{3}{7}$

as common sense dictates.

4.2 Maximum entropy distributions

From the above, we see that evaluation of propositions given evidence is entirely determined by the sampling distribution, $p(y | \theta)$, or likelihood function, $\ell(\theta | y)$, and the prior distribution for the proposition, $p(\theta)$. Consequently, assignment of these probabilities is a matter of some considerable import. How do we proceed? Jaynes suggests we take account of our background knowledge, \Im , and evaluate the evidence in a manner consistent with both background knowledge and evidence. That is, the posterior likelihood (or distribution) is more aptly represented by

 $p(\theta \mid y, \Im) \propto p(y \mid \theta, \Im) p(\theta \mid \Im)$

Now, we're looking for a mathematical statement of what we know and only what we know. For this idea to be properly grounded requires a sense of complete ignorance (even though this may never represent our state of background knowledge). For instance, if we think that μ_1 is more likely the mean or expected value than μ_2 then we must not be completely ignorant about the location of the random variable and consistency demands that our probability assignment reflect this knowledge. Further, if the order of events or outcomes is not *exchangeable* (if one permutation is more plausible than another), then the events are not seen as stochastically independent or identically distributed.¹ The mathematical statement of our background knowledge is defined in terms of Shannon's entropy (or sense of diffusion or uncertainty).

¹Exchangeability is foundational for independent and identically distributed events (iid), a cornerstone of inference. However, exchangeability is often invoked in a conditional sense. That is, conditional on a set of variables exchangeability applies — a foundational idea of conditional expectations or regression.

4.3 Entropy

Shannon defines entropy as^2

$$h = -\sum_{i=1}^{n} p_i * \log\left(p_i\right)$$

for discrete events where

$$\sum_{i=1}^{n} p_i = 1$$

or differential entropy as^3

$$h = -\int p(x)\log p(x) \, dx$$

for events with continuous support where

$$\int p(x) \, dx = 1$$

Shannon derived a measure of entropy so that five conditions are satisfied: (1) a measure h exists, (2) the measure is smooth, (3) the measure is monotonically increasing in uncertainty, (4) the measure is consistent in the sense that if different measures exist they lead to the same conclusions, and (5) the measure is additive. Additivity says joint entropy equals the entropy of the signals (y) plus the probability weighted average of entropy conditional on the signals.

$$h(x,y) = H(y) + \Pr(y_1)H(x|y_1) + \dots + \Pr(y_n)H(x|y_n)$$

This latter term, $\Pr(y_1)H(x|y_1)+\cdots+\Pr(y_n)H(x|y_n)$, is called conditional entropy. Internal logical consistency of entropy is maintained primarily via additivity — the analog to the workings of Bayes' theorem for probabilities.

$$h = -\int p(x)\log \frac{p(x)}{m(x)}dx$$

 $^{^{2}}$ The axiomatic development for this measure of entropy can be found in Jaynes [2003] or Accounting and Causal Effects: Econometric Challenges, ch. 13.

³ Jaynes [2003] argues that Shannon's differential entropy logically includes an invariance measure m(x) such that differential entropy is defined as

4.4 Discrete maximum entropy examples⁴

4.4.1 Discrete uniform

Example 2 (discrete uniform) Suppose we know only that there are three possible (exchangeable) events, $x_1 = 1$, $x_2 = 2$, and $x_3 = 3$. The maximum entropy probability assignment is found by solving the Lagrangian

$$\mathcal{L} \equiv \max_{p_i} \left[-\sum_{i=1}^3 p_i \log p_i - (\lambda_0 - 1) \left(\sum_{i=1}^3 p_i - 1 \right) \right]$$

First order conditions yield

$$p_i = e^{-\lambda_0}$$
 for $i = 1, 2, 3$

and

$$\lambda_0 = \log 3$$

Hence, as expected, the maximum entropy probability assignment is a discrete uniform distribution with $p_i = \frac{1}{3}$ for i = 1, 2, 3.

4.4.2 Discrete nonuniform

Example 3 (discrete nonuniform) Now suppose we know a little more. We know the mean is $2.5.^5$ The Lagrangian is now

$$\mathcal{L} \equiv \max_{p_i} \left[-\sum_{i=1}^3 p_i \log p_i - (\lambda_0 - 1) \left(\sum_{i=1}^3 p_i - 1 \right) - \lambda_1 \left(\sum_{i=1}^3 p_i x_i - 2.5 \right) \right]$$

First order conditions yield

$$p_i = e^{-\lambda_0 - x_i \lambda_1}$$
 for $i = 1, 2, 3$

and

$$\lambda_0 = 2.987$$
$$\lambda_1 = -0.834$$

The maximum entropy probability assignment is

$$p_1 = 0.116$$

 $p_2 = 0.268$
 $p_3 = 0.616$

⁴A table summarizing some maximum entropy probability assignments is found at the end of the chapter (see Park and Bera [2009] for additional maxent assignments).

 $^{^5\}mathrm{Clearly},$ if we knew the mean is 2 then we would assign the uniform discrete distribution above.

4.5 Normalization and partition functions

The above analysis suggests a general approach for assigning probabilities where normalization is absorbed into a denominator. Since

$$\exp\left[-\lambda_{0}\right]\sum_{k=1}^{n}\exp\left[-\sum_{j=1}^{m}\lambda_{j}f_{j}\left(x_{i}\right)\right] = 1,$$

$$p\left(x_{i}\right) = \frac{\exp\left[-\lambda_{0}\right]\exp\left[-\sum_{j=1}^{m}\lambda_{j}f_{j}\left(x_{i}\right)\right]}{1}$$

$$= \frac{\exp\left[-\lambda_{0}\right]\exp\left[-\sum_{j=1}^{m}\lambda_{j}f_{j}\left(x_{i}\right)\right]}{\exp\left[-\lambda_{0}\right]\sum_{k=1}^{n}\exp\left[-\sum_{j=1}^{m}\lambda_{j}f_{j}\left(x_{i}\right)\right]}$$

$$\exp\left[-\sum_{i=1}^{m}\lambda_{i}f_{i}\left(x_{i}\right)\right]$$

$$= \frac{\exp\left[-\sum_{j=1}^{m} \lambda_j f_j(x_i)\right]}{\sum_{k=1}^{n} \exp\left[-\sum_{j=1}^{m} \lambda_j f_j(x_i)\right]}$$
$$= \frac{k(x_i)}{Z(\lambda_1, \dots, \lambda_m)}$$

where $f_j(x_i)$ is a function of the random variable, x_i , reflecting what we know,⁷

$$k(x_i) = \exp\left[-\sum_{j=1}^m \lambda_j f_j(x_i)\right]$$

is a kernel, and

$$Z(\lambda_1, \dots, \lambda_m) = \sum_{k=1}^n \exp\left[-\sum_{j=1}^m \lambda_j f_j(x_k)\right]$$

is a normalizing factor, called a partition function.⁸ Probability assignment is completed by determining the Lagrange multipliers, λ_j , $j = 1, \ldots, m$, from the *m* constraints which are function of the random variables.

$$-\frac{\partial \log Z}{\partial \lambda_1} = \mu$$

⁷Since λ_0 simply ensures the probabilities sum to unity and the partition function assures this, we can define the partition function without λ_0 . That is, λ_0 cancels as demonstrated above.

⁸In physical statistical mechanics, the partition function describes the partitioning among different microstates and serves as a generator function for all manner of results regarding a process. The notation, Z, refers to the German word for sum over states, zustandssumme. An example with relevance for our purposes is

Return to the example above. Since we know support and the mean, n = 3 and the function $f(x_i) = x_i$. This implies

$$Z(\lambda_1) = \sum_{i=1}^{3} \exp\left[-\lambda_1 x_i\right]$$

and

$$p_{i} = \frac{k(x_{i})}{Z(\lambda_{1})}$$
$$= \frac{\exp\left[-\lambda_{1}x_{i}\right]}{\sum_{k=1}^{3}\exp\left[-\lambda_{1}x_{k}\right]}$$

where $x_1 = 1$, $x_2 = 2$, and $x_3 = 3$. Now, solving the constraint

$$\sum_{i=1}^{3} p_i x_i - 2.5 = 0$$
$$\sum_{i=1}^{3} \frac{\exp\left[-\lambda_1 x_i\right]}{\sum_{k=1}^{3} \exp\left[-\lambda_1 x_k\right]} x_i - 2.5 = 0$$

produces the multiplier, $\lambda_1 = -0.834$, and identifies the probability assignments

$$p_1 = 0.116$$

 $p_2 = 0.268$
 $p_3 = 0.616$

We utilize the analog to the above partition function approach next to address continuous density assignment as well as a special case involving binary (Bernoulli) probability assignment which takes the shape of a logistic distribution.

4.5.1 Special case: Logistic shape, Bernoulli distribution

Example 4 (Bernoulli) Suppose we know a binary (0,1) variable has mean equal to $\frac{3}{5}$. The maximum entropy probability assignment by the par-

$$-\frac{\partial \log Z}{\partial \lambda_1} = \frac{3 + 2e^{\lambda_1} + e^{2\lambda_1}}{1 + e^{\lambda_1} + e^{2\lambda_1}} = 2.5$$

Solving gives $\lambda_1 = -0.834$ — consistent with the approach below.

where μ is the mean of the distribution. For the example below, we find

tition function approach is

$$p(x) = \frac{\exp \left[-\lambda\left(1\right)\right]}{\exp \left[-\lambda\left(0\right)\right] + \exp \left[-\lambda\left(1\right)\right]}$$
$$= \frac{\exp \left[-\lambda\right]}{1 + \exp \left[-\lambda\right]}$$
$$= \frac{1}{1 + \exp \left[\lambda\right]}$$

This is the shape of the density for a logistic distribution (and would be a logistic density if support were unbounded rather than binary).⁹ Solving for $\lambda = \log \frac{2}{3}$ or $p(x = 1) = \frac{3}{5}$, reveals the assigned probability of success or characteristic parameter for a Bernoulli distribution.

4.6 Combinatoric maximum entropy examples

Combinatorics, the number of exchangeable ways events occur, plays a key role in discrete (countable) probability assignment. The binomial distribution is a fundamental building block.

