

REVISITED BOX COUNTING TECHNIQUE IN BAYESIAN SENSE

Václav Hubata–Vacek¹, Jaromír Kukal¹

¹CTU in Prague, Faculty of Nuclear Sciences and Physical Engineering
Department of Software Engineering in Economics
Břehová 7, 115 19 Prague 1
Czech Republic
hubatvac@fjfi.cvut.cz

Abstract:

Keywords: *unbiased estimation*

1 Introduction

2 Multinomic Distribution and Naive Entropy Estimates

Multinomic distribution model plays main role in investigation of point set structures. Let $n \in \mathbb{N}$ be number of distinguish events. Let $p_j > 0$ be probability of j^{th} event for $j = 1, \dots, n$ satisfying $\sum_{j=1}^n p_j = 1$. Then random variable j has multinomic distribution $\text{Mul}(p_1, \dots, p_n)$. After realization of multinomic distribution sample of size $N \in \mathbb{N}$, we can count the events and obtain $N_j \in \mathbb{N}_0$ as number of j^{th} event occurrences for $j = 1, \dots, n$ satisfying $\sum_{j=1}^n N_j = N$. Therefore, we define number of various events in sample as $K = \sum_{N_j > 0} 1 \leq \min(n, N)$. Remembering Hartley and Shannon entropies definitions as

$$H_0 = \ln n, \quad (1)$$

$$H_1 = - \sum_{j=1}^n p_j \ln p_j, \quad (2)$$

we can perform direct but naive estimation of them as

$$H_{0,\text{NAIVE}} = \ln K, \quad (3)$$

$$H_{1,\text{NAIVE}} = - \sum_{N_j > 0} \frac{N_j}{N} \ln \frac{N_j}{N}. \quad (4)$$

The main disadvantage of naive estimates is their biasness. Random variable $K = \{1, \dots, n\}$ is upper constrained by n , then $\mathbb{E}H_{0,\text{NAIVE}} = \mathbb{E} \ln K < \mathbb{E} \ln n = \ln n = H_0$. Therefore, naive estimate of Hartley entropy $H_{0,\text{NAIVE}}$ is negative biased. On the other hand, traditional Box Counting Technique is based on this estimate because we plot logarithm of covering element number $C(a) \in \mathbb{N}$ against logarithm of covering element size $a > 0$ and then estimate their dependency in linear form $\ln C(a) = A_0 - D_{0,\text{NAIVE}} \ln a$. Recognizing equivalence $C(a) = K$, we obtain $\ln C(a) = \ln K = H_{0,\text{NAIVE}}$ and then $H_{0,\text{NAIVE}} = A_0 - D_{0,\text{NAIVE}} \ln a$. Defining $D_{0,\text{NAIVE}}$ as estimate of capacity dimension and recognizing the occurrence of $H_{0,\text{NAIVE}}$ in Box Counting procedure, we are not surprised to be victims of the bias of Hartley entropy estimate.

Similar situation is the case of Shannon entropy estimation. There are several approaches how to declare the bias of $H_{1,\text{NAIVE}}$ to be closer to Shannon entropy H_1 . Miller [2] modified naive estimate $H_{1,\text{NAIVE}}$ using first order Taylor expansion, which produces

$$H_{1,\text{M}} = H_{1,\text{NAIVE}} + \frac{K-1}{2N}. \quad (5)$$

Lately, Harris [2] improved the formula to

$$H_{1,\text{H}} = H_{1,\text{NAIVE}} + \frac{K-1}{2N} + \frac{1}{12N^2} \left(1 - \sum_{p_j > 0} \frac{1}{p_j} \right) \quad (6)$$

Finally, we can estimate information dimension according to relation

$$H_{1,EST} = A_1 - D_{1,EST} \ln a \quad (7)$$

where $H_{1,EST}$ is any estimate of H_1 . Therefore, we can also estimate Hausdorff dimension D_H using inequalities $D_1 \leq D_H \leq D_0$ and then also supposing $D_{1,EST} \leq D_H \leq D_{0,EST}$ for any "good" estimates $D_{0,EST}$, $D_{1,EST}$ of capacity and information dimensions. Next section is oriented to Bayesian estimation of H_0 , H_1 for $D_{0,EST}$ and $D_{1,EST}$ evaluations.

