

An Introduction to Maximum Entropy Methods

Marcos Costa Santos Carreira

Global Uncertainty Reading Group

20-Apr-2023

Contents

- 1 What do we want today?
- 2 Entropy
- 3 Maximum Entropy - Discrete
- 4 Maximum Entropy - Continuous
- 5 Maximum Ignorance Probability
- 6 Tail Risk Constraints and Maximum Entropy
- 7 Multivariate extensions
- 8 To do
- 9 Conclusions

A new season

This time there's not much previous material

- Nassim's paper with Helyette and Donald Geman
- papers and books on Maximum Entropy

Let's work it out as we go

Definitions

- Discrete (always non-negative):

$$h(p) = - \sum_{i=1}^n p_i \cdot \log(p_i)$$

- Continuous (Differential Entropy, can be negative):

$$h(p) = - \int_I p(x) \cdot \log(p(x))$$

Information gain

- There is a “right” choice/answer
- We want to reduce uncertainty with the lower effort (number of features) possible
- We want the decision chain that reduces entropy by the most according to some cost function (average number of questions, max number of questions, weighted cost of asking questions, etc.)

Wordle

- One choice of known words, test words reduce set of possibilities
- Easier to look at the maximum information gain for each round
- Best test word (specially on the 3rd round) might not be among the possible answers
- Useful for familiarization with these concepts (program a game solver is more fun than playing it afterwards)

Principle of Maximum Entropy

- Choose between (families of) possible densities
- Constraints (mean, variance, ...)
- The answer
 - Follows these constraints
 - And has the maximum possible entropy

Brandeis Dice (Jaynes)

- Repeated throws of a die that might be weighted
- Given the average x of a large number of throws, what is the expected probability for each number on the next toss?

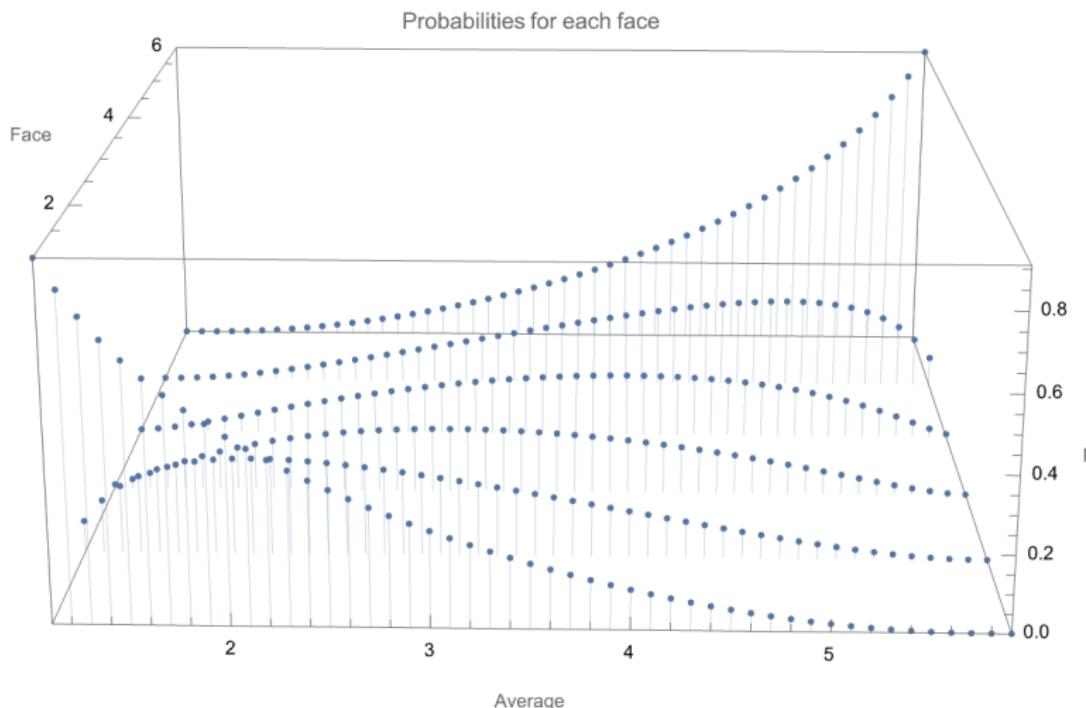
Mathematica code

```
: H[x_] := Maximize[{- \sum_{i=1}^6 (p_i Log[p_i]), \sum_{i=1}^6 (p_i) == 1, \sum_{i=1}^6 (i p_i) == x}, Table[p_i, {i, 1, 6}]];

: H[4.5]
: {1.61359, {p1 \rightarrow 0.054354, p2 \rightarrow 0.0787725, p3 \rightarrow 0.114161, p4 \rightarrow 0.165447, p5 \rightarrow 0.239774, p6 \rightarrow 0.347492}}

: N[H[3]]
: {1.74851, {p1 \rightarrow 0.246782, p2 \rightarrow 0.20724, p3 \rightarrow 0.174034, p4 \rightarrow 0.146148, p5 \rightarrow 0.122731, p6 \rightarrow 0.103065}}
```