4.6.1 Binomial

Suppose we know there are binary (Bernoulli or "success"/"failure") outcomes associated with each of n draws and the expected value of success equals np. The expected value of failure is redundant, hence there is only one moment condition. Then, the maximum entropy assignment includes the number of combinations which produce x_1 "successes" and x_0 "failures". Of course, this is the binomial operator, $\binom{n}{x_1,x_0} = \frac{n!}{x_1!x_0!}$ where $x_1 + x_0 = n$. Let $m_i \equiv \binom{n}{x_{1i},x_{0i}}$ and generalize the entropy measure to account for m,

$$s = -\sum_{i=1}^{n} p_i \log\left(\frac{p_i}{m_i}\right)$$

Now, the kernel (drawn from the Lagrangian) is

$$k_i = m_i \exp\left[-\lambda_1 x_i\right]$$

Satisfying the moment condition yields the maximum entropy or binomial probability assignment.

$$p(x,n) = \frac{k_i}{Z} = \begin{pmatrix} n \\ x_{1i}, x_{0i} \end{pmatrix} p^{x_1} (1-p)^{x_0}, \quad x_1 + x_0 = n \\ 0 \quad \text{otherwise}$$

 $^{^9\,{\}rm This}$ suggests logistic regression is a natural (as well as the most common) strategy for modeling discrete choice.

Example 5 (binomial; balanced coin) Suppose we assign $x_1 = 1$ for each coin flip resulting in a head, $x_0 = 1$ for a tail, and the expected value is $E[x_1] = E[x_0] = 6$ in 12 coin flips. We know there are $2^{12} = \sum_{\substack{x_1+x_0=12\\ x_1+x_0=12}} {\binom{12}{x_1,x_0}} = 4,096$ combinations of heads and tails. Solving $\lambda_1 = 0$,

heads and x_0 tails in 12 coin flips is

$$p(x, n = 12) = \begin{pmatrix} 12\\ x_1 \end{pmatrix} \frac{1}{2} \frac{x_1}{2} \frac{1}{2} x_0, \quad x_1 + x_0 = 12\\ 0 \quad otherwise$$

Example 6 (binomial; unbalanced coin) Continue the coin flip example above except heads are twice as likely as tails. In other words, the expected values are $E[x_1] = 8$ and $E[x_0] = 4$ in 12 coin flips. Solving $\lambda_1 = -\log 2$, the maximum entropy probability assignment associated with s heads in 12 coin flips is

$$p(x, n = 12) = \begin{pmatrix} 12\\ x_1 \end{pmatrix} \frac{2}{3}^{x_1} \frac{1}{3}^{x_0}, \quad x_1 + x_0 = 12\\ 0 \quad otherwise$$

	$p\left(x,n=12 ight)$		
$[x_1, x_0]$	balanced coin	unbalanced coin	
[0, 12]	$\frac{1}{4,096}$	$\frac{1}{531,441}$	
[1, 11]	$\frac{12}{4,096}$	$\frac{24}{531,441}$	
[2, 10]	$\frac{66}{4,096}$	$\frac{264}{531,441}$	
[3,9]	$\frac{220}{4,096}$	$\frac{1,760}{531,441}$	
[4, 8]	$\frac{495}{4,096}$	$\frac{7,920}{531,441}$	
[5,7]	$\frac{792}{4,096}$	$\frac{25,344}{531,441}$	
[6,6]	$\frac{924}{4,096}$	$\frac{59,136}{531,441}$	
[7,5]	$\frac{792}{4,096}$	$\frac{101,376}{531,441}$	
[8, 4]	$\frac{495}{4,096}$	$\frac{126,720}{531,441}$	
[9,3]	$\tfrac{220}{4,096}$	$\frac{112,640}{531,441}$	
[10, 2]	$\frac{66}{4,096}$	$\frac{67,584}{531,441}$	
[11, 1]	$\frac{12}{4,096}$	$\frac{24,576}{531,441}$	
[12, 0]	$\frac{1}{4,096}$	$\frac{4,096}{531,441}$	

Canonical analysis

The above is consistent with a canonical analysis based on relative entropy $\sum_{i} p_i \log\left(\frac{p_i}{p_{old,i}}\right)$ where $p_{old,i}$ reflects probabilities assigned based purely on the number of exchangeable ways of generating x_i ; hence, $p_{old,i} = \frac{\binom{n}{x_i}}{\sum_i \binom{n}{x_i}} = \frac{m_i}{\sum_i m_i}$. Since the denominator of $p_{old,i}$ is absorbed via normalization it can be dropped, then entropy reduces to $\sum_i p_i \log\left(\frac{p_i}{m_i}\right)$ and the kernel is the same as above

$$k_i = m_i \exp\left[-\lambda_1 x_{1i}\right]$$

4.6.2 Multinomial

The multinomial is the multivariate analog to the binomial accommodating k rather than two nominal outcomes. Like the binomial, the sum of the outcomes equals n, $\sum_{i=1}^{k} x_i = n$ where $x_k = 1$ for each occurrence of event k. We know there are

$$k^n = \sum_{x_1 + \dots + x_k = n} \binom{n}{x_1, \dots, x_k} = \sum_{x_1 + \dots + x_k = n} \frac{n!}{x_1! \cdots x_k!}$$

possible combinations of k outcomes, $x = [x_1, \dots, x_k]$, in n trials. From here, probability assignment follows in analogous fashion to that for the binomial case. That is, knowledge of the expected values associated with the k events leads to the multinomial distribution when the $\frac{n!}{x_1!\cdots x_k!}$ exchangeable ways for each occurrence x is taken into account, $E[x_j] = np_j$ for $j = 1, \dots, k-1$. Only k-1 moment conditions are employed as the kth moment is a linear combination of the others, $E[x_k] = n - \sum_{j=1}^{k-1} E[x_j] =$ $n\left(1 - \sum_{j=1}^{k-1} p_j\right)$. The kernel is $\frac{n!}{x_1!\cdots x_k!} \exp\left[-\lambda_1 x_1 - \cdots - \lambda_{k-1} x_{k-1}\right]$

which leads to the standard multinomial probability assignment when the moment conditions are resolved

$$p(x,n) = \begin{array}{c} p(x,n) = \frac{n!}{x_1!\cdots x_k!} p_1^{x_1} \cdots p_k^{x_k}, \quad \sum_{i=1}^k x_i = n, \quad \sum_{i=1}^k p_i = 1\\ 0 \quad \text{otherwise} \end{array}$$

Example 7 (multinomial; one balanced die) Suppose we roll a balanced die (k = 6) one time (n = 1), the moment conditions are $E[x_1] = \cdots E[x_6] = \frac{1}{6}$. Incorporating the number of exchangeable ways to generate n = 1 results, the kernel is

$$\frac{1!}{x_1!\cdots x_6!} \exp\left[-\lambda_1 x_1 - \cdots - \lambda_5 x_5\right]$$

Solving for $\lambda_1 = \cdots = \lambda_5 = 0$, yields the multinomial distribution

$$p(x, n = 1) = \frac{1!}{x_1! \cdots x_6!} \frac{1}{6}^{x_1} \cdots \frac{1}{6}^{x_6} = \frac{1}{6}, \quad \sum_{j=1}^k x_j = 1$$

otherwise

Example 8 (multinomial; one unbalanced die) Suppose we roll an unbalanced die (k = 6) one time (n = 1) where the moment conditions are $E[x_1] = \frac{1}{21}, E[x_2] = \frac{2}{21}, E[x_3] = \frac{3}{21}, E[x_4] = \frac{4}{21}, E[x_5] = \frac{5}{21}, E[x_6] = \frac{6}{21}$. Incorporating the number of exchangeable ways to generate n = 1 results, the kernel is

$$\frac{1!}{x_1!\cdots x_6!} \exp\left[-\lambda_1 x_1 - \cdots - \lambda_5 x_5\right]$$

Solving for $\lambda_1 = 1.79176$, $\lambda_2 = 1.09861$, $\lambda_3 = 0.693147$, $\lambda_4 = 0.405465$, and $\lambda_5 = 0.182322$, yields the multinomial distribution

$$p(x, n = 1) = \frac{1!}{x_1! \cdots x_6!} \frac{1}{21} \frac{x_1}{21} \frac{2}{21} \frac{x_2}{21} \frac{3}{21} \frac{x_3}{21} \frac{4}{21} \frac{x_4}{21} \frac{5}{21} \frac{x_5}{21} \frac{6}{21} \frac{x_6}{6}, \quad \sum_{j=1}^k x_j = 1$$
otherwise

Example 9 (multinomial; two balanced dice) Suppose we roll two balanced dice (k = 6) one time (n = 2), the moment conditions are $E[x_1] = \cdots E[x_6] = 2(\frac{1}{6}) = \frac{1}{3}$. The number of combinations is $k^n = 6^2 = 6 + 2 * 15 = 36$, that is, 6 permutations of [2, 0, 0, 0, 0, 0] and 15 permutations of [1, 1, 0, 0, 0, 0] times two orders of the dice. Incorporating the number of exchangeable ways to generate n = 2 results, the kernel is

$$\frac{2!}{x_1!\cdots x_6!}\exp\left[-\lambda_1 x_1-\cdots-\lambda_5 x_5\right]$$

Solving for $\lambda_1 = \cdots = \lambda_5 = 0$, yields the multinomial distribution

$$p(x, n = 2) = \frac{2!}{x_1! \cdots x_6!} \frac{1}{6} \frac{x_1}{16} \cdots \frac{1}{6} \frac{x_6}{6} = \frac{1}{6}, \quad \sum_{j=1}^k x_j = 2$$

otherwise

Example 10 (multinomial; two unbalanced dice) Suppose we roll two unbalanced dice (k = 6) one time (n = 2) where the moment conditions are $E[x_1] = \frac{2}{21}, E[x_2] = \frac{4}{21}, E[x_3] = \frac{6}{21}, E[x_4] = \frac{8}{21}, E[x_5] = \frac{10}{21}, E[x_6] = \frac{12}{21}$. Incorporating the number of exchangeable ways to generate n = 2 results, the kernel is

$$\frac{2!}{x_1!\cdots x_6!} \exp\left[-\lambda_1 x_1 - \cdots - \lambda_5 x_5\right]$$

Solving for $\lambda_1 = 1.79176$, $\lambda_2 = 1.09861$, $\lambda_3 = 0.693147$, $\lambda_4 = 0.405465$, and $\lambda_5 = 0.182322$, yields the multinomial distribution

$$p\left(x,n=2\right) = \begin{array}{cc} \frac{2!}{x_{1}!\cdots x_{6}!} \frac{1}{21} x_{1} \frac{2}{21} x_{2} \frac{3}{21} x_{3} \frac{4}{21} x_{4} \frac{5}{21} x_{5} \frac{6}{21} x_{6}, & \sum_{j=1}^{k} x_{j} = 2 \\ 0 & otherwise \end{array}$$