3 Bayesian Estimation of Hartley Entropy

We suppose uniform distribution of random vector $\vec{p} = (p_1, \dots, p_n)$ satisfying $p_j \geq 0$, $\sum_{j=1}^n p_j = 1$. Using properties of multinomic and Dirichlet distributions, we can calculate density $p(K|n, N)$ of random variable $K \in \mathbb{N}$ for $K \leq \min(n, N)$ as

$$p(K|n, N) = \text{prob} \left(\sum_{N_j > 0} 1 = K \middle| n, \sum_{j=1}^n N_j = N \right) = \frac{\binom{n}{K} \binom{N-1}{K-1}}{\binom{N+n-1}{n-1}}. \quad (8)$$

When $N \geq K + 2$, we can calculate

$$S_{K,N} = \sum_{n=K}^{\infty} p(K|n, N). \quad (9)$$

Using inequality

$$\begin{aligned} p(K|n, N) &= \frac{N!(N-1)!}{K!(K-1)!(N-K)!} \frac{n!(n-1)!}{(n-K)!(n+N-1)!} = \\ &= q(K, N) \frac{n(n-1)\dots(n-K+1)}{(n+N-1)(n+N-2)\dots n} \leq q(K, N) \frac{n^K}{n^N} \end{aligned} \quad (10)$$

we can overestimate

$$S_{K,N} \leq \sum_{n=K}^{\infty} q(K, N) n^{K-N} = q(K, N) \sum_{n=K}^{\infty} n^{K-N} < +\infty \quad (11)$$

and then recognize the convergence of infinite series (11). Having a knowledge of K, N where $N \geq K + 2$, we can calculate bayesian density

$$p(n|K, N) = \frac{p(K|n, N)}{S_{K,N}} \quad (12)$$

for $n \geq K$. Therefore, Bayesian estimate of Hartley entropy is

$$H_{0,BAYES} = EH_0 = \sum_{n=K}^{\infty} p(n|K, N) \ln n > \ln k \quad (13)$$

which is also convergent sum. Substituing $n = K + j$ we obtain equivalent formula

$$H_{0,BAYES} = \frac{\sum_{j=0}^{\infty} b_j \ln(K+j)}{\sum_{j=0}^{\infty} b_j} \quad (14)$$

where $b_j = \frac{\binom{K+j}{j} \binom{K+j-1}{j}}{\binom{K+j+N-1}{j}}$.

Then

$$\begin{aligned} b_0 &= 1, \\ b_1 &= \frac{(K+1)K}{K+N}, \\ b_2 &= \frac{1}{2} \frac{(K+2)(K+1)^2 K}{(K+N+1)(K+N)}, \\ b_3 &= \frac{1}{6} \frac{(K+3)(K+2)^2 (K+1)^2 K}{(K+N+2)(K+N+1)(K+N)}, \\ b_4 &= \frac{1}{24} \frac{(K+4)(K+3)^2 (K+2)^2 (K+1)^2 K}{(K+N+3)(K+N+2)(K+N+1)(K+N)}. \end{aligned} \quad (15)$$

Asymptotic properties of Bayesian estimate for $N \rightarrow +\infty$ can be investigated via limits

$$\begin{aligned}
\lim_{N \rightarrow +\infty} H_{0,\text{BAYES}} &= \ln K, \\
\lim_{N \rightarrow +\infty} (H_{0,\text{BAYES}} - \ln K)N &= K(K+1) \ln(1+1/K), \\
\lim_{N \rightarrow +\infty} \left(H_{0,\text{BAYES}} - \ln K - \frac{K(K+1) \ln(1+1/K)}{N} \right) N^2 &= \\
\frac{1}{2} (K(K+2)(K+1) (\ln(K+2) - \ln(K) - 2K \ln(K+1) + K \ln(K+2) + K \ln(K))), &
\end{aligned} \tag{16}$$

Therefore

$$\begin{aligned}
H_{0,\text{BAYES}} \approx \ln K + \frac{K(K+1) \ln(1+1/K)}{N} + \\
\frac{(K(K+2)(K+1) (\ln(K+2) - \ln(K) - 2K \ln(K+1) + K \ln(K+2) + K \ln(K)))}{2N^2}
\end{aligned} \tag{17}$$

When K is also large, we can roughly approximate Hartley entropy as

$$H_{0,\text{BAYES}} \approx \ln K + \frac{K+1}{N} \tag{18}$$

which is very similar to Miller correction [2] in the case of Shannon entropy estimation.

4 Bayesian estimation of Shannon entropy

In the case when the number of events n is known, we can perform Bayesian estimation of Shannon entropy as

$$H_{1,n} = EH_1(K=m) = - \sum_{j=1}^m \left(\frac{N_j+1}{N+m} (\psi(N_j+2) - \psi(N+m+1)) \right) \tag{19}$$

as derived in [TODO]. But when the number of events n is unknown, we can use K as lower estimate of n and perform final Bayesian estimation as

$$H_{1,\text{BAYES}} = \sum_{n=K}^{\infty} p(n|K, N) H_{1,n} \tag{20}$$

which is also convergent sum for $N \geq K+2$. Asymptotic properties of (20) for $N \rightarrow +\infty$ can be investigated by the same technique as in previous section. Resulting asymptotic formula is

$$H_{1,\text{BAYES}} \approx c_0 + \frac{c_1}{N} + \frac{c_2}{N^2} \tag{21}$$

where

$$\begin{aligned}
c_0 &= \sum_{N_j > 0} \frac{N_j}{N} \ln \frac{N}{N_j}, \\
c_1 &= \text{TODO}, \\
c_2 &= \text{TODO}.
\end{aligned} \tag{22}$$

