

Introduction to Bayesian Methods - 6

Edoardo Milotti

Università di Trieste and INFN-Sezione di Trieste

2. The Li&Ma method (ctd.)

- **MaxL ratio**

$$\lambda_{\max} = \frac{L(D|H_0)|_{\max}}{L(D|H_1)|_{\max}} = \left(\frac{\alpha}{\alpha + 1} \frac{N_{\text{on}} + N_{\text{off}}}{N_{\text{on}}} \right)^{N_{\text{on}}} \left(\frac{1}{\alpha + 1} \frac{N_{\text{on}} + N_{\text{off}}}{N_{\text{off}}} \right)^{N_{\text{off}}}$$

therefore the significance can be obtained from $-2 \ln \lambda_{\max}$ because $-2 \ln \lambda$ has a chi-square distribution with 1 degree of freedom (only one parameter – the background rate – matters in the case of null hypothesis, while the alternative hypothesis has two parameters – background rate and source rate).

- if $x^2 \sim \chi^2(1)$ then $|x| \sim \chi(1)$, and we estimate the significance as

$$S \approx \sqrt{-2 \ln \lambda_{\max}} = \sqrt{2} \left\{ N_{\text{on}} \ln \left[\frac{\alpha + 1}{\alpha} \left(\frac{N_{\text{on}}}{N_{\text{on}} + N_{\text{off}}} \right) \right] + N_{\text{off}} \ln \left[(\alpha + 1) \left(\frac{N_{\text{off}}}{N_{\text{on}} + N_{\text{off}}} \right) \right] \right\}$$

(a perfect match with exp. data gives a vanishing chi, the actual value of chi is an estimate of the size of the fluctuation in terms of standard deviations).

- to be continued ...

Bayesian approach

(see M.L. Knoetig, ApJ 790:106, 2014)

- broader definition of α (Berge & al., A&A 466, 1219–1229 (2007))

$$\alpha = \frac{A_{\text{on}} t_{\text{on}}}{A_{\text{off}} t_{\text{off}}} = \frac{\int_{\text{on}} A_{\text{on}}^{\gamma}(\psi_x, \psi_y, \phi_z, E, t) d\psi_x d\psi_y d\phi_z dE dt}{\int_{\text{off}} A_{\text{off}}^{\gamma}(\psi_x, \psi_y, \phi_z, E, t) d\psi_x d\psi_y d\phi_z dE dt}$$

Acceptance FOV coords. zenithal angle Energy time

- comparison between competing hypotheses:
 - H_0 : the observed counts are due to background only
 - H_1 : a signal process contributes to the counts

$$P(H_i | N_{\text{on}}, N_{\text{off}}) = \frac{P(N_{\text{on}}, N_{\text{off}} | H_i) P_0(H_i)}{P(N_{\text{on}}, N_{\text{off}})}$$

- **evidence**

$$P(N_{\text{on}}, N_{\text{off}}) = \sum_i P(N_{\text{on}}, N_{\text{off}} | H_i) P_0(H_i)$$

- since the likelihood is determined by Poisson probabilities defined by continuous parameters, this must be extended to the following expression

$$P(N_{\text{on}}, N_{\text{off}}) = \sum_i \int_{\Lambda_i} P(N_{\text{on}}, N_{\text{off}} | \boldsymbol{\lambda}_i, H_i) P_0(\boldsymbol{\lambda}_i | H_i) d\boldsymbol{\lambda}_i P_0(H_i)$$

where the vectors of mean counts can be further specified in terms of signal count and background count

$$\begin{aligned} P(N_{\text{on}}, N_{\text{off}}) &= \int_{\Lambda_0} P(N_{\text{on}}, N_{\text{off}} | \lambda_{\text{bkg}}, H_0) P_0(\lambda_{\text{bkg}} | H_0) d\lambda_{\text{bkg}} P_0(H_0) \\ &+ \int_{\Lambda_1} P(N_{\text{on}}, N_{\text{off}} | \lambda_s, \lambda_{\text{bkg}}, H_1) P_0(\lambda_s, \lambda_{\text{bkg}} | H_1) d\lambda_s d\lambda_{\text{bkg}} P_0(H_1) \end{aligned}$$