$[x_1, x_2, x_3, x_4, x_5, x_6]$	p(x, n = 2) balanced dice unbalanced d	
$[x_1, x_2, x_3, x_4, x_5, x_6]$ $[2, 0, 0, 0, 0, 0]$		$\frac{1}{441}$
	$\frac{1}{36}$	
$\left[0,2,0,0,0,0\right]$	$\frac{1}{36}$	$\frac{4}{441}$
$\left[0,0,2,0,0,0\right]$	$\frac{1}{36}$	$\frac{9}{441}$
$\left[0,0,0,2,0,0\right]$	$\frac{1}{36}$	$\frac{16}{441}$
$\left[0,0,0,0,2,0\right]$	$\frac{1}{36}$	$\frac{25}{441}$
$\left[0,0,0,0,0,2\right]$	$\frac{1}{36}$	$\frac{36}{441}$
$\left[1,1,0,0,0,0\right]$	$\frac{1}{18}$	$\frac{4}{441}$
$\left[1,0,1,0,0,0\right]$	$\frac{1}{18}$	$\frac{6}{441}$
$\left[1,0,0,1,0,0\right]$	$\frac{1}{18}$	$\frac{8}{441}$
$\left[1,0,0,0,1,0\right]$	$\frac{1}{18}$	$\frac{10}{441}$
$\left[1,0,0,0,0,1\right]$	$\frac{1}{18}$	$\frac{12}{441}$
$\left[0,1,1,0,0,0\right]$	$\frac{1}{18}$	$\frac{12}{441}$
$\left[0,1,0,1,0,0\right]$	$\frac{1}{18}$	$\frac{16}{441}$
$\left[0,1,0,0,1,0\right]$	$\frac{1}{18}$	$\frac{20}{441}$
$\left[0,1,0,0,0,1\right]$	$\frac{1}{18}$	$\frac{24}{441}$
$\left[0,0,1,1,0,0\right]$	$\frac{1}{18}$	$\frac{24}{441}$
$\left[0,0,1,0,1,0\right]$	$\frac{1}{18}$	$\frac{30}{441}$
$\left[0,0,1,0,0,1\right]$	$\frac{1}{18}$	$\frac{36}{441}$
$\left[0,0,0,1,1,0\right]$	$\frac{1}{18}$	$\frac{40}{441}$
$\left[0,0,0,1,0,1\right]$	$\frac{1}{18}$	$\frac{48}{441}$
$\left[0,0,0,0,1,1\right]$	$\frac{1}{18}$	$\frac{60}{441}$

4.6.3 Hypergeometric

The hypergeometric probability assignment is a purely combinatorics exercise. Suppose we take n draws without replacement from a finite population of N items of which m are the target items (or events) and x is the number of target items drawn. There are $\binom{m}{x}$ ways to draw the targets times $\binom{N-m}{n-x}$ ways to draw nontargets out of $\binom{N}{n}$ ways to make n draws. Hence, our combinatoric measure is

$$\frac{\binom{m}{x}\binom{N-m}{n-x}}{\binom{N}{n}}, \quad x \in \{\max\left(0, m+n-N\right), \min\left(m, n\right)\}$$

and

$$\sum_{x=\max(0,m+n-N)}^{\min(m,n)} \frac{\binom{m}{x}\binom{N-m}{n-x}}{\binom{N}{n}} = 1$$

which completes the probability assignment as there is no scope for moment conditions or maximizing entropy. 10

Example 11 (hypergeometric) Suppose we have an inventory of six items (N = 6) of which one is red, two are blue, and three are yellow. We wish to determine the likelihood of drawing x = 0, 1, or 2 blue items in two draws without replacement (m = 2, n = 2).

$$p(x = 0, n = 2) = \frac{\binom{2}{0}\binom{6-2}{2-0}}{\binom{6}{2}} = \frac{4}{6} \cdot \frac{3}{5} = \frac{6}{15}$$

$$p(x = 1, n = 2) = \frac{\binom{2}{1}\binom{6-2}{2-1}}{\binom{6}{2}}$$
$$= \frac{2}{6} \cdot \frac{4}{5} + \frac{4}{6} \cdot \frac{2}{5} = \frac{8}{15}$$
$$p(x = 2, n = 2) = \frac{\binom{2}{2}\binom{6-2}{2-2}}{\binom{6}{2}}$$

$$= \frac{2}{6} \cdot \frac{1}{5} = \frac{1}{15}$$

4.7 Continuous maximum entropy examples

The partition function approach for continuous support involves density assignment

$$p(x) = \frac{\exp\left[-\sum_{j=1}^{m} \lambda_j f_j(x)\right]}{\int_a^b \exp\left[-\sum_{j=1}^{m} \lambda_j f_j(x)\right] dx}$$

$$k = \binom{m}{x} \binom{N-m}{n-x} \exp\left[0\right]$$

and the partition function is

$$Z = \sum_{x=\max(0,m+n-N)}^{\min(m,n)} k = \binom{N}{n}$$

 $^{^{10}\,{\}rm If}$ we omit $\binom{N}{n}$ in the denominator, it would be captured via normalization. In other words, the kernel is

where support is between a and b.

4.7.1 Continuous uniform

Example 12 (continuous uniform) Suppose we only know support is between zero and three. The above partition function density assignment is simply (there are no constraints so there are no multipliers to identify)

$$p(x) = \frac{\exp[0]}{\int_0^3 \exp[0] \, dx} = \frac{1}{3}$$

Of course, this is the density function for a uniform with support from 0 to 3.

4.7.2 Known mean

Example 13 (truncated exponential) Continue the example above but with known mean equal to 1.35. The partition function density assignment is

$$p(x) = \frac{\exp\left[-\lambda_1 x\right]}{\int_0^3 \exp\left[-\lambda_1 x\right] dx}$$

and the mean constraint is

$$\int_{0}^{3} xp(x) dx - 1.35 = 0$$
$$\int_{0}^{3} x \frac{\exp[-\lambda_{1}x]}{\int_{0}^{3} \exp[-\lambda_{1}x] dx} dx - 1.35 = 0$$

so that $\lambda_1 = 0.2012$, $Z = \int_0^3 \exp[-\lambda_1 x] dx = 2.25225$, and the density function is a truncated exponential distribution with support from 0 to 3.

$$p(x) = 0.444 \exp\left[-0.2012x\right], \quad 0 \le x \le 3$$

The base (non-truncated) distribution is exponential with mean approximately equal to 5 (4.9699).

$$p(x) = \frac{1}{4.9699} \exp\left[-\frac{x}{4.9699}\right], \quad 0 \le x < \infty$$
$$p(x) = 0.2012 \exp\left[-0.2012x\right]$$

4.7.3 Gaussian (normal) distribution and known variance

Example 14 (Gaussian or normal) Suppose we know the average dispersion or variance is $\sigma^2 = 100$. Then, a finite mean must exist, but even if we don't know it, we can find the maximum entropy density function for arbitrary mean, μ . Using the partition function approach above we have

$$p(x) = \frac{\exp\left[-\lambda_2 \left(x-\mu\right)^2\right]}{\int_{-\infty}^{\infty} \exp\left[-\lambda_2 \left(x-\mu\right)^2\right] dx}$$

and the average dispersion constraint is

$$\int_{-\infty}^{\infty} (x-\mu)^2 p(x) dx - 100 = 0$$
$$\int_{-\infty}^{\infty} (x-\mu)^2 \frac{\exp\left[-\lambda_2 (x-\mu)^2\right]}{\int_{-\infty}^{\infty} \exp\left[-\lambda_2 (x-\mu)^2\right] dx} - 100 = 0$$

so that $\lambda_2 = \frac{1}{2\sigma^2}$ and the density function is

$$p(x) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right]$$
$$= \frac{1}{\sqrt{2\pi}10} \exp\left[-\frac{(x-\mu)^2}{200}\right]$$

Of course, this is a Gaussian or normal density function. Strikingly, the Gaussian distribution has greatest entropy of any probability assignment with the same variance.

4.7.4 Multivariate normal distribution

As with the univariate normal distribution, knowledge of the variances (and covariances) of a collection of n random variables leads to a natural (maximum entropy) multivariate normal probability assignment. The kernel accounts for all n(n+1)/2 known moment conditions

$$k(x) = \exp \begin{bmatrix} -\lambda_1 (x_1 - \mu_1)^2 - \dots - \lambda_n (x_n - \mu_n)^2 \\ -\lambda_{n+1} (x_1 - \mu_1) (x_2 - \mu_2) - \dots \\ -\lambda_{n(n+1)/2} (x_{n-1} - \mu_{n-1}) (x_n - \mu_n) \end{bmatrix}$$

and the density function is

$$f\left(x\right) = \frac{k\left(x\right)}{Z\left(x\right)}$$

where

$$Z\left(x\right) = \int k\left(x\right) dx_{1} \cdots dx_{n}$$

Then, the multipliers are found by solving the known moment conditions in the usual manner.

However, there is a computationally simpler approach. Let $Var[x] \equiv \Sigma = \Gamma \Gamma^T$ where the latter is determined via Cholesky decomposition. Let

$$z \equiv \Gamma^{-1} \left(x - \mu \right)$$

$$Var[z] = E\left[\Gamma^{-1}(x-\mu)(x-\mu)^{T}(\Gamma^{-1})^{T}\right]$$
$$= \Gamma^{-1}E\left[(x-\mu)(x-\mu)^{T}\right](\Gamma^{-1})^{T}$$
$$= \Gamma^{-1}\Sigma(\Gamma^{-1})^{T} = I$$

In other words, all random variables in vector z have variance one and covariance zero along with zero mean. Since maximum entropy corresponds to zero covariance, the covariance terms drop out (the multipliers are zero) and the kernel is simply

$$k(z) = \exp\left[-\lambda_1 z_1^2 - \dots - \lambda_n z_n^2\right]$$

As the variances are all equal to one, their multipliers are all equal and we can utilize the univariate normal result indicating the multipliers are $\frac{1}{2}$.