Here c_0 is naive estimate of Shannon entropy, c_1/N corresponds with Miller estimate [TODO] but the term c_2/N^2 differs from Harris estimate [TODO]. The main advantage of formulas TODO,TODO is in absence of theoretical probability knowledge.

5 Revisited Box Counting

6 Experimental parth

7 Conclusion

Acknowledgement: The paper was created with the support of CTU in Prague, Grant SGS11/165/OHK4/3T/14.

References

- [1] todo t.todo. todo.
- [2] Harris, B., *The statistical estimation of entropy in the non-parametric case*. MRC Technical Summary Report, 1975

Let $\mathbb{Q}_n = \{\vec{q} \in (\mathbb{R}_0^+)^n \mid \sum_{j=1}^n q_j = 1\}$ be support set for uniform random variable $\vec{p} \in \mathbb{Q}_n$. Then for integer K satisfying $1 \leq K \leq \min(n, N)$, the conditional probability is

$$p(K|n, N) = \text{prob} \left(\sum_{j=1}^n (N_j > 0) = K \middle| n, \sum_{j=1}^n N_j = N \right). \quad (23)$$

The vector of N_j can be reorganized to begin with positive values. Therefore

$$p(K|n, N) = \binom{n}{K} \text{prob} \left(N_j > 0 \Leftrightarrow j \leq K \middle| n, \sum_{j=1}^K N_j = N \right). \quad (24)$$

Let $\mathbb{D}_{K,N} = \{\vec{x} \in \mathbb{N}^K \mid \sum_{j=1}^K x_j = N\}$ be domain of $\vec{N} = (N_1, \dots, N_K) \in \mathbb{D}_{K,N}$. Using mean value of multinomic distribution over \mathbb{Q}_n , we obtain

$$p(K|n, N) = \binom{n}{K} \mathbb{E} \left(\sum_{\vec{N} \in \mathbb{D}_{K,N}} \binom{N}{N_1, \dots, N_K} \prod_{j=1}^K p_j^{N_j} \prod_{j=K+1}^n p_j^0 \right) = \binom{n}{K} \mathbb{E} \left(\sum_{\vec{N} \in \mathbb{D}_{K,N}} \binom{N}{N_1, \dots, N_K} \mathbb{E} \prod_{j=1}^K p_j^{N_j} \right). \quad (25)$$

Using generalized Beta function

$$B(\vec{x}) = \int_{\vec{p} \in \mathbb{Q}_m} \prod_{j=1}^m p_j^{x_j-1} d\vec{p} = \frac{\prod_{j=1}^m \Gamma(x_j)}{\Gamma(\sum_{j=1}^m x_j)}, \quad (26)$$

we can calculate

$$\mathbb{E} \left(\prod_{j=1}^K p_j^{N_j} \right) = \frac{\int_{\vec{p} \in \mathbb{Q}_n} \prod_{j=1}^K p_j^{N_j} d\vec{p}}{\int_{\vec{p} \in \mathbb{Q}_n} d\vec{p}} = \frac{\prod_{j=1}^K \Gamma(N_j + 1)}{\Gamma(N + n)} \quad (27)$$

Therefore

$$\text{prob}(K|n, N) = \binom{n}{K} \sum_{\vec{N} \in \mathbb{D}_{K,N}} \frac{N!(n-1)!}{\prod_{j=1}^K N_j!} \frac{\prod_{j=1}^K N_j!}{(N+n-1)!} = \binom{n}{K} \sum_{\vec{N} \in \mathbb{D}_{K,N}} \frac{N!(n-1)!}{(N+n-1)!} = \frac{\binom{n}{K} \text{card}(\mathbb{D}_{K,N})}{\binom{N+n-1}{n-1}}. \quad (28)$$

The last question is about the cardinality of $\mathbb{D}_{K,N}$, which corresponds with number of possibilities, how to place N balls into k boxes under assumptions that no box is empty and the balls are identic. We place k balls into k different boxes in the first phase. The rest of $N - k$ balls can be distributed without any constraints. Therefore

$$\text{card}(\mathbb{D}_{K,N}) = \binom{(N-K) + K - 1}{K-1} = \binom{N-1}{K-1}. \quad (29)$$

Resulting formula is (8).

A Hash function

function