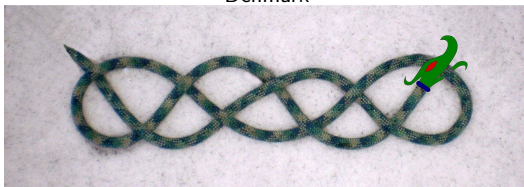


Inequalities for the Binomial Distributions

Peter Harremoës

Copenhagen Business College
Denmark



Workshop on Probability and Information Theory, Hong Kong
2019

Thanks to my coauthors



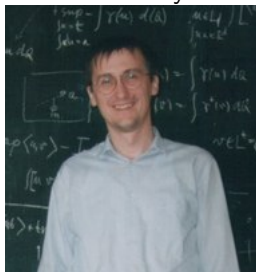
László Györfi



Oliver Johnson



Ioannis Kontoyiannis



František Matuš



Pavel Ruzankin



Gábor Tusnády

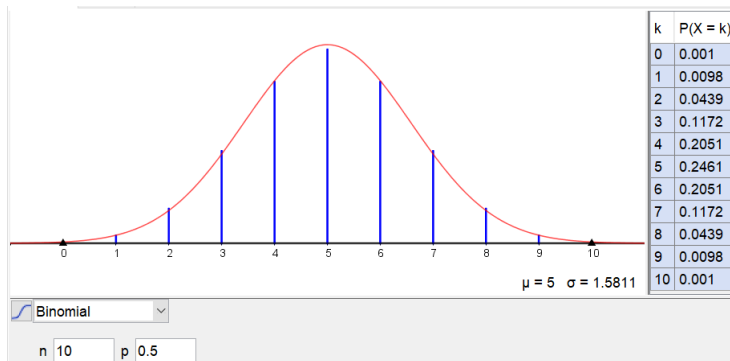


Christophe Vignat

What is the problem?

The random variable X is binomial if

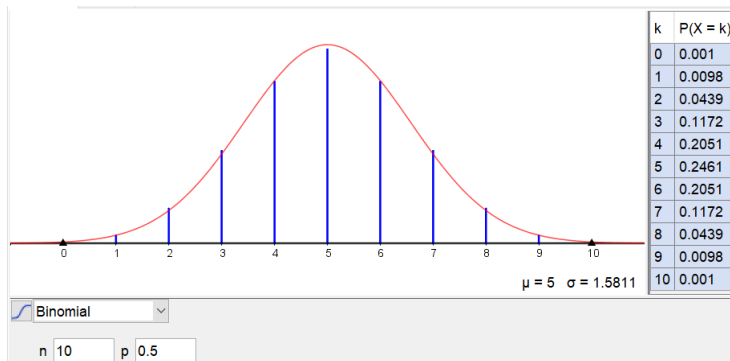
$$\Pr(X = j) = \binom{n}{j} p^j (1 - p)^{n-j}.$$



What is the problem?

The random variable X is binomial if

$$\Pr(X = j) = \binom{n}{j} p^j (1 - p)^{n-j}.$$



- Often n or p are not known.

The binomial distribution and its cousins

- Hypergeometric distribution.
- Bernoulli sum.
- Poisson distribution.
- Negative binomial distribtuion.
- Gaussian distribution.
- Multinomial distribution.

Maximum entropy

Let $B_n(\lambda)$ denote the set of distributions of sums $S_n = X_1 + X_2 + \cdots + X_n$ with mean λ where X_i is a Bernoulli random variable with $\Pr(X_i = 1) = p_i$.

Lemma (Shepp and Olkin 1978, E. Hillion and O. Johnson 2015)

The map $(p_1, p_2, \dots, p_n) \rightarrow H(S_n)$ is concave.

Theorem (PH 2001)

The $H(P)$ entropy restricted to $P \in B_n(\lambda)$ has maximum when $p_i = \lambda/n$, i.e. when P is $\text{bin}(n, \lambda/n)$.

Let $B_\infty(\lambda) = \text{cl}(\bigcup B_n(\lambda))$.

Corollary (PH 2001)

The entropy restricted to $B_\infty(\lambda)$ has maximum at $\text{Po}(\lambda)$. Further $H(\text{bin}(n, \lambda/n)) \rightarrow H(\text{Po}(\lambda))$ for $n \rightarrow \infty$.