From here, we have the choice of normalizing to find the density function f(z) and transform z to x to find f(x) = |J| f(z) where $|J| = \left| \frac{\partial z_1}{\partial x_1} \cdots \frac{\partial z_1}{\partial x_n} \right|$

 $\begin{vmatrix} \overline{\partial x_1} & \cdots & \overline{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_n}{\partial x_1} & \cdots & \frac{\partial z_n}{\partial x_n} \end{vmatrix}$ refers to the Jacobian

$$f(x) = |J| (2\pi)^{-n/2} \exp \left[-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right]$$
$$= (2\pi)^{-n/2} |\Sigma|^{-\frac{1}{2}} \exp \left[-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu) \right]$$

or rewrite the kernel in terms of x and normalize to recover f(x).

$$k(z) = \exp\left[-\frac{1}{2}z^{T}z\right]$$

=
$$\exp\left[-\frac{1}{2}(x-\mu)^{T}(\Gamma^{-1})^{T}\Gamma^{-1}(x-\mu)\right]$$

=
$$\exp\left[-\frac{1}{2}(x-\mu)^{T}\Sigma^{-1}(x-\mu)\right]$$

This is recognized as the kernel for a multivariate normal distribution with normalizing constant as above $(2\pi)^{-n/2} |\Sigma|^{-\frac{1}{2}}$. We complete this discussion with a bivariate example.

Example 15 (bivarate normal distribution) Suppose our background knowledge indicates

$$Var\left[\begin{array}{c} x_1\\ x_2 \end{array}\right] = \left[\begin{array}{cc} 3 & 1\\ 1 & 2 \end{array}\right]$$

Then, the kernel is

$$k(x) = \exp\left[-\lambda_1 (x_1 - \mu_1)^2 - \lambda_2 (x_2 - \mu_2)^2 - \lambda_3 (x_1 - \mu_1) (x_2 - \mu_2)\right]$$

and the partition function is

$$Z\left(x\right) = \int k\left(x\right) dx_1 dx_2$$

Together, they give

$$f\left(x\right) = \frac{k\left(x\right)}{Z\left(x\right)}$$

Now, we solve for the multipliers utilizing $f(x) = \frac{k(x)}{Z(x)}$.

$$\int f(x) (x_1 - \mu_1)^2 dx_1 dx_2 = 3$$
$$\int f(x) (x_2 - \mu_2)^2 dx_1 dx_2 = 2$$
$$\int f(x) (x_1 - \mu_1) (x_2 - \mu_2) dx_1 dx_2 = 1$$

This yields

$$\lambda_1 = \frac{1}{5}, \lambda_2 = \frac{3}{10}, \lambda_3 = -\frac{1}{5}$$

Substituting the multipliers into f(x) gives the bivariate normal density function.

$$f(x) = \left(2\pi\sqrt{5}\right)^{-1} \exp\left\{-\frac{1}{2(5)} \left[\begin{array}{cc} x_1 - \mu_1 \\ x_2 - \mu_2 \end{array}\right]^T \left[\begin{array}{cc} 2 & -1 \\ -1 & 3 \end{array}\right] \left[\begin{array}{cc} x_1 - \mu_1 \\ x_2 - \mu_2 \end{array}\right]\right\}$$

Alternatively, employ Cholesky decomposition of $\Sigma = \Gamma \Gamma^T$ to transform x into uncorrelated, unit variance (mean zero) random variables z.

$$\begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = \Gamma^{-1} \begin{bmatrix} x_1 - \mu_1 \\ x_2 - \mu_2 \end{bmatrix}$$
$$= \begin{bmatrix} \frac{x_1 - \mu_1}{\sqrt{3}} \\ \sqrt{\frac{3}{5}} (x_2 - \mu_2) - \frac{x_1 - \mu_1}{\sqrt{15}} \end{bmatrix}$$

The kernel for z is

$$k(z) = \exp\left[-\lambda_1 z_1^2 - \lambda_2 z_2^2 - \lambda_3 z_1 z_2\right]$$

and the partition function is

$$Z\left(z\right) = \int k\left(z\right) dz_1 dz_2$$

then the density function for z is

$$f\left(z\right) = \frac{k\left(z\right)}{Z\left(z\right)}$$

Hence, the multipliers are determined from

$$\int f(z) z_1^2 dz_1 dz_2 = 1$$
$$\int f(z) z_2^2 dz_1 dz_2 = 1$$
$$\int f(z) z_1 z_2 dz_1 dz_2 = 0$$

This gives

$$\lambda_1 = \frac{1}{2}, \lambda_2 = \frac{1}{2}, \lambda_3 = 0$$

This is the density function for independent, bivariate standard normal random variables.

$$f(z) = (2\pi)^{-1} \exp\left\{-\frac{1}{2} \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}^T \begin{bmatrix} z_1 \\ z_2 \end{bmatrix}\right\}$$

Transformation of variables yields the density for x.

$$f(x) = \left\{ -\frac{1}{2} \left[\begin{array}{c} \frac{x_1 - \mu_1}{\sqrt{3}} \\ \sqrt{\frac{3}{5}} \left(x_2 - \mu_2\right) - \frac{x_1 - \mu_1}{\sqrt{15}} \end{array} \right]^T \left[\begin{array}{c} \frac{x_1 - \mu_1}{\sqrt{3}} \\ \sqrt{\frac{3}{5}} \left(x_2 - \mu_2\right) - \frac{x_1 - \mu_1}{\sqrt{15}} \end{array} \right] \right\}$$

where the Jacobian is

$$|J| = \begin{vmatrix} \frac{\partial z_1}{\partial x_1} & \frac{\partial z_1}{\partial x_2} \\ \frac{\partial z_2}{\partial x_1} & \frac{\partial z_2}{\partial x_2} \end{vmatrix}$$
$$= \begin{vmatrix} \frac{1}{\sqrt{3}} & 0 \\ -\frac{1}{\sqrt{15}} & \sqrt{\frac{3}{5}} \end{vmatrix} = \sqrt{\frac{1}{5}}$$

Of course, this is the same density function as above.

$$f(x) = \left(2\pi\sqrt{5}\right)^{-1} \exp\left\{-\frac{1}{2(5)} \left[\begin{array}{cc} x_1 - \mu_1 \\ x_2 - \mu_2 \end{array}\right]^T \left[\begin{array}{cc} 2 & -1 \\ -1 & 3 \end{array}\right] \left[\begin{array}{cc} x_1 - \mu_1 \\ x_2 - \mu_2 \end{array}\right]\right\}$$

4.7.5 Lognormal distribution

Example 16 (lognormal) Suppose we know the random variable has positive support (x > 0) then it is natural to work with (natural) logarithms. If we also know $E[\log x] = \mu = 1$ as well as $E[(\log x)^2] = \sigma^2 = 10$, then the maximum entropy probability assignment is the lognormal distribution

$$p(x) = \frac{1}{x\sqrt{2\pi\sigma}} \exp\left[-\frac{(\log x - \mu)^2}{2\sigma^2}\right], \quad 0 < x, \mu, \sigma < \infty$$

Again, we utilize the partition function approach to demonstrate.

$$p(x) = \frac{\exp\left[-\lambda_1 \log x - \lambda_2 \left(\log x\right)^2\right]}{\int_0^\infty \exp\left[-\lambda_1 \log x - \lambda_2 \left(\log x\right)^2\right] dx}$$

and the constraints are

$$E\left[\log x\right] = \int_{0}^{\infty} \log xp\left(x\right) dx - 1 = 0$$

and

$$E\left[\left(\log x\right)^{2}\right] = \int_{0}^{\infty} \left(\log x\right)^{2} p(x) \, dx - 10 = 0$$

so that $\lambda_1 = 0.9$ and $\lambda_2 = \frac{1}{2\sigma^2} = 0.05$. Substitution gives

$$p(x) \propto \exp\left[-\frac{2(9)}{20}\log x - \frac{1}{20}(\log x)^2\right]$$

Completing the square and adding in the constant (from normalization) $\exp\left[-\frac{1}{20}\right]$ gives

$$p(x) \propto \exp\left[-\log x\right] \exp\left[-\frac{1}{20} + \frac{2}{20}\log x - \frac{1}{20}\left(\log x\right)^2\right]$$

Rewriting produces

$$p(x) \propto \exp\left[\log x^{-1}\right] \exp\left[-\frac{\left(\log x - 1\right)^2}{20}\right]$$

which simplifies as

$$p(x) \propto x^{-1} \exp\left[-\frac{(\log x - 1)^2}{20}\right]$$

 $\label{eq:constants} Including \ the \ normalizing \ constants \ yields \ the \ probability \ assignment \ asserted \ above$

$$p(x) = \frac{1}{x\sqrt{2\pi 10}} \exp\left[-\frac{(\log x - 1)^2}{2(10)}\right], \quad 0 < x < \infty$$

4.7.6 Logistic distribution

The logistic or extreme-value distribution is symmetric (1 - F(z) = F(-z)) with cumulative distribution function (cdf)

$$F(z) = \frac{1}{1 + \exp[-z]}$$
$$= \frac{\exp[z]}{1 + \exp[z]}$$

and

$$F(x) = \frac{1}{1 + \exp\left[-\frac{x-\mu}{s}\right]}$$

where $z = \frac{x-\mu}{s}$ with mean μ and variance $\frac{\pi^2 s^2}{3}$, and probability density function (pdf)

$$f(z) = \frac{1}{s}F(z)F(-z)$$

=
$$\frac{\exp[z]}{s(1 + \exp[z])^2}$$

=
$$\frac{1}{s\left(\exp\left[-\frac{z}{2}\right] + \exp\left[\frac{z}{2}\right]\right)^2}$$

and

$$f(x) = \frac{1}{s\left(\exp\left[-\frac{x-\mu}{2s}\right] + \exp\left[\frac{x-\mu}{2s}\right]\right)^2}$$

logistic as posterior distribution

The logistic distribution is the posterior distribution (cdf) following from a binary state and Normally distributed evidence.

$$\Pr(s_1 \mid y = y_0) = \frac{pf(y_0 \mid s_1)}{pf(y_0 \mid s_1) + (1 - p)f(y_0 \mid s_0)}$$

where $p = \Pr(s_1)$, $f(y | s_j) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left[-\frac{(y-\mu_j)^2}{2\sigma^2}\right]$, j = 0, 1. Since the variances of the evidence conditional on the state are the same the normalizing constant, $\frac{1}{\sqrt{2\pi\sigma}}$, cancels in the ratio. Accordingly, substitution produces

$$\Pr(s_1 \mid y = y_0) = \frac{p \exp\left[-\frac{(y_0 - \mu_1)^2}{2\sigma^2}\right]}{p \exp\left[-\frac{(y_0 - \mu_1)^2}{2\sigma^2}\right] + (1 - p) \exp\left[-\frac{(y_0 - \mu_0)^2}{2\sigma^2}\right]}$$