- with these assumptions, the expected numbers of events are
 - in the OFF region (H_0 only)

$$\mathbf{E}(N_{\text{off}}) = \lambda_{\text{bkg}}$$

- in the ON region (H_0)

$$\mathbf{E}(N_{\text{on}}) = \alpha \lambda_{\text{bkg}}$$

- in the ON region (H_1)

$$\mathbf{E}(N_{\text{on}}) = \lambda_s + \alpha \lambda_{\text{bkg}}$$

- corresponding **likelihoods**

$$P(N_{\text{on}}, N_{\text{off}} | \lambda_{\text{bkg}}, H_0) = \frac{\lambda_{\text{bkg}}^{N_{\text{off}}}}{N_{\text{off}}!} e^{-\lambda_{\text{bkg}}} \times \frac{(\alpha \lambda_{\text{bkg}})^{N_{\text{on}}}}{N_{\text{on}}!} e^{-\alpha \lambda_{\text{bkg}}}$$

$$P(N_{\text{on}}, N_{\text{off}} | \lambda_s, \lambda_{\text{bkg}}, H_1) = \frac{\lambda_{\text{bkg}}^{N_{\text{off}}}}{N_{\text{off}}!} e^{-\lambda_{\text{bkg}}} \times \frac{(\lambda_s + \alpha \lambda_{\text{bkg}})^{N_{\text{on}}}}{N_{\text{on}}!} e^{-(\lambda_s + \alpha \lambda_{\text{bkg}})}$$

- **priors** from Jeffreys' rule

$$P_0(\boldsymbol{\lambda}_i | H_i) \propto \sqrt{\det [I(\boldsymbol{\lambda}_i | H_i)]}$$

$$I_{kl}(\boldsymbol{\lambda}_i | H_i) = -\mathbf{E} \left[\frac{\partial^2 \ln P(N_{\text{on}}, N_{\text{off}}) | \boldsymbol{\lambda}_i, H_i}{\partial \lambda_k \partial \lambda_l} \right]$$

- null hypothesis (expand the expression to prove the result)

$$P_0(\lambda_{\text{bkg}} | H_0) =$$

$$= \left(- \sum_{N_{\text{off}}, N_{\text{on}}=0}^{\infty} \left\{ \partial_{\lambda_{\text{bkg}}}^2 \ln \left[\frac{\lambda_{\text{bkg}}^{N_{\text{off}}}}{N_{\text{off}}!} e^{-\lambda_{\text{bkg}}} \times \frac{(\alpha \lambda_{\text{bkg}})^{N_{\text{on}}}}{N_{\text{on}}!} e^{-\alpha \lambda_{\text{bkg}}} \right] \right\} \frac{\lambda_{\text{bkg}}^{N_{\text{off}}}}{N_{\text{off}}!} e^{-\lambda_{\text{bkg}}} \times \frac{(\alpha \lambda_{\text{bkg}})^{N_{\text{on}}}}{N_{\text{on}}!} e^{-\alpha \lambda_{\text{bkg}}} \right)^{1/2}$$

$$= \sqrt{\frac{1 + \alpha}{\lambda_{\text{bkg}}}}$$

- alternative hypothesis (signal + background; results only, for details see the Knoetig 2014)

$$\begin{aligned}
 I_{s,s} &= \frac{1}{\lambda_s + \alpha \lambda_{\text{bkg}}} \\
 I_{s,\text{bkg}} &= I_{\text{bkg},s} = \frac{\alpha}{\lambda_s + \alpha \lambda_{\text{bkg}}} \\
 I_{\text{bkg},\text{bkg}} &= \frac{\lambda_s + \alpha \lambda_{\text{bkg}} + \alpha^2 \lambda_{\text{bkg}}}{\lambda_{\text{bkg}}(\lambda_s + \alpha \lambda_{\text{bkg}})}
 \end{aligned}
 \quad \Rightarrow \quad
 P_0(\lambda_s, \lambda_{\text{bkg}} | H_1) = \sqrt{\frac{1}{\lambda_{\text{bkg}}(\lambda_s + \alpha \lambda_{\text{bkg}})}}$$