Universal coding interpretation

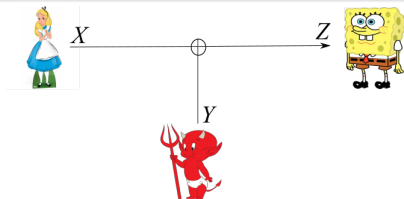
Assume that we are going to code a data point in \mathbb{N} that are generated by some $P \in B_n(\lambda)$, but the exact distribution P is unknown. The code $\kappa : \mathbb{N} \rightarrow A^*$ is characterized by a the code length function $j \rightarrow |\kappa(j)|$ satisfying Kraft's inequality $\sum_j a^{|\kappa(j)|} \leq 1$ where $a = |A|$. The goal is to minimize the maximum mean code length.

$$\min_{\kappa} \max_P E_P(|\kappa(j)|).$$

The solution is $|\kappa(n)| = -\log(\text{bin}(n, p, j))$, i.e. use the code that is optimal if we knew $P = \text{bin}(n, \lambda/n)$.

Similarly, assume that we are going to code a data point in \mathbb{N} that are generated by some $P \in B_n(\lambda)$, but both P and n are unknown. The it is optimal to code as if $P = P_o(\lambda)$.

Relation to the Poisson channel



- The goal for Alice is to maximize $I(X, Z)$ over $X \in B_\infty(\lambda)$.
- The goal for the devil is to minimize $I(X, Z)$ over $Y \in B_\infty(\mu)$.

$$\begin{aligned} I(X, Z) &= H(X + Y) - H(X + Y | X) \\ &= H(X + Y) - H(Y | X) \\ &= H(X + Y) - H(Y). \end{aligned}$$

For any Y it is optimal for Alice to choose $X \sim Po(\lambda)$. If $X \sim Po(\lambda)$ then it is optimal for the devil to choose $Z \sim Po(\mu)$ [PH and C. Vignat, 2003].

Entropy power inequality

Theorem ([PH and C. Vignat 2004])

Assume that $X \sim \text{bin}(m, 1/2)$ and $Y \sim \text{bin}(n, 1/2)$. Then

$$e^{2H(X)} + e^{2H(Y)} \leq e^{2H(X+Y)}.$$

For $X \sim \text{bin}(m, p)$ and $Y \sim \text{bin}(n, q)$ the inequality does not hold for small values of m, n 😞

but it holds for sufficiently large values of m, n [N. Sharma, S. Das, S. Muthukrishnan, 2010].

Entropy power inequality

Theorem ([PH and C. Vignat 2004])

Assume that $X \sim \text{bin}(m, 1/2)$ and $Y \sim \text{bin}(n, 1/2)$. Then

$$e^{2H(X)} + e^{2H(Y)} \leq e^{2H(X+Y)}.$$

For $X \sim \text{bin}(m, p)$ and $Y \sim \text{bin}(n, q)$ the inequality does not hold for small values of m, n 😞

but it holds for sufficiently large values of m, n [N. Sharma, S. Das, S. Muthukrishnan, 2010].



Bernoulli sum and hypergeometric distributions

For $P \in B_n(\lambda)$ we have

$$H(P) + D(P \| \text{bin}(n, \lambda/n)) \leq H(\text{bin}(n, \lambda/n))$$

so if $H(P_k) \rightarrow H_{\max}(B_n(\lambda))$ for $k \rightarrow \infty$ then
 $D(P_n \| \text{bin}(n, \lambda/n)) \rightarrow 0$ for $k \rightarrow \infty$.

Law of small numbers

Since $\text{bin}(n, \lambda/n) \in B_{\infty}(\lambda)$ we have

$$H(\text{bin}(n, \lambda/n)) + D(\text{bin}(n, \lambda/n) \| \text{Po}(\lambda)) \leq H(\text{Po}(\lambda))$$

so

$$H(\text{bin}(n, \lambda/n)) = H_{\max}(B_n(\lambda)) \rightarrow H_{\max}(B_{\infty}(\lambda))$$

for $k \rightarrow \infty$ then $D(\text{bin}(n, \lambda/n) \| \text{Po}(\lambda)) \rightarrow 0$ for $k \rightarrow \infty$.