Plot



Two approaches (Conrad)

- Given a distribution q , which constraints / extra information on p are needed
- Maximize the entropy with Lagrange multipliers and find the (family of) distributions

First approach

- Look at the integral

$$-\int_I p(x) \cdot \log(q(x))$$

- Substitute the desired distribution in place of q (e.g.

$$q(x) = \frac{\exp(-x/\lambda)}{\lambda}):$$

$$-\int_0^\infty [p(x) \cdot \log(q(x))] dx = \log(\lambda) + \frac{1}{\lambda} \cdot \int_0^\infty [x \cdot p(x)] dx$$

- The last term is the mean of p

Examples

- From Conrad

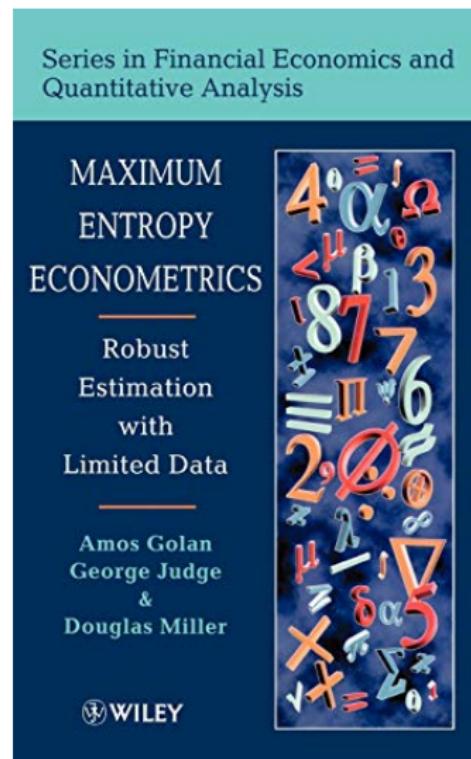
Distribution	Domain	Fixed Value
Uniform	Finite Set	None
Normal with mean μ	\mathbf{R}	$\int_{\mathbf{R}} (x - \mu)^2 p(x) dx$
Exponential	$(0, \infty)$	$\int_0^\infty x p(x) dx$

TABLE 1. Extra constraints

- Wikipedia: https://en.wikipedia.org/wiki/Maximum_entropy_probability_distribution

Lagrange multipliers

- This is the approach of Golan:



Lagrange multipliers

$$H(p) = -\sum_{k=1}^K p_k \ln p_k \rightarrow -p' \ln p$$

$p' \ln p = 0 \text{ for } p \neq 0$

$$p' \cdot 1 = 1 \quad y = x_p$$

Max $H(p)$

$$\sum_{k=1}^K p_k f_k(x_k) = y_k \quad \{f_1(x), f_2(x), \dots, f_K(x)\}$$

$$L = -\sum_{k=1}^K p_k \ln p_k + \sum_{k=1}^K \lambda_k \left[y_k - \sum_{i=1}^K p_i f_i(x_k) \right] + \mu \left(1 - \sum_{k=1}^K p_k \right)$$

$$\frac{\partial L}{\partial p_k} = -\ln p_k - 1 - \sum_{i=1}^K \lambda_i f_i(x_k) - \mu = 0 \quad k \in \{1, \dots, K\}$$

$$\frac{\partial L}{\partial \lambda_k} = y_k - \sum_{i=1}^K \lambda_i f_i(x_k) = 0 \quad k \in \{1, \dots, K\}$$

$$\frac{\partial L}{\partial \mu} = 1 - \sum_{k=1}^K p_k = 0$$

$$\hat{p}_k = e^{-\sum_{i=1}^K \lambda_i f_i(x_k)} / \sum_{i=1}^K e^{-\sum_{j=1}^K \lambda_j f_j(x_k)}$$

$$\sum_{k=1}^K e^{-\sum_{i=1}^K \lambda_i f_i(x_k)} = y_k$$

$$\sum_{k=1}^K e^{-\sum_{i=1}^K \lambda_i f_i(x_k)} / \sum_{i=1}^K e^{-\sum_{j=1}^K \lambda_j f_j(x_k)} = 1$$

CE: minimize $\int p(x) \ln \left[\frac{p(x)}{q(x)} \right] dx = \int p(x) \ln [p(x)] dx - \int p(x) \ln [q(x)] dx$

$$\int p(x) f_k(x) dx = y_k \quad ; \quad \int p(x) dx = 1$$

Lagrange multipliers

$$\begin{aligned}
 H(x) &= - \int f(x) \ln[f(x)] dx \\
 \int f(x) dx &= 1 \\
 L &= - \int f(x) \ln[f(x)] dx + \mu (1 - \int f(x) dx) \\
 &= \mu + \int \left[-f(x) \ln[f(x)] - \mu f(x) \right] dx \\
 &\quad \xrightarrow{\frac{d}{dx}} 0 \\
 &+ \int \left[-f'(x) \ln[f(x)] - \mu f'(x) \right] dx = 0 \\
 &\left(\ln[f(x)] + 1 + \frac{\mu}{\mu} \right) f'(x) = 0 \\
 &\xrightarrow{-1-\frac{\mu}{\mu}} f'(x) = e
 \end{aligned}$$