Now, mulitply by $\frac{\frac{1}{p} \exp\left[\frac{(y_0-\mu_1)^2}{2\sigma^2}\right]}{\frac{1}{p} \exp\left[\frac{(y_0-\mu_1)^2}{2\sigma^2}\right]}$ to find

$$\Pr(s_1 \mid y = y_0) = \frac{1}{1 + \frac{1-p}{p} \exp\left[-\frac{(y_0 - \mu_0)^2}{2\sigma^2} + \frac{(y_0 - \mu_1)^2}{2\sigma^2}\right]}$$

Rewrite as log-odds ratio and expand and simplify the exponential term

$$\Pr(s_1 \mid y = y_0) = \frac{1}{1 + \exp\left[\log\left(\frac{1-p}{p}\right)\right] \exp\left[\frac{y_0 - \frac{(\mu_0 + \mu_1)}{2}}{\frac{\sigma^2}{\mu_0 - \mu_1}}\right]}$$

Collecting terms gives

$$\Pr(s_1 \mid y = y_0) = \frac{1}{1 + \exp\left[\frac{y_0 - \frac{(\mu_0 + \mu_1)}{2} + \frac{\sigma^2}{\mu_0 - \mu_1}\log\left(\frac{1-p}{p}\right)}{\frac{\sigma^2}{\mu_0 - \mu_1}}\right]}$$
$$= \frac{1}{1 + \exp\left[\frac{y_0 - \left\{\frac{(\mu_0 + \mu_1)}{2} - \frac{\sigma^2}{\mu_0 - \mu_1}\log\left(\frac{1-p}{p}\right)\right\}}{\frac{\sigma^2}{\mu_0 - \mu_1}}\right]}$$

which is the cdf for a logistic random variable with mean $\frac{(\mu_0 + \mu_1)}{2} - \frac{\sigma^2}{\mu_0 - \mu_1} \log\left(\frac{1-p}{p}\right)$ and scale parameter $s = \frac{\sigma^2}{\mu_0 - \mu_1}$.

Example 17 (logistic distribution) Suppose we know

$$E\left[\log\left(\exp\left[-\frac{z}{2}\right] + \exp\left[\frac{z}{2}\right]\right)\right] = 1$$

where $z = \frac{x-\mu}{s}$ then the kernel is

$$k(z) = \exp\left[-\lambda \log\left(\exp\left[-\frac{z}{2}\right] + \exp\left[\frac{z}{2}\right]\right)\right]$$

the partition function is

$$Z = \int_{-\infty}^{\infty} \exp\left[-\lambda \log\left(\exp\left[-\frac{z}{2}\right] + \exp\left[\frac{z}{2}\right]\right)\right] dz$$

The moment constraint

$$E\left[\log\left(\exp\left[-\frac{z}{2}\right] + \exp\left[\frac{z}{2}\right]\right)\right] = \int_{-\infty}^{\infty} \log\left(\exp\left[-\frac{z}{2}\right] + \exp\left[\frac{z}{2}\right]\right) \frac{k(z)}{Z} dz = 1$$

yields $\lambda = 2$. Hence,

$$f(z) = \frac{k(z; \lambda = 2)}{Z(\lambda = 2)} = \frac{1}{s\left(\exp\left[-\frac{z}{2}\right] + \exp\left[\frac{z}{2}\right]\right)^2}$$

and by transformation of variables we have

$$f(x) = \frac{1}{s\left(\exp\left[-\frac{x-\mu}{2s}\right] + \exp\left[\frac{x-\mu}{2s}\right]\right)^2}$$

the density function for a logistic random variables with mean μ and variance $\frac{\pi^2 s^2}{3}$.

4.7.7 Logistic regression as maximum entropy assignment

For binary choice, the log-odds ratio is

$$\log \frac{p}{1-p}$$

If we employ a logistic link function along with an index function, $X^T \gamma$,¹¹ to describe the conditional probabilities we have

$$\log \frac{\left(1 + \exp\left[-X^T\gamma\right]\right)^{-1}}{\left(1 + \exp\left[X^T\gamma\right]\right)^{-1}} = \log\left(\exp\left[X^T\gamma\right]\right)$$
$$= X^T\gamma$$

which is suggestive of a natural connection for the logistic distribution to binary choice.

From here we can provide a more rigorous argument for the frequent utilization of logistic regression when faced with discrete choice analysis. The logit model for discrete choice D conditional on (regime differences in) covariates X is

$$\Pr(D \mid X) = \frac{1}{1 + \exp[-Y]}$$
$$= \frac{1}{1 + \exp[-X^T \gamma]}$$

Following Blower [2004], we develop this model specification from the maximum entropy principle.

Bayesian revision yields

$$\Pr\left(D \mid X\right) = \frac{\Pr\left(D, X\right)}{\Pr\left(X\right)}$$

and for treatment selection

$$\Pr(D = 1 \mid X) = \frac{\Pr(D = 1, X)}{\Pr(D = 1, X) + \Pr(D = 0, X)}$$

 $^{{}^{11}}X^T\gamma = (X_1 - X_0)^T\gamma$ where X_j describes obervable characteristics associated with individuals when choosing option j while γ describes the relative weight on the various characteristics. X is the difference in the characteristic between the choices.

_

Rewrite this expression as

$$\Pr\left(D=1\mid X\right) = \frac{1}{1 + \frac{\Pr(D=0,X)}{\Pr(D=1,X)}}$$

The maximum entropy probability assignments, denoted \hbar , for the joint likelihoods, $\Pr(D = 1, X)$ and $\Pr(D = 0, X)$, are

$$\Pr\left(D=1, X, \hbar\right) = \frac{\exp\left[\sum_{j=1}^{m} \lambda_j f_j\left(X_1\right)\right]}{Z\left(\lambda_1, \dots, \lambda_m\right)}$$

and

$$\Pr\left(D=0, X, \hbar\right) = \frac{\exp\left[\sum_{j=1}^{m} \lambda_j f_j\left(X_0\right)\right]}{Z\left(\lambda_1, \dots, \lambda_m\right)}$$

_

The likelihood ratio is

$$\frac{\Pr\left(D=0,X,\hbar\right)}{\Pr\left(D=1,X,\hbar\right)} = \frac{\exp\left[\sum_{j=1}^{m}\lambda_{j}f_{j}\left(X_{0}\right)\right]}{\exp\left[\sum_{j=1}^{m}\lambda_{j}f_{j}\left(X_{1}\right)\right]}$$
$$= \exp\left[-Y\right]$$

where

$$Y = \sum_{j=1}^{m} \lambda_j \{ f_j (X_1) - f_j (X_0) \}$$

Hence, we have the logistic regression specification as a maximum entropy probability assignment where the $m f_j(X_1) - f_j(X_0)$ and multipliers λ_j identify observable characteristics related to choice and their regression weights.¹²

$$\Pr(D = 1 \mid X, \hbar) = \frac{1}{1 + \frac{\Pr(D = 0, X, \hbar)}{\Pr(D = 1, X, \hbar)}}$$
$$= \frac{1}{1 + \exp[-Y]}$$

$$D^* = X^T \gamma - \nu$$

¹²For the latent variable random utility model with index structure $(Y = (X_1 - X_0)^T \gamma = X^T \gamma)$ and unobservable ν a logistic random variable we have

Example 18 (logistic regression) Suppose the data generating process (DGP) is described by the following joint distribution for choice, D = 0, 1, and conditions, B = 1, 2, 3 and C = L, M, H.

	$\Pr\left(D=1,B,C\right)$				
	B = 1	B=2	B=3	$\Pr\left(D=1,C\right)$	
C = L	0.008	0.022	0.02	0.05	
C = M	0.009	0.027	0.024	0.06	
C = H	0.013	0.041	0.036	0.09	
$\Pr\left(D=1,B\right)$	0.03	0.09	0.08	1	
$\Pr\left(D=1\right)$					0.2

	B = 1	B=2	B=3	$\Pr\left(D=0,C\right)$	$\Pr(C)$
C = L	0.060	0.145	0.078	0.283	0.333
C = M	0.058	0.140	0.075	0.273	0.333
C = H	0.052	0.125	0.067	0.244	0.334
$\Pr\left(D=0,B\right)$	0.17	0.41	0.22	,	
$\Pr\left(D=0\right)$					0.8
$\Pr\left(B\right)$	0.2	0.5	0.3		

Since we only observe D = 0, 1, we write

$$\Pr(D = 1 \mid X) = \Pr(D^* \ge 0 \mid X)$$
$$= \Pr\left(\nu \le X^T \gamma \mid X\right)$$
$$= F_{\nu}\left(X^T \gamma\right)$$
$$= \frac{1}{1 + \exp\left[-X^T \gamma\right]}$$
$$= \frac{1}{1 + \exp\left[-Y\right]}$$

Suppose we have the following background knowledge.¹³

$$E[D] = 0.2$$

$$E[B_1] = 0.2$$

$$E[B_2] = 0.5$$

$$E[C_L] = 0.333$$

$$E[C_M] = 0.333$$

$$E[DB_1] = 0.03$$

$$E[DB_2] = 0.09$$

$$E[DC_L] = 0.05$$

$$E[DC_M] = 0.06$$

With this background knowledge the kernel determined by maximum entropy is

$$k = \exp \begin{bmatrix} \lambda_1 D + \lambda_2 B_1 + \lambda_3 B_2 + \lambda_4 C_L + \lambda_5 C_M \\ + \lambda_6 D B_1 + \lambda_7 D B_2 + \lambda_8 D C_L + \lambda_9 D C_M \end{bmatrix}$$

where all variables, D, B_1, B_2, C_L , and C_M , are indicator variables. Scaling k by the partition function, $Z = \sum k$, and solving for the multipliers that satisfy the moment conditions reveals

 $^{^{13}\,{\}rm Complete}$ knowledge includes all the interactions as well as the nine moments indicated. The other eight moments are

This produces the following probability assignment.¹⁴

$\Pr\left(D=1,B,C\right)$						
B = 1 $B = 2$ $B = 3$						
C = L	0.0075	0.0225	0.02			
C = M	0.009	0.027	0.024			
C = H	0.0135	0.0405	0.036			

$\Pr\left(D=0,B,C\right)$			
	B = 1	B=2	B=3
C = L	0.0601375	0.1450375	0.077825
C = M	0.0580125	0.1399125	0.075075
C = H	0.05185	0.12505	0.0671

Hence, the maximum entropy conditional probabilities are

$$\Pr(D = 1 \mid B, C) = \frac{\Pr(D = 1, B, C)}{\Pr(D = 1, B, C) + \Pr(D = 0, B, C)}$$

A linear probability model supplies effective starting values for maximum likelihood estimation of a logistic regression. The linear probability model is^{15}

$$Pr(D = 1 | B, C) = X\beta$$

= 0.3342 - 0.1147B₁ - 0.0855B₂ - 0.1178C_L - 0.0881C_H

Logistic regression is

$$\Pr\left(D=1 \mid B, C\right) = \frac{1}{1 + \exp\left[X\gamma\right]}$$

Equating the two expressions and solving gives

$$X\gamma = \log \frac{1 - X\beta}{X\beta}$$

or starting values

$$\gamma_0 = (X^T X)^{-1} X^T \log \frac{1 - X\beta}{X\beta}$$

$$\gamma_0^T = \begin{bmatrix} -0.5991 & -0.7463 & -0.5185 & -0.7649 & -0.5330 \end{bmatrix}$$

Maximum likelihood estimation with a logistic link function gives

$$\gamma^T = \begin{bmatrix} -0.6227 & -0.7230 & -0.5047 & -0.7361 & -0.5178 \end{bmatrix}$$

 $^{^{14}}$ Of course, this doesn't match the DGP as our background knowledge is incomplete. 15 Since the model involves only indicator variables, the predicted values are bounded between zero and one.