Upper bounds on total variation

[Babour and Hall, 1984] has

$$\begin{aligned}\frac{1}{16} \min \{p, np^2\} &\leq V(\text{bin}(n, p), \text{Po}(\lambda)) \\ &\leq 2 \min \{p, np^2\}\end{aligned}$$

Upper bounds on total variation

[Babour and Hall, 1984] has

$$\begin{aligned}\frac{1}{16} \min \{p, np^2\} &\leq V(\text{bin}(n, p), \text{Po}(\lambda)) \\ &\leq 2 \min \{p, np^2\}\end{aligned}$$

A factor of **32** in difference between upper and lower bound ☹

Bounds on divergence

We have $D(P\|Q) = \sum f\left(\frac{p_i}{q_i}\right) \cdot q_i$ where $f(x) = x \ln(x)$. For

$$x - 1 \leq f(x) \leq x - 1 + (x - 1)^2.$$

Some better bound

$$\begin{aligned} x - 1 + \frac{1}{2}(x - 1)^2 - \frac{1}{6}(x - 1)^3 &\leq f(x) \\ &\leq x - 1 + \frac{1}{2}(x - 1)^2 - \frac{1}{6}(x - 1)^3 + \frac{1}{3}(x - 1)^4. \end{aligned}$$

$$D(P\|Q) \leq \chi^2(P, Q),$$

$$D(P\|Q) \approx \frac{1}{2}\chi^2(P, Q).$$

Orthogonal polynomials

Assume that f_0, f_1, f_2, \dots are orthogonal normalized polynomials with respect to Q . Then

$$\begin{aligned}\frac{dP}{dQ}(x) &= \sum_{i=0}^{\infty} f_i(x) \cdot \left\langle f_i \middle| \frac{dP}{dQ} \right\rangle, \\ \left\langle f_i \middle| \frac{dP}{dQ} \right\rangle &= \int f_i(x) \frac{dP}{dQ}(x) dQx \\ &= E_P[f_i(X)].\end{aligned}$$

Therefore

$$\chi^2(P, Q) = \sum_{i=1}^{\infty} (E_P[f_i(X)])^2.$$

Upper bounds on divergence

We have

$$\begin{aligned} D(\text{bin}(n, p) \| \text{Po}(\lambda)) &= \sum_{j=0}^n \ln \left(\frac{\text{bin}(n, p, j)}{\text{Po}(\lambda, j)} \right) \cdot \text{bin}(n, p, j) \\ &= \sum_{j=0}^n \ln \left(\frac{\binom{n}{j} p^j (1-p)^{n-j}}{\frac{\lambda^j}{j!} e^{-\lambda}} \right) \cdot \text{bin}(n, p, j) \\ &= \sum_{j=0}^n \left(\lambda + (n-j) \ln(1-p) + \ln \left(\frac{n!}{n^j} \right) \right) \cdot \text{bin}(n, p, j) \\ &= \lambda + (n-\lambda) \ln(1-p) + \sum_{j=0}^n \left(\ln \left(\prod_{i=0}^{j-1} \left(1 - \frac{i}{n} \right) \right) \right) \cdot \text{bin}(n, p, j) . \end{aligned}$$

Stirling numbers

Expand

$$\begin{aligned}\ln \left(\prod_{i=0}^{j-1} \left(1 - \frac{i}{n} \right) \right) &= \sum_{i=0}^j \ln \left(1 - \frac{i}{n} \right) \\ &= - \sum_{i=0}^j \sum_{k=1}^{\infty} \frac{1}{k} \cdot \left(\frac{j}{n} \right)^k.\end{aligned}$$

Introduce Stirling numbers

$$\begin{aligned}j^\ell &= \sum_{m=1}^{\ell} j^{\underline{\ell}} \left\{ \begin{matrix} \ell \\ m \end{matrix} \right\}, \\ j_{[m]} &= \sum_{m=0}^{\ell} j^{\underline{\ell}} \left[\begin{matrix} \ell \\ m \end{matrix} \right].\end{aligned}$$

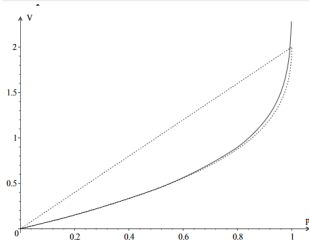
Truncations of these identities leads to inequalities.