$$\begin{aligned}
 &x \geq 0: \\
 &\int x f(x) dx = a \\
 L &= - \int f(x) \ln[f(x)] dx + \mu (1 - \int f(x) dx) + \lambda (a - \int x f(x) dx) \\
 L &= \lambda a + \mu - \int_0^\infty \left[f(x) \ln[f(x)] + \mu f(x) + \lambda x f(x) \right] dx \\
 &- \frac{d}{dx} \left\{ f(x) \ln[f(x)] + \mu f(x) + \lambda x f(x) \right\} = 0 \\
 &- \left(\ln f(x) + 1 + \frac{\mu}{\mu} + \lambda x \right) = 0 \\
 f(x) &= \frac{e^{-x}}{a} \quad \text{EXPON. DIST.} \quad \text{MEAN } a \\
 &\text{VARI. } a^2
 \end{aligned}$$

Missing data

- How can we estimate failure rates if we have never met one?
- Can we have enough flexibility in Bayesian priors (ie, not constrained by the properties of a beta distribution)?
- The probability p of success has a probability distribution - we want it to be the maximum entropy

Nassim's choice

- For MaxEnt q must be equal to $1/2$:

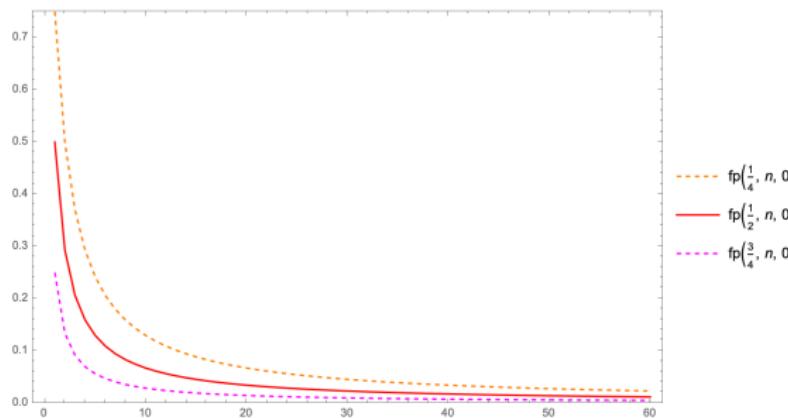
```
: Refine[CDF[BinomialDistribution[n, p], x], 0 ≤ x < n && x ∈ Integers]
: BetaRegularized[1 - p, n - x, 1 + x]
: f[p_] := BetaRegularized[1 - p, n - x, 1 + x];
: InverseFunction[f][q]
  ↪ InverseFunction: Inverse functions are being used. Values may be lost for multivalued inverses. ⓘ
: 1 - InverseBetaRegularized[q, n - x, 1 + x]
: InverseBetaRegularized[ $\frac{1}{2}$ , n - x, 1 + x] // TraditionalForm
'retionalForm=

$$\frac{\Gamma_1^1(n - x, x + 1)}{2}$$

: fp[α_, n_, x_] := 1 - InverseBetaRegularized[α, n - x, 1 + x];
```

Plot

- Bands for different values of q :



Other possibilities

- Bayesian methods with grids (non-parametric)?

High water marks

- If you want to build a house, you don't ask what's the average level of the river (or the variance, etc.)
- You look for high water marks and for how long the region has been inhabited



Geman, Geman and Taleb (2014)

- VaR Constraints
- Add one or some of the following constraints:
 - Global Mean
 - Absolute Mean
 - Power Laws for the tail
- Multi-Period?

Other possibilities

- Price of OTM options
- Multi-period (the Joseph problem)

Blind corners

- Two words no one likes to read: “Previously uncorrelated”
- Bring in the insurers!
 - NFL games cancellation
- 60-40 portfolio in 2022

To do

- Work through basic distributions with the 2 methods (Exponential, Gaussian, StudentT)
- Look more closely at Geman, Geman and Taleb (Var or CVaR?)
- Think about time series (Is the data censored? Is regime switching good?)
- The doctor problem: Can non-parametrical Bayesian methods work?
- Mutual Information
- Cross entropy
- Multivariate

What this talk was about anyway?

An introduction to maximum entropy methods

- Understand entropy and **use it in practice**
- Understand properties of distributions and which **maximum entropy constraint they relate to**
- Understand **the value of the absence of a failure**

Books, papers, posts

- “Maximum Ignorance Probability, with application to surgery’s error rates” (Taleb’s blog)
- Geman, Geman and Taleb (2014), “Tail Risk Constraints and Maximum Entropy” arXiv:1412.7647v1
- Golan, Judge and Miller (1996), “Maximum Entropy Econometrics: Robust Estimation with Limited Data”
- Keith Conrad: “Probability Distributions and Maximum Entropy”
<https://kconrad.math.uconn.edu/blurbs/analysis/entropyppost.pdf>