As the previous discussion indicates maximum entropy probability assignment and logistic regression give the same probability of D conditional on B and C while the linear probability model differs somewhat.

$\Pr\left(D=1 \mid B, C\right)$	maxent or mle	linear probability model
B_1, C_L	0.1109	0.1018
B_2, C_L	0.1343	0.1310
B_3, C_L	0.2044	0.2164
B_1, C_M	0.1343	0.1314
B_2, C_M	0.1618	0.1607
B_3, C_M	0.2422	0.2461
B_1, C_H	0.2066	0.2196
B_2, C_H	0.2446	0.2488
B_3, C_H	0.3492	0.3342

On the other hand, the complete background information model (with 17 moment conditions) implies a saturated, fully interactive regression model for which maximum entropy probability assignment, maximum likelihood logistic regression, and a linear (ordinary least squares) probability model each reproduce the DGP joint (and conditional) probabilities. That is, illustration of the equivalence of maximum entropy and logistic regression draws from limited background information as all consistent models are equivalent with full background knowledge.

4.8 Maximum entropy posterior distributions

We employ maximum entropy to choose among a priori probability distributions subject to our knowledge of moments (of functions) of the distribution (e.g., the mean). When new evidence is collected, we're typically interested in how this data impacts our posterior beliefs regarding the parameter space. Of course, Bayes' theorem provides guidance. For some problems, we simply combine our maximum entropy prior distribution with the likelihood function to determine the posterior distribution. Equivalently, we can find (again, by Lagrangian methods) the maximum relative entropy posterior distribution conditional first on the moment conditions then conditional on the data so that the data eventually outweighs the moment conditions.

When sequential or simultaneous processing of information produces the same inferences we say the constraints commute, as in standard Bayesian updating. However, when one order of the constraints reflects different information (addresses different questions) than a permutation, the constraints are noncommuting. For noncommuting constraint problems, consistent reasoning demands we modify our approach from the standard one above (Giffin and Caticha [2007]). Moment constraints along with data

constraints are typically noncommuting. Therefore, we augment standard Bayesian analysis with a "canonical" factor. We illustrate the difference between sequential and simultaneous processing of moment and data conditions via a simple three state example.

4.8.1 Sequential processing for three states

Suppose we have a three state process $(\theta_1, \theta_2, \theta_3)$ where, for simplicity,

$$\theta_i \in \{0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8\}$$

and

$$\sum_{i=1}^{3} \theta_i = 1$$

We'll refer to this a die generating process where θ_1 refers to outcome one or two, θ_2 to outcome three or four, and θ_3 corresponds to die face five or six. The maximum entropy prior with no moment conditions is

$$p_{old}(\theta_1, \theta_2, \theta_3) = \frac{1}{\sum_{j=1}^8 j} = \frac{1}{36}$$

for all valid combinations of $(\theta_1, \theta_2, \theta_3)$ and the likelihood function is multinomial

$$p(x \mid \theta) = \frac{(n_1 + n_2 + n_3)!}{n_1! n_2! n_3!} \theta_1^{n_1} \theta_2^{n_2} \theta_3^{n_3}$$

where n_i represents the number of *i* draws in the sample of $n_1 + n_2 + n_3$ total draws.

Now, suppose we know, on average, state one (θ_1) is twice as likely as state three (θ_3) . That is, the process produces "die" with these average properties. The maximum entropy prior given this moment condition is generated from solving

$$\max_{p(\theta)} \quad h = -\sum_{\theta} p(\theta_1, \theta_2, \theta_3) \log p(\theta_1, \theta_2, \theta_3)$$

s.t.
$$\sum_{\theta} p(\theta_1, \theta_2, \theta_3) (\theta_1 - 2\theta_3) = 0$$

Lagrangian methods yield¹⁶

$$p(\theta) = \frac{\exp\left[\lambda\left(\theta_1 - 2\theta_3\right)\right]}{Z}$$

where the partition function is

$$Z = \sum_{\theta} \exp\left[\lambda \left(\theta_1 - 2\theta_3\right)\right]$$

¹⁶ Frequently, we write $p(\theta)$ in place of $p(\theta_1, \theta_2, \theta_3)$ to conserve space.

and λ is the solution to

$$-\frac{\partial \log Z}{\partial \lambda} = 0$$

In other words,

$$p(\theta) = \frac{\exp\left[1.44756\left(\theta_1 - 2\theta_3\right)\right]}{28.7313}$$

and

$$p(x,\theta) = p(x \mid \theta) p(\theta)$$

= $\frac{(n_1 + n_2 + n_3)!}{n_1! n_2! n_3!} \theta_1^{n_1} \theta_2^{n_2} \theta_3^{n_3} \frac{\exp\left[1.44756\left(\theta_1 - 2\theta_3\right)\right]}{28.7313}$

Hence, prior to collection and evaluation of evidence expected values of θ are

$$E[\theta_1] = \sum_{x} \sum_{\theta} \theta_1 p(x, \theta)$$
$$= \sum_{\theta} \theta_1 p(\theta) = 0.444554$$

$$E \left[\theta_{2}\right] = \sum_{x} \sum_{\theta} \theta_{2} p\left(x,\theta\right)$$
$$= \sum_{\theta} \theta_{2} p\left(\theta\right) = 0.333169$$

$$E [\theta_3] = \sum_{x} \sum_{\theta} \theta_3 p(x, \theta)$$
$$= \sum_{\theta} \theta_3 p(\theta) = 0.222277$$

where $E[\theta_1 + \theta_2 + \theta_3] = 1$ and $E[\theta_1] = 2E[\theta_3]$. In other words, probability assignment matches the known moment condition for the process.

Next, we roll one "die" ten times to learn about this specific die and the outcome x is $m = \{n_1 = 5, n_2 = 3, n_3 = 2\}$. The joint probability is

$$p(\theta, x = m) = 2520\theta_1^5\theta_2^3\theta_3^2 \frac{\exp\left[1.44756\left(\theta_1 - 2\theta_3\right)\right]}{28.7313}$$

the probability of the data is

$$p(x) = \sum_{\theta} p(\theta, x = m) = 0.0286946$$

and the posterior probability of θ is

$$p(\theta \mid x = m) = 3056.65\theta_1^5\theta_2^3\theta_3^2 \exp\left[1.44756\left(\theta_1 - 2\theta_3\right)\right]$$

Hence,

$$E \left[\theta_1 \mid x = m \right] = \sum_{\theta} p \left(\theta \mid x = m \right) \theta_1$$
$$= 0.505373$$

$$E \left[\theta_2 \mid x = m\right] = \sum_{\theta} p\left(\theta \mid x = m\right) \theta_2$$
$$= 0.302243$$

$$E \left[\theta_3 \mid x = m\right] = \sum_{\theta} p\left(\theta \mid x = m\right) \theta_3$$

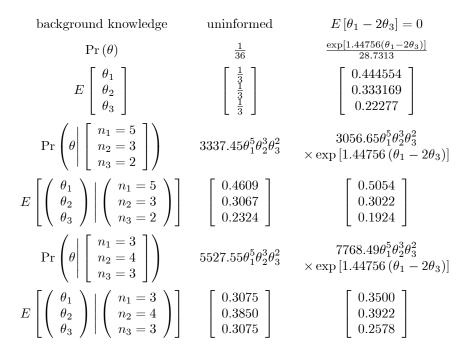
= 0.192384

and

$$E\left[\theta_1 + \theta_2 + \theta_3 \mid x = m\right] = 1$$

However, notice $E \left[\theta_1 - 2\theta_3 \mid x = m\right] \neq 0$. This is because our priors including the moment condition refer to the process and here we're investigating a specific "die".

To help gauge the relative impact of priors and likelihoods on posterior beliefs we tabulate the above results and contrast with uninformed priors as well as alternative evidence. It is important to bear in mind we subscribe to the notion that probability beliefs represent states of knowledge.



Next, we repeat the above question utilizing maximum relative entropy for the joint distribution for θ and x. Then, we consider evaluation of the process; this involves maximization of relative entropy for the joint distribution again but with simultaneous consideration of the moment and data conditions.