Upper bounds

Theorem (PH and P. Ruzankin 2005)

For all

$$D(\text{bin}(n, p) \| \text{Po}(\lambda)) \leq -\frac{\ln(1-p) + p}{2} + \frac{p^2}{12n(1-p)} + \frac{p^2(2 + 11p + 11p^2)}{12n^2(1-p)^5}.$$



Observe that $\limsup n^2 \cdot D(\text{bin}(n, p) \| \text{Po}(\lambda)) \leq \lambda^2/4$.

Theorem

If $\lambda = np$ then

$$D(\text{bin}(n, p) \| \text{Po}(\lambda)) \geq \frac{p^2}{4}.$$

Key observation: Assume that $S_n \sim \text{bin}(n, p)$ and $Y \sim \text{Po}(\lambda)$ where $\lambda = np$. Then

$$E[S_n] = E[Y]$$

and

$$\begin{aligned} \text{Var}(S_n) &= np(1 - p) \\ &< np \\ &= \text{Var}(Y). \end{aligned}$$

Improved rate of convergence

Theorem

Let $Po_{\beta}(\lambda)$ denote the information projection of $Po(\lambda)$ on the set of distributions with the same 1st and 2nd moment as $\text{bin}(n, \lambda/n)$. Then

$$n^2 \cdot D(\text{bin}(n, \lambda/n) \| Po_{\beta}(\lambda)) \rightarrow 0$$

for $n \rightarrow \infty$.

Proof.

We have

$$\begin{aligned} D(\text{bin}(n, p) \| Po(\lambda)) &= D(\text{bin}(n, p) \| Po_{\beta}(\lambda)) + D(Po_{\beta}(\lambda) \| Po(\lambda)) \\ &\geq D(\text{bin}(n, p) \| Po_{\beta}(\lambda)) + \frac{p^2}{4} \end{aligned}$$

Multiply both sides by n^2 . □

Poisson Charlier polynomials

The orthogonal polynomials with respect to $Po(\lambda)$ are

$$C_k^\lambda(x) = (\lambda k!)^{-1/2} \sum_{\ell=0}^k \binom{k}{\ell} (-\lambda)^{k-\ell} x^\ell$$

If $E[X] = \lambda$ then

$$E[C_2^\lambda(X)] = \frac{\text{Var}(X) - \lambda}{2^{1/2}\lambda}$$

Conjecture For any random variable with $E[C_k^\lambda(X)] \leq 0$ we have

$$D(X \| Po(\lambda)) \geq \frac{1}{2} \left(E[C_k^\lambda(X)] \right)^2$$

The conjecture has been proved for $k = 1, 2$ and for any value of k when $E[C_k^\lambda(X)]$ is small [PH, Johnson and Kontoyiannis 2015].

Hypergeometric distributions and Bernoulli sums

A hypergeometric distribution is given by

$$\Pr(X = j) = \frac{\binom{K}{j} \binom{N-K}{n-j}}{\binom{N}{n}}$$

Then there exist p_1, p_2, \dots, p_n such that

$$\Pr(X = j) = \Pr(S_n = j)$$

where $S_n = \sum_{i=1}^n X_i$ is a Bernoulli sum and $\Pr(X_i = 1) = p_i$. The mean is $E[S_n] = \sum p_i$. Then $\text{bin}(n, \bar{p})$ has the same means as S_n if $\bar{p} = \frac{\sum p_i}{n}$. The variance is

$$\begin{aligned} \text{Var}(S_n) &= \sum p_i (1 - p_i) \\ &\leq n\bar{p}(1 - \bar{p}) \\ &= \text{Var}(\text{bin}(n, \bar{p})) \end{aligned}$$