- the Jeffreys' priors are improper, they are determined up to unknown proportionality constants and – taking equal prior probabilities $\frac{1}{2}$ for the two hypotheses – one finds

$$P(H_0 | N_{\text{on}}, N_{\text{off}}) = \frac{c_0 \gamma'}{c_0 \gamma' + c_1 \delta'}$$

with

$$\gamma' := \int_0^\infty P(N_{\text{on}}, N_{\text{off}} | \lambda_{\text{bg}}, H_0) P_0(\lambda_{\text{bg}} | H_0) d\lambda_{\text{bg}} P_0(H_0),$$

$$\delta' := \int_0^\infty \int_0^\infty P(N_{\text{on}}, N_{\text{off}} | \lambda_s, \lambda_{\text{bg}}, H_1) \times P_0(\lambda_s, \lambda_{\text{bg}} | H_1) d\lambda_s d\lambda_{\text{bg}} P_0(H_1).$$

- a possible choice for the c constants is to take them equal, Knoetig 2014 advocates a different choice

$$\frac{c_1}{c_0} = \frac{\gamma'}{\delta'} \Big|_{N_{\text{on}}, N_{\text{off}}=0}$$

and, finally,

$$\gamma := (1 + 2N_{\text{off}})\alpha^{1/2+N_{\text{on}}+N_{\text{off}}}\Gamma(1/2 + N_{\text{on}} + N_{\text{off}})$$

$$P(H_0|N_{\text{on}}, N_{\text{off}}) = \frac{\gamma}{\gamma + c_1/c_0\delta}$$

where

$$\delta := 2(1 + \alpha)^{N_{\text{on}}+N_{\text{off}}}\Gamma(1 + N_{\text{on}} + N_{\text{off}}) \times {}_2F_1(1/2 + N_{\text{off}}, 1 + N_{\text{on}} + N_{\text{off}}; 3/2 + N_{\text{off}}; -1/\alpha)$$

$$\frac{c_1}{c_0} = \frac{\sqrt{\pi}}{2 \arctan(1/\sqrt{\alpha})}$$

- since

$$P(H_1|N_{\text{on}}, N_{\text{off}}) = 1 - P(H_0|N_{\text{on}}, N_{\text{off}})$$

we determine the significance of the alternative hypothesis in terms of Gaussian standard deviations from

$$S_B = \sqrt{2} \operatorname{erf}^{-1} [1 - P(H_0|N_{\text{on}}, N_{\text{off}})]$$

- if the signal hypothesis holds, we can compute the posterior pdf for the mean signal counts from Bayes law