4.8.2 Another frame: sequential maximum relative entropy

Now, we repeat the sequential processing of the moment condition for the process followed by processing of the data for a specific die but utilize relative entropy for the joint distribution (as compared with maximization of entropy for the prior distribution of θ above). That is, the problem can be formulated as

$$\max_{p(x,\theta)} \quad s = -\sum_{x} \sum_{\theta} p(x,\theta) \log \frac{p(x,\theta)}{p_{old}(x,\theta)}$$

s.t.
$$\sum_{x} \sum_{\theta} p(x,\theta) (\theta_1 - 2\theta_3) = 0$$

where $p_{old}(x,\theta)$ equals the likelihood $p_{old}(x \mid \theta) = \frac{(n_1+n_2+n_3)!}{n_1!n_2!n_3!} \theta_1^{n_1} \theta_2^{n_2} \theta_3^{n_3}$ times $p_{old}(\theta) = \frac{1}{36}$, or, $p_{old}(x,\theta) = \frac{(n_1+n_2+n_3)!}{36n_1!n_2!n_3!} \theta_1^{n_1} \theta_2^{n_2} \theta_3^{n_3}$. Lagrangian methods yield

$$p(x,\theta) = \frac{\exp\left[\lambda_{\theta} \left(\theta_{1} - 2\theta_{3}\right)\right]}{Z\left(\theta\right)} p_{old}\left(x,\theta\right)$$

where the partition function is

$$Z\left(\theta\right) = \sum_{\theta} \exp\left[\lambda_{\theta} \left(\theta_{1} - 2\theta_{3}\right)\right] p_{old}\left(\theta\right)$$

and λ_{θ} is the solution to

$$-\frac{\partial \log Z\left(\theta\right)}{\partial \lambda_{\theta}} = 0$$

Since $\lambda_{\theta} = \lambda$ in the original frame, we have the same joint distribution given the moment condition as for the original frame

$$p(x,\theta) = \frac{\exp\left[1.44756\left(\theta_1 - 2\theta_3\right)\right]}{28.7313} \frac{(n_1 + n_2 + n_3)!}{n_1! n_2! n_3!} \theta_1^{n_1} \theta_2^{n_2} \theta_3^{n_3}$$

In turn, this implies data for a specific die where outcome x is $m = \{n_1 = 5, n_2 = 3, n_3 = 2\}$ produces identical inferences as above. The joint probability is

$$p(\theta, x = m) = 2520\theta_1^5\theta_2^3\theta_3^2 \frac{\exp\left[1.44756\left(\theta_1 - 2\theta_3\right)\right]}{28.7313}$$

the probability of the data is

$$p(x) = \sum_{\theta} p(\theta, x = m) = 0.0286946$$

and the posterior probability of θ is

$$p(\theta \mid x = m) = 3056.65\theta_1^5\theta_2^3\theta_3^2 \exp\left[1.44756\left(\theta_1 - 2\theta_3\right)\right]$$

Hence,

$$E \left[\theta_1 \mid x = m\right] = \sum_{\theta} p\left(\theta \mid x = m\right) \theta_1$$
$$= 0.505373$$

$$E \left[\theta_2 \mid x = m\right] = \sum_{\theta} p\left(\theta \mid x = m\right) \theta_2$$
$$= 0.302243$$

4.8 Maximum entropy posterior distributions 33

$$E \left[\theta_3 \mid x = m\right] = \sum_{\theta} p\left(\theta \mid x = m\right) \theta_3$$

= 0.192384

and

$$E\left[\theta_1 + \theta_2 + \theta_3 \mid x = m\right] = 1$$

as above.

This maximum relative entropy analysis for the joint distribution helps set up simultaneous evaluation of moment and data conditions when we're evaluating the process rather than a specific die.

4.8.3 Simultaneous processing of moment and data conditions

Suppose we know, on average, the process produces die where, on average, $\theta_1 = 2\theta_3$ and we randomly sample die producing $m = \{n_1 = 5, n_2 = 3, n_3 = 2\}$ to learn more about the process. This is a different question and calls for a different, simultaneous analysis.

The analysis is deceptively similar to the previous analysis but the key difference is the data are considered simultaneously with the moment conditions. That is,

$$\max_{p(x,\theta)} \quad s = -\sum_{x} \sum_{\theta} p(x,\theta) \log \frac{p(x,\theta)}{p_{old}(x,\theta)}$$

s.t.
$$\sum_{x} \sum_{\theta} p(x,\theta) (\theta_1 - 2\theta_3) = 0$$

$$x = m$$

Since the likelihood function remains the same, Lagrangian methods yield

$$p(x = m, \theta) = \frac{\exp\left[\lambda_{(x,\theta)} \left(\theta_1 - 2\theta_3\right)\right]}{Z(m,\theta)} p_{old}(x = m, \theta)$$

where the partition function is

$$Z(m,\theta) = \sum_{\theta} \exp \left[\lambda_{(x,\theta)} \left(\theta_1 - 2\theta_3\right)\right] p_{old} \left(x = m, \theta\right)$$

and $\lambda_{(x,\theta)}$ is the solution to

$$-\frac{\partial \log Z(m,\theta)}{\partial \lambda_{(x,\theta)}} = 0$$

In other words, the joint probability given the moment and data conditions is

$$p(x = m, \theta) = \exp\left[0.0420198(\theta_1 - 2\theta_3)\right] \frac{2520}{36} \theta_1^5 \theta_2^3 \theta_3^2$$

Since $p(x = m) = \sum_{\theta} p(x = m, \theta) = 0.0209722$, the posterior probability of θ is

$$p(\theta \mid x = m, E[\theta_1 - 2\theta_3] = 0) = 70 \frac{\exp\left[0.0420198\left(\theta_1 - 2\theta_3\right)\right]}{0.0209722}$$

Hence,

$$E [\theta_1 | x = m, E [\theta_1 - 2\theta_3] = 0] = \sum_{\theta} p(\theta | x = m, E [\theta_1 - 2\theta_3] = 0) \theta_1$$

= 0.462225

$$E \left[\theta_2 \mid x = m, E \left[\theta_1 - 2\theta_3\right] = 0\right] = \sum_{\theta} p\left(\theta \mid x = m, E \left[\theta_1 - 2\theta_3\right] = 0\right) \theta_2$$

= 0.306663

$$E [\theta_3 | x = m, E [\theta_1 - 2\theta_3] = 0] = \sum_{\theta} p(\theta | x = m, E [\theta_1 - 2\theta_3] = 0) \theta_3$$

= 0.2311125

$$E[\theta_1 + \theta_2 + \theta_3 \mid x = m, E[\theta_1 - 2\theta_3] = 0] = 1$$

and unlike the previous analysis of a specific die, for the process the moment condition is maintained

$$E[\theta_1 - 2\theta_3 \mid x = m, E[\theta_1 - 2\theta_3] = 0] = 0$$

What we've done here is add a canonical term, $\frac{\exp[\lambda_{(x,\theta)}(\theta_1-2\theta_3)]}{Z(m,\theta)}$, to the standard Bayesian posterior for θ given the data, $p_{old}(\theta \mid x = m)$, to account for the moment condition. The partition function, $Z(m,\theta)$, serves to normalize the moment conditioned-posterior distribution.

4.9 Convergence or divergence in beliefs

Probability beliefs derive from background knowledge \Im mapped into prior beliefs and likelihood functions to produce posterior beliefs.

$$\Pr(\theta \mid y, \Im) \propto \ell(\theta \mid y, \Im) \Pr(\theta, \Im)$$

A specific sample or draw maps into a sampling distribution $\Pr(y \mid \theta, \Im)$ to produce a likelihood function $\ell(\theta \mid y, \Im)$.

4.9.1 diverse beliefs

Now, imagine beginning with complete ignorance (uninformed priors) then diversity of posterior beliefs derives entirely from differences in likelihood functions say due to differences in interpretation of the veracity of the evidence and/or asymmetry of information.

4.9.2 complete ignorance

Complete ignorance refers to lack of knowledge regarding location and scale of the random variables of interest.

location

Complete ignorance regarding location translates into f(x) = 1 uniform or rectangular. Uninformed or constant priors means likelihood functions completely determine posterior beliefs regarding location.

scale

Complete ignorance regarding scale translates into $f(x) = \frac{1}{x}$ for x > 0or assigning a uniform distribution to the log of scale, f(y) = 1 where $y = \log(x)$ (note: $\log \sigma^2 = 2 \log \sigma$ so it matters little whether we speak, for instance, of standard deviation, $x = \sigma$, or variance, $x = \sigma^2$, which makes sense as both are indicators of scale and we're depicting ignorance regarding scale). Intuition for this is probability assignment is invariant to choice of units. This translates into f(x) dx = f(bx) d(bx) for some constant b > 0. Since $d(bx) = bdx (\frac{d(bx)}{dx} = b$ which implies d(bx) = bdx), we require f(x) dx = bf(bx) dx and this is only true for $f(x) = \frac{1}{x}$.

$$f(x) dx = \frac{dx}{x}$$
$$bf(bx) dx = \frac{b}{bx} dx = \frac{dx}{x}$$

For example, f(x) = f(bx) = 1 (uniform) leads to $f(x) dx \neq bf(bx) dx$ or $1dx \neq b1dx$, which is not scale invariant. Nor does $f(x) = \exp(-x)$ satisfy scale invariance as $f(x) dx = \exp(-x) dx \neq bf(bx) dx = b \exp(-bx) dx$.

However, it is intuitively appealing to consider scale ignorance as assigning a uniform probability to the log of scale $x, y = \log x, f(y) = 1$. Then, $f(x) = f(y) \left| \frac{dy}{dx} \right| = \frac{1}{x}$. Uninformed priors means likelihood functions completely determine posterior beliefs regarding scale.

The point here is consistency suggests, or rather demands, that divergence in probability beliefs builds from different likelihood functions (processing the information differently and/or processing different information).

4.9.3 convergence in beliefs

Why do probability beliefs converge? Individuals share the same information, individuals agree on the information's veracity, and/or information cascades form. The first two are straightforward extensions of the above discussion. Information cascades arise when individuals regard public processing of information as so compelling that their own private information

leaves their posterior beliefs largely unaffected. Then, collectively individuals are dissuaded from engaging in (private) information search and an information cascade results characterized by herding behavior. Such bubbles are sustained until some compelling information event bursts it. Changing herding behavior or bursting an information bubble involves a likelihood so powerful it overwhelms the public (common) information prior.