Kravchuk polynomials

The Kravchuk polynomials $\tilde{K}(n, x)$ are orthogonal with respect to $\text{bin}(n, p)$. are

$$C_k^\lambda(x) = (\lambda k!)^{-1/2} \sum_{\ell=0}^k \binom{k}{\ell} (-\lambda)^{k-\ell} x^\ell$$

If $E[X] = \lambda$ then

$$E[C_2^\lambda(X)] = \frac{\text{Var}(X) - \lambda}{2^{1/2} \lambda}$$

Conjecture For any random variable with $E[\tilde{K}_k(X)] \leq 0$ we have

$$D(X \| \text{bin}(n, p)) \geq \frac{1}{2} \left(E[\tilde{K}_k(X)] \right)^2$$

The conjecture has been proved for $k = 1, 2$ and for any value of k when $E[C_k^\lambda(X)]$ is small [PH and F. Matúš, 2019].

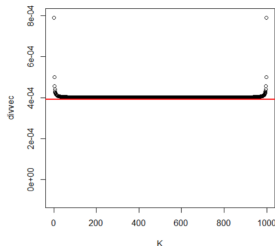
Lower bound for hypergeometric distributions

The hypergeometric distribution satisfies

$$D\left(\text{hyp}(N, K, n) \parallel \text{bin}\left(n, \frac{K}{N}\right)\right) \geq \frac{n(n-1)}{4(N-1)^2}.$$

This result confirms the rule of thumb:

Assume independence when sample size is less than 5 % of population size.



Upper bound for hypergeometric distributions

Stam 1978 proved

$$D\left(\text{hyp}(N, K, n) \parallel \text{bin}\left(n, \frac{K}{N}\right)\right) \leq \frac{n(n-1)}{2(N-1)(N-n+1)}.$$

Upper bound for hypergeometric distributions

Stam 1978 proved

$$D\left(\text{hyp}(N, K, n) \parallel \text{bin}\left(n, \frac{K}{N}\right)\right) \leq \frac{n(n-1)}{2(N-1)(N-n+1)}.$$

By taking higher order terms into account we get

$$D\left(\text{hyp}(N, K, n) \parallel \text{bin}\left(n, \frac{K}{N}\right)\right) \leq \frac{N \ln \frac{N-1/2}{N-n-3/2} - n + \frac{N}{N-n-1}}{N-1}.$$

Let $N(\mu, \sigma^2)$ denote a Gaussian with mean μ and standard deviation σ . Then

$$D\left(N(\lambda, \sigma^2) \parallel N(\mu, \sigma^2)\right) = \frac{(\lambda - \mu)^2}{2\sigma^2}.$$

For the binomial distributions we have

$$D\left(\text{bin}(n, p) \parallel \text{bin}(n, q)\right) = n \left(p \ln \frac{p}{q} + (1 - p) \ln \frac{1 - p}{1 - q} \right).$$

Let (P^λ) denote elements of an exponential family in its mean value parametrization. Define

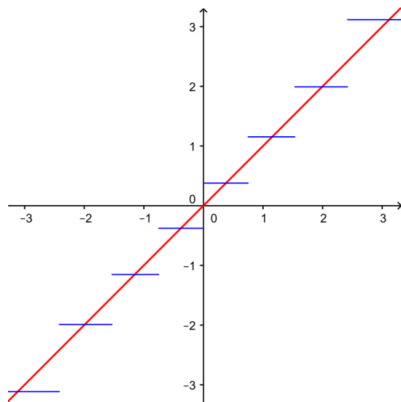
$$G(x) = \begin{cases} + (2D(P^x \| P^\mu))^{1/2}, & \text{for } \lambda \geq \mu; \\ - (2D(P^x \| P^\mu))^{1/2}, & \text{for } \lambda < \mu. \end{cases}$$

If $P^\lambda = N(\lambda, \sigma^2)$ then $G(x) = \frac{x-\mu}{\sigma}$.

For any exponential family $G(x) = \frac{x-\mu}{\sigma}$ is the first part of the Taylor expansion of G around $x = \mu$.

QQ-plot for binomial

Assume that $X \sim \text{bin}(n, p)$. For each $q \in (0, 1)$ plot the q -quantile of a standard Gaussian against the q -quantile of $G(X)$.



$$\Pr(X < j) \leq \Pr(Z \leq G(j)) \leq \Pr(X \leq j).$$

[Serov and Zubkov, 2013]

Where do they intersect?

The intersection point is approximately given by the following result. If $X \sim \text{bin}(n, p)$ then if nq is an integer we have

$$\Pr(X \leq nq) = \Phi(G(j + c_q)) \cdot \left(1 + O\left(\frac{1}{n}\right)\right)$$

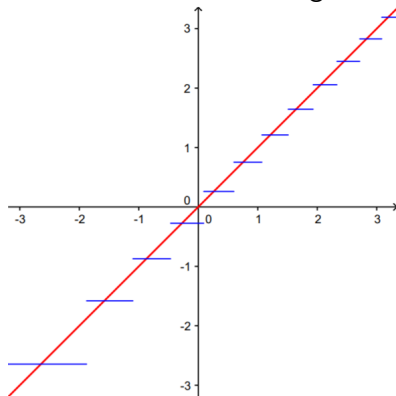
where

$$c_q = \frac{1}{2} + \frac{\ln\left(\frac{2D(q\|p)}{(q-p)^2} p(1-p)\right)}{2 \ln\left(\frac{q(1-p)}{p(1-q)}\right)}.$$

[PH, L. Györfi and G. Tusnády, 2012]

QQ-plot for Poisson

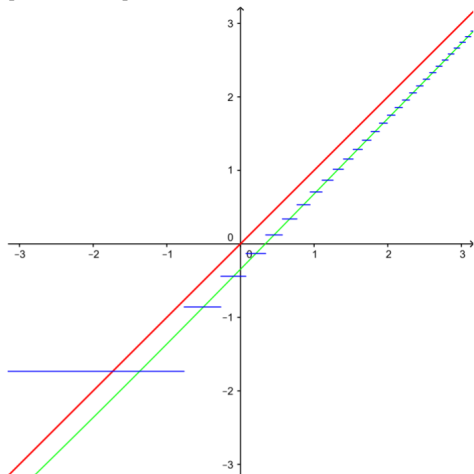
Assume that $X \sim \text{Po}(\lambda)$. For each $q \in (0, 1)$ plot the q -quantile of a standard Gaussian against the q -quantile of $G(X)$.



[PH and Tusnády, 2011]

QQ-plot for negative binomial

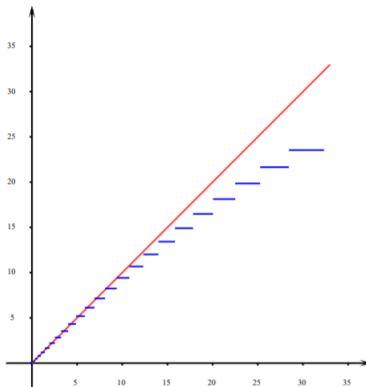
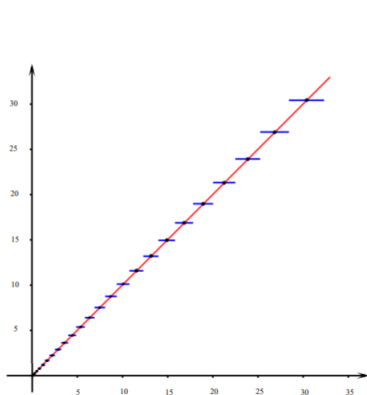
Assume that $X \sim \text{negbin}(k, p)$. For each $q \in (0, 1)$ plot the q -quantile of a standard Gaussian against the q -quantile of $G(X)$.
[PH 2016]



- 1 Prove majorization for Gamma distributions.
- 2 Prove intersection for negative binomial and Gamma distributions.
- 3 Combine to get upper bound for binomial.
- 4 Use upper bound on the binomial variable $n - X$ to get a lower bound for X .

Application

Information divergence is more χ^2 -distributed than the χ^2 -statistic.



- If you expand too little you will get punished by a factor of 2.
- Lower bounds can be found using orthogonal polynomials.
- Saddlepoint approximations can often be replaced by powerful inequalities.
- Use information divergence rather than total variation or χ^2 -divergence.

Work in progress:

- Simplify upper bounds.
- Bounds on moment generating functions.
- Generalizations to multivariate distributions.