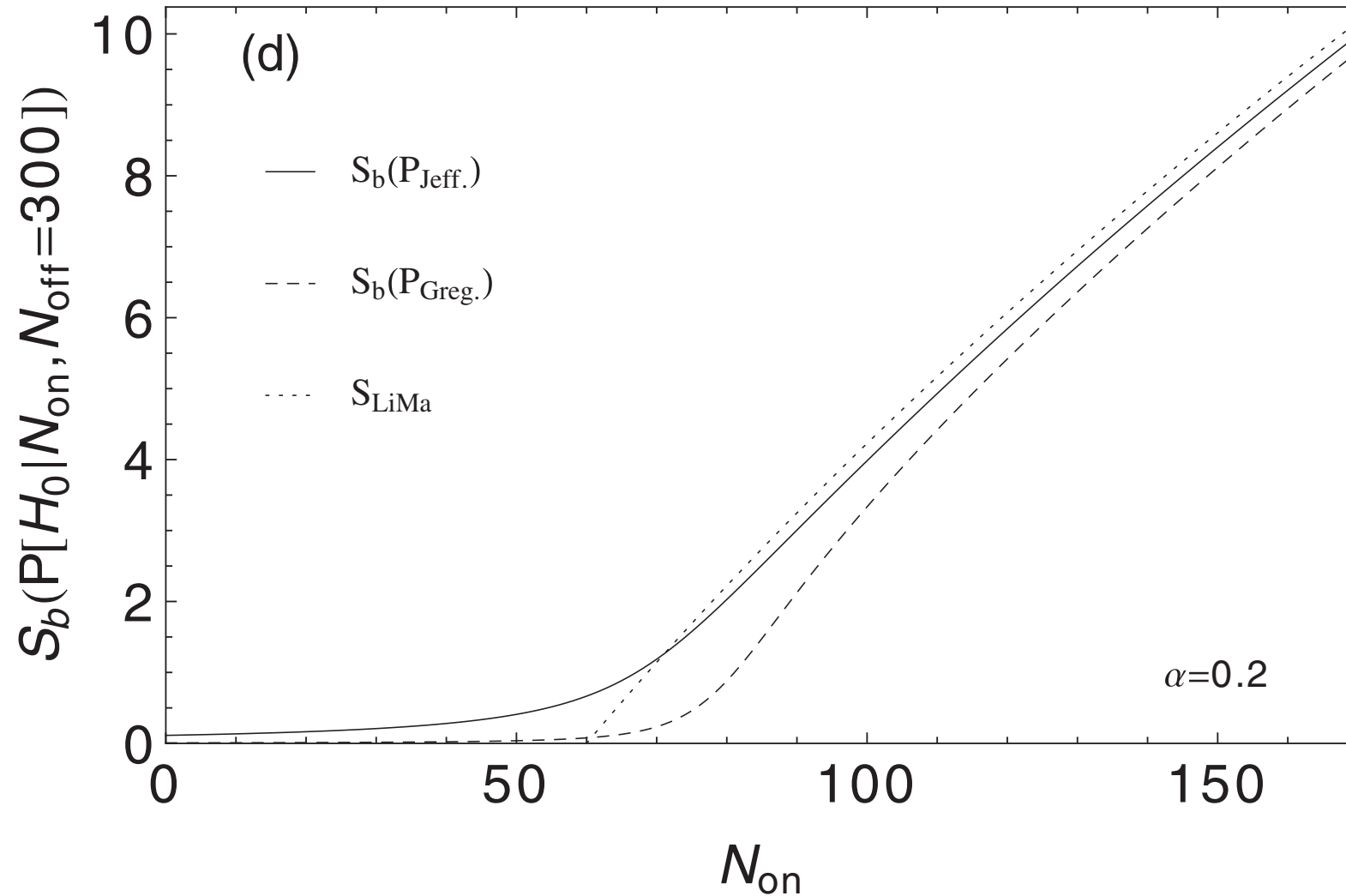
$$p(\lambda_s, \lambda_{\text{bkg}} | N_{\text{on}}, N_{\text{off}}, H_1) = \frac{p(N_{\text{on}}, N_{\text{off}} | \lambda_s, \lambda_{\text{bkg}}, H_1) p_0(\lambda_s, \lambda_{\text{bkg}} | H_1)}{\int_0^\infty \int_0^\infty p(N_{\text{on}}, N_{\text{off}} | \lambda'_s, \lambda'_{\text{bkg}}, H_1) p_0(\lambda'_s, \lambda'_{\text{bkg}} | H_1) d\lambda'_s d\lambda'_{\text{bkg}}}$$

and marginalizing with respect to the mean background counts

$$p(\lambda_s | N_{\text{on}}, N_{\text{off}}, H_1) = \int_0^\infty p(\lambda_s, \lambda'_{\text{bkg}} | N_{\text{on}}, N_{\text{off}}, H_1) d\lambda'_{\text{bkg}}$$

In Knoetig 2014, it is shown that the integral can be evaluated analytically in terms of the regularized hypergeometric function and of the Tricomi confluent hypergeometric function

$$P(\lambda_s | N_{\text{on}}, N_{\text{off}}, H_1) = P_{\text{P}}(N_{\text{on}} + N_{\text{off}} | \lambda_s) \times \frac{U[1/2 + N_{\text{off}}, 1 + N_{\text{off}} + N_{\text{on}}, (1 + 1/\alpha)\lambda_s]}{{}_2\tilde{F}_1(1/2 + N_{\text{off}}, 1 + N_{\text{off}} + N_{\text{on}}; 3/2 + N_{\text{off}}; -1/\alpha)}$$



Comparison of significance obtained in the large counts regime between two priors and with Li & Ma (from Knoetig 2024).

3. Model selection

The generic purpose of a model selection statistic is to set up a tension between the predictiveness of a model (for instance indicated by the number of free parameters) and its ability to fit observational data. Oversimplistic models offering a poor fit should of course be thrown out, but so should more complex models that offer poor predictive power.

There are two main types of model selection statistic that have been used in the literature so far. Information criteria look at the best-fitting parameter values and attach a penalty for the number of parameters; they are essentially a technical formulation of "chi-squared per degrees of freedom" arguments. By contrast, the Bayesian evidence applies the same type of likelihood analysis familiar from parameter estimation, but at the level of models rather than parameters. It depends on goodness of fit across the entire model parameter space.

(Liddle & al., 2006 – Astronomy & Geophysics, Volume 47, Issue 4, pp. 4.30-4.33)

Information criteria for astrophysical model selection

Andrew R. Liddle¹ ²★

¹*Astronomy Centre, University of Sussex, Brighton BN1 9QH*

²*Institute for Astronomy, University of Hawai'i, 2680 Woodlawn Drive, Honolulu, Hawai'i 96822, USA*

Accepted 2007 February 19. Received 2007 February 16; in original form 2007 January 8

ABSTRACT

Model selection is the problem of distinguishing competing models, perhaps featuring different numbers of parameters. The statistics literature contains two distinct sets of tools, those based on information theory such as the Akaike Information Criterion (AIC), and those on Bayesian inference such as the Bayesian evidence and Bayesian Information Criterion (BIC). The Deviance Information Criterion combines ideas from both heritages; it is readily computed from Monte Carlo posterior samples and, unlike the AIC and BIC, allows for parameter degeneracy. I describe the properties of the information criteria, and as an example compute them from *Wilkinson Microwave Anisotropy Probe* 3-yr data for several cosmological models. I find that at present the information theory and Bayesian approaches give significantly different conclusions from that data.

Akaike Information Criterion (AIC).

This was derived by Hirotugu Akaike in 1974, and takes the form

$$\text{AIC} = -2 \ln \mathcal{L}_{\max} + 2k$$

where k is the number of parameters in the model. The subscript “max” indicates that one should find the parameter values yielding the highest possible likelihood within the model. This second term acts as a kind of “Occam factor”; initially, as parameters are added, the fit to data improves rapidly until a reasonable fit is achieved, but further parameters then add little and the penalty term $2k$ takes over. The generic shape of the AIC as a function of number of parameters is a rapid fall, a minimum, and then a rise. The preferred model sits at the minimum.

The AIC was derived from information-theoretic considerations, specifically an approximate minimization of the Kullback–Leibler information entropy which measures the distance between two probability distributions.

(Liddle & al., 2006)

Outline of Akaike's derivation

1. max log-likelihood ratio between conjectured model (k -dimensional parameter vector) and true model (L -dimensional parameter vector)

$$\ln \frac{f(x|\hat{\theta}^{(k)})}{f(x|\theta)}$$

2. this depends on the dataset, which is distributed according to the true model; in order to get rid of the fluctuations, we average the max log-likelihood over the true distribution

$$\mathbf{E} \left[\ln \frac{f(x|\hat{\theta}^{(k)})}{f(x|\theta)} \right] = \int_{\Theta} f(x|\theta) \ln \frac{f(x|\hat{\theta}^{(k)})}{f(x|\theta)} dx$$

3. here we remark that:

- this is purely theoretical, since we do not know the true pdf
- the r.h.s. expression is the negative of the Kullback-Leibler divergence between the conjectured and the true pdf
- the r.h.s. expression can be written as

$$\int_{\Theta} f(x|\theta) \ln \frac{f(x|\hat{\theta}^{(k)})}{f(x|\theta)} dx = \int_{\Theta} f(x|\theta) \ln f(x|\hat{\theta}^{(k)}) - \int_{\Theta} f(x|\theta) \ln f(x|\theta)$$

Outline of Akaike's derivation

- the second term in the expansion is unknown, but it is a constant and we can get rid of it, and change sign as well (with an additional factor 2, see later), so that by minimizing the first term we actually minimize the KL divergence

$$\int_{\Theta} f(x|\theta) \ln \frac{f(x|\hat{\theta}^{(k)})}{f(x|\theta)} dx = \int_{\Theta} f(x|\theta) \ln f(x|\hat{\theta}^{(k)}) - \int_{\Theta} f(x|\theta) \ln f(x|\theta) \rightarrow -2 \int_{\Theta} f(x|\theta) \ln f(x|\hat{\theta}^{(k)})$$

- going back to Wilks' theorem, we know that the remaining $L-k$ degrees of freedom in the likelihood ratio are (asymptotically) normally distributed, therefore the $-2\log$ has a chi-square distribution with $L-k$ degrees of freedom, with mean value $L-k$, and therefore the required mean value has an asymptotic bias $2(L-k)$; using the max likelihood as an estimator of the mean, we find that the discrepancy expressed by the equation above can be written as

$$-2 \ln f(x|\hat{\theta}^{(k)}) + 2k$$

after dropping the constant L

Bayesian Information Criterion (BIC).

This was derived by Gideon Schwarz in 1978, and strongly resembles the AIC. It is given by

$$\text{BIC} = -2 \ln \mathcal{L}_{\max} + k \ln N$$

where N is the number of datapoints. Since a typical dataset will have $\ln N > 2$, the BIC imposes a stricter penalty against extra parameters than the AIC.

It was derived as an approximation to the Bayesian evidence, to be discussed next, but the assumptions required are very restrictive and unlikely to hold in practice, rendering the approximation quite crude.

(Liddle & al., 2006)

Bayesian evidence

Model selection aims to determine which theoretical models are most plausible given some data, without necessarily considering preferred values of model parameters.

Ideally, we would like to estimate posterior probabilities on the set of all competing models using Bayes' theorem:

$$P(M_i|D, I) = \frac{P(D|M_i, I)P(M_i|I)}{\sum_k P(D|M_k, I)P(M_k|I)}$$

and select the best model using the odds ratio

$$\mathcal{O}_{i,j} = \frac{P(M_i|D, I)}{P(M_j|D, I)} = \frac{P(D|M_i, I)P(M_i|I)}{P(D|M_j, I)P(M_j|I)}$$

or the Bayes factor, if we assume equal prior probabilities for the different models:

$$B_{i,j} = \frac{P(D|M_i, I)}{P(D|M_j, I)}$$

Thus, we see that the Bayes factor is a ratio of evidences

$$B_{i,j} = \frac{P(D|M_i, I)}{P(D|M_j, I)}$$

As usual, each evidence is obtained by marginalizing the likelihood with respect to the (potentially different) parameters:

$$P(D|M_i, I) = \int_{\Theta_i} P(D|\boldsymbol{\theta}_i, M_i, I)p(\boldsymbol{\theta}_i|M_i, I)d\boldsymbol{\theta}_i$$

The evidence of a model is thus the average likelihood of the model in the prior.

Unlike the AIC and BIC, it does not focus on the best-fitting parameters of the model but asks “of all the parameter values you thought were viable before the data came along, how well on average did they fit the data?”. Literally, it is the likelihood of the model given the data.

The evidence rewards predictability of models, provided they give a good fit to the data, and hence gives an axiomatic realization of Occam's razor.

A model with little parameter freedom is likely to fit data over much of its parameter space, whereas a model that could match pretty much any data that might have cropped up will give a better fit to the actual data but only in a small region of its larger parameter space, pulling the average likelihood down.

(Liddle & al., 2006)

Which statistics?

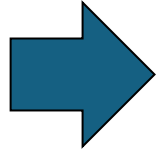
Of these statistics, we would advocate using – wherever possible – the Bayesian evidence, which is a full implementation of Bayesian inference and can be directly interpreted in terms of model probabilities. It is computationally challenging to compute, being a highly peaked multidimensional integral, but recent algorithm development has made it feasible in cosmological contexts.

*If the Bayesian evidence cannot be computed, the BIC can be deployed as a substitute. It is much simpler to compute as one need only find the point of maximum likelihood for each model. **However, interpreting it can be difficult. Its main usefulness is as an approximation to the evidence, but this holds only for gaussian likelihoods and provided the datapoints are independent and identically distributed.** The latter condition holds poorly for the current global cosmological dataset, though it can potentially be improved by binning of the data, hence decreasing the N in the penalty term.*

*The AIC has been widely used outside astrophysics but is of debatable utility. It has been shown to be “dimensionally inconsistent”, meaning that it is not guaranteed to give the right result even in the limit of infinite unbiased data. It may be useful for checking the robustness of conclusions drawn using the BIC. **The evidence and BIC are dimensionally consistent.***

(Liddle & al., 2006)

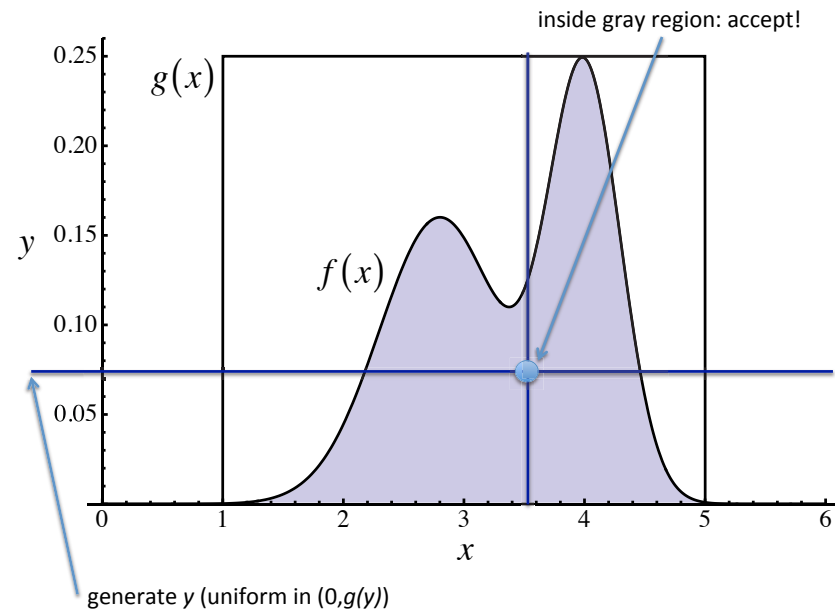
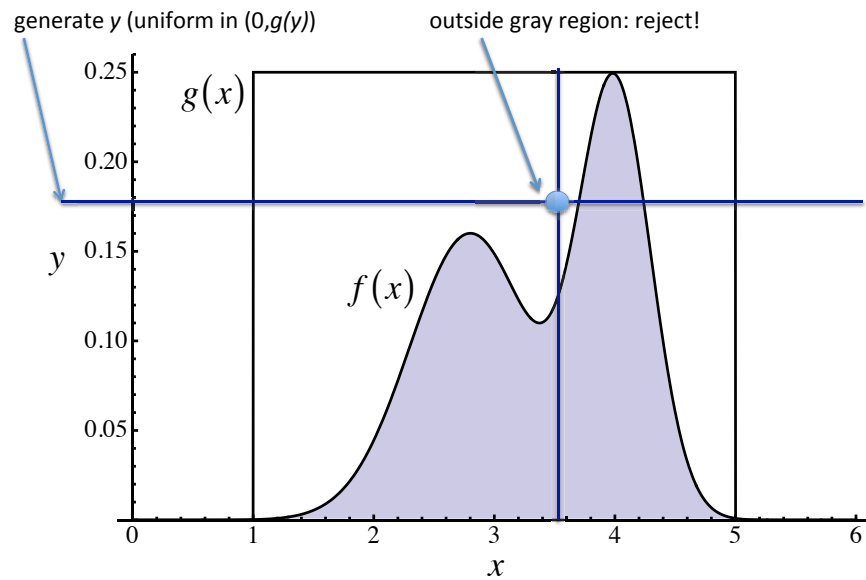
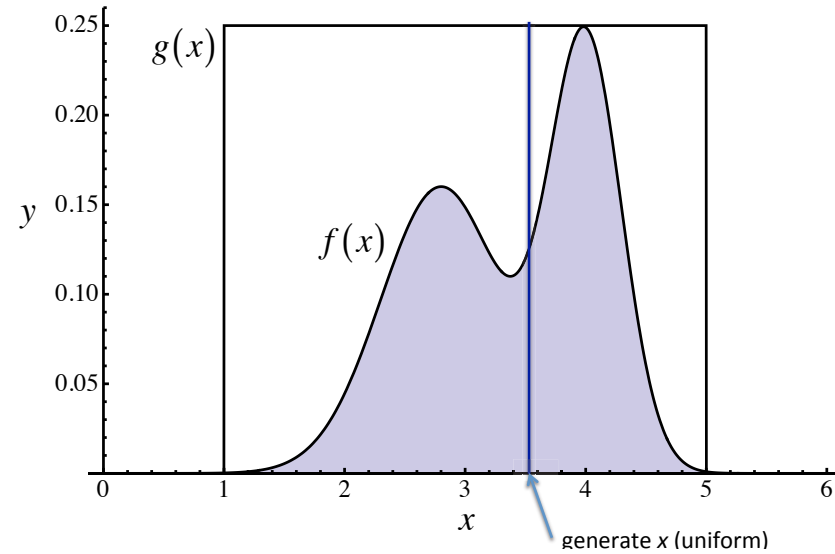