distribution	moment constraints	kernel form ¹⁷	mass function
discrete:			
uniform	none	$\exp\left[0\right] = 1$	$\Pr(x_i = i) = \frac{1}{n}, i = 1, \dots, n$
nonuniform	$E\left[x\right]=\mu$	$\exp\left[-\lambda x_i\right]$	$\Pr(x_i = i) = p_i, \\ i = 1, \dots, n$
Bernoulli	$E\left[x\right] = p$	$\exp\left[-\lambda ight]$	$\Pr (x = 1) = p,$ x = (0, 1)
binomial ¹⁸	$E[x \mid n] = np,$ Pr (x = 1 n = 1) = p	$\binom{n}{x} \exp\left[-\lambda x\right]$	$\Pr \left(x = s \mid n \right) = \binom{n}{s} p^s \left(1 - p \right)^{n-s}, s = (0, \dots, n)$
${ m multinomial}^{19}$	$E[x_i \mid n] = np_i,$ $\Pr(x_i = 1 \mid n = 1) = p_i,$ $i = 1, \dots, k - 1$	$\exp\left[\frac{\frac{n!}{x_1!\cdots x_k!}\times}{\left[-\sum_{i=1}^{k-1}\lambda_i x_i\right]}\right]$	$\Pr \left(x_1, \dots, x_k \mid n \right) = \frac{n!}{x_1! \cdots x_k!} p_1^{x_1} \cdots p_k^{x_k}, \\ x_i = (0, \dots, n)$
Poisson ²⁰	$E\left[x\right]=\mu$	$\frac{1}{x!} \exp\left[-\lambda x\right]$	$\Pr(x = s) = \frac{\mu^s}{s!} \exp[-\mu],$ $s = (0, 1, \ldots)$
$geometric^{21}$	$E[x] = \frac{1}{p},$ Pr(x = 1) = p	$\exp\left[-\lambda x ight]$	$\Pr (x = r) = p (1 - p)^{r-1} r = 1, 2, \dots$
negative binomial ²²	$E[x] = \frac{pr}{1-p},$ Pr (x = 1) = p	$\binom{x+r-1}{x} \times \exp\left[-\lambda x\right]$	$\Pr (x = s; r) = {\binom{s+r-1}{s}} \\ \times p^s (1-p)^r, \\ s = (0, 1,)$
$logarithmic^{23}$	$E[x] = -\frac{p}{(1-p)\log(1-p)},$ $E[\log x] = \sum_{\substack{x=1\\ \delta t}}^{\infty} -\frac{p^x \log x}{x\log(1-p)}$ $= \frac{\frac{\partial}{\partial t} \left(\sum_{\substack{k=1\\ k^t}}^{\infty} \frac{p^k}{k^t}\right) _{t=1}}{\log(1-p)}$	$\exp\left[-\lambda_1 x - \lambda_2 \log x\right]$	$\Pr(x) = -\frac{p^x}{x \log(1-p)},$ $x = (1, 2, \ldots)$
hyper- geometric ²⁴	none (note: $E[x] = \frac{nm}{N}$)	$\frac{\binom{m}{x}\binom{N-m}{n-x}}{\binom{N}{n}}$	$\Pr \left({x = s;m,n,N} \right) = \frac{{\binom{m}{s}\binom{N-m}{n-s}}}{\binom{N}{n}}$

Appendix: summary of maximum entropy 4.10 probability assignments

 $^{^{17}}$ excluding partition function

$$\begin{split} \Gamma\left(z\right) &= \int_{0}^{\infty} e^{-t} t^{z-1} dt \\ \Gamma\left(n\right) &= (n-1)! \text{ for } n \text{ a positive integer} \\ \Gamma'\left(z\right) &= \Gamma\left(z\right) \frac{dLog(\Gamma(z))}{dz} \\ B\left(a,b\right) &= \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)} \end{split}$$

¹⁸ The kernel for the binomial includes a measure of the number of exchangeable ways to generate x successes in n trials, $\binom{n}{x}$. The measure, say m(x), derives from generalized entropy, $S = -\sum_{i} p \log \frac{p_i}{m_i}$, where m_i reflects a measure that ensures entropy is invariant under transformation (or change in units) as required to consistently capture background information including complete ignorance (see Jeffreys [1939], Jaynes [2003] or Sivia and Skilling [2006]). The first order condition from the Lagrangian with mean moment condition yields the kernel $m_i \exp [-\lambda x_i]$, as reflected for the binomial probability assignment. Since $m_i \neq \frac{1}{n}$ for the binomial distribution, m_i is not absorbed through normalization.

¹⁹ Analogous to the binomial, the kernel for the multinomial includes a measure of the number of exchangeable ways to generate x_1, \ldots, x_k events in *n* trials, $\frac{n!}{x_1!\cdots x_k!}$, where events are mutually exclusive (as with the binomial) and draws are with replacement.

²⁰Like the binomial, the kernel for the Poisson includes an invariance measure based on the number of exchangeable permutations for generating x occurrences in a given interval, $\frac{n^x}{x!}$. Since a Poisson process resembles a binomial process with a large number of trials within a fixed interval, we can think of n^x as the large n approximation to $\frac{n!}{(n-x)!}$. n^x is absorbed via $n^x p^x = \left(\frac{n\lambda}{n}\right)^x$ where expected values for the binomial and Poisson are equated, $np = \lambda$ implies $p = \frac{\lambda}{n}$. The Poisson distribution is a Taylor series expansion of $\exp[\mu]$ around zero (which is equal to $\sum_{x=0}^{\infty} \frac{\lambda^x}{x!}$) and normalized via $\exp[-\mu]$ (which is equal to $\lim_{n\to\infty} \left(1-\frac{\lambda}{n}\right)^{n-x}$). Hence, $\lim_{n\to\infty} \binom{n}{n} p^x (1-p)^{n-x} = \frac{e^{-\lambda}\lambda^x}{x!}$. ²¹The geometric distribution indicates the likelihood success occurs in the *r*th trial. Hence, the measure for the number of ways for this to occur is one.

 22 Like the binomial, the kernel for the negative binomial includes a measure of the number of exchangeable ways to generate x successes before r failures occur, $\binom{x+r-1}{x}$. 23 Like the Poisson distribution, the logarithmic distribution is a Taylor series expan-

sion. The expansion involves $-\log [1-p]$ around zero and normalized via $-\frac{1}{\log[1-p]}$. ²⁴Like the binomial, the hypergeometric distribution includes a measure of the number

of exchangeable ways to generate x successes in n trials from a finite population of size N containing m successes without replacement, $\frac{\binom{m}{x}\binom{N-m}{n-x}}{\binom{N}{n}}$. However, no moment conditions are needed.

distribution	moment constraints	kernel form	density function
continuous:			
uniform	none	$\exp\left[0\right] = 1$	$f(x) = \frac{1}{b-a},$ a < x < b
exponential	$E\left[x\right] = \mu > 0$	$\exp\left[-\lambda x ight]$	$f(x) = \frac{1}{\mu} \exp\left[-\frac{x}{\mu}\right], \\ 0 < x < \infty$
gamma	E[x] = a > 0, $E[\log x] = \frac{\Gamma'(a)}{\Gamma(a)}$	$\exp\left[-\lambda_1 x - \lambda_2 \log x\right]$	$f(x) = \frac{1}{\Gamma(a)} \exp\left[-x\right] x^{a-1},$ $0 < x < \infty$
$\begin{array}{c} \text{chi-squared}^{25} \\ \nu \text{ d.f.} \end{array}$	$E[x] = \nu > 0,$ $E[\log x] = \frac{\Gamma'(\frac{1}{2})}{\Gamma(\frac{1}{2})}$ $+ \log 2$	$\exp\left[-\lambda_1 x - \lambda_2 \log x\right]$	$\begin{split} f\left(x\right) &= \frac{1}{2^{\nu/2} \Gamma\left(\nu/2\right)} \\ \exp\left[-\frac{x}{2}\right] x^{\nu/2-1}, \\ 0 &< x < \infty \end{split}$
$beta^{26}$	$E [\log x] = \frac{\Gamma'(a)}{\Gamma[a]} - \frac{\Gamma'(a+b)}{\Gamma[a+b]}, E [\log (1-x)] = \frac{\Gamma'(b)}{\Gamma(b)} - \frac{\Gamma'(a+b)}{\Gamma[a+b]}, a, b > 0$	$\exp\left[\begin{array}{c} -\lambda_1 \log x\\ -\lambda_2 \log \left(1-x\right) \end{array}\right]$	$f(x) = \frac{1}{B(a,b)} x^{a-1} (1-x)^{b-1}, 0 < x < 1$
normal or Gaussian	$E\left[\left(x-\mu\right)^2\right] = \sigma^2$	$\exp\left[-\lambda\left(x-\mu\right)^2\right]$	$f(x) = \frac{1}{\sqrt{2\pi\sigma}} \exp\left[-\frac{(x-\mu)^2}{2\sigma^2}\right],$ $-\infty < x < \infty$
student's t	$E\left[\log\left(\nu+x^2 ight) ight]$	$\left(\nu + x^2\right)^{-\lambda}$	$f(x) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\sqrt{\nu\pi\Gamma\left(\frac{\nu}{2}\right)}} \\ \left(1 + \frac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}} \\ -\infty < x < \infty$
lognormal	$E\left[\log x\right] = \mu$ $E\left[\left(\log x\right)^2\right] = \sigma^2 + \mu^2$	$\exp\left[\begin{array}{c}-\lambda_1\log x\\-\lambda_2\left(\log x\right)^2\end{array}\right]$	$f(x) = \frac{1}{x\sqrt{2\pi}\sigma}$ $\exp\left[-\frac{(\log x - \mu)^2}{2\sigma^2}\right],$ $0 < x < \infty$
Pareto	$E\left[\log x\right] = \frac{1 + \alpha \log x_m}{\alpha},$ $\alpha > 0$	$\exp\left[-\lambda\log x\right]$	$f(x) = \frac{\alpha x_m^{\alpha}}{x^{\alpha+1}}, \\ 0 < x_m \le x < \infty$
Laplace	$\begin{split} E\left[x \right] &= \frac{1}{\beta}, \\ \beta &> 0 \end{split}$	$\exp\left[-\lambda\left x\right \right]$	$f(x) = \frac{\beta}{2} \exp \left[-\beta x \right], \\ -\infty < x < \infty$

 $^{^{25}\}mathrm{special}$ case of gamma distribution

distribution	moment constraints	kernel form	density function
logistic	$E\left[\log\left(\begin{array}{c} \exp\left[-\frac{z}{2}\right]\\ +\exp\left[\frac{z}{2}\right] \end{array}\right)\right] = 1$	$\exp\left[-\lambda \log \left(\begin{array}{c} \exp\left[-\frac{z}{2}\right] \\ +\exp\left[\frac{z}{2}\right] \end{array}\right)\right]$	$\frac{f(x) = \frac{1}{s(\exp[-\frac{z}{2}] + \exp[\frac{z}{2}])^2}}{z = \frac{x - \mu}{s}}$
Wishart $n ext{ d.f.}$	$E [tr (\Sigma)] = ntr (\Psi)$ n > p - 1 $\Sigma \text{ symmetric,}$ positive definite $E [\log \Sigma] = \log \Psi $ $+\Omega (n), \text{ real}$ where $\Omega (n)$ is a function of n	$\exp\left[\begin{array}{c} -\lambda_1 tr\left(\Sigma\right)\\ -\lambda_2 \log \Sigma \end{array}\right]$	$f(\Sigma) = \frac{ \Psi ^{-\frac{n}{2}} \Sigma ^{\frac{n-p-1}{2}}}{2^{\frac{np}{2}} \Gamma(\frac{n}{2})}$ $\exp\left[-\frac{tr(\Psi^{-1}\Sigma)}{2}\right]$

²⁶ The Dirichlet distribution is the multivariate extension of the beta distribution with maxent moment conditions involving $E [\log x_j]$ for all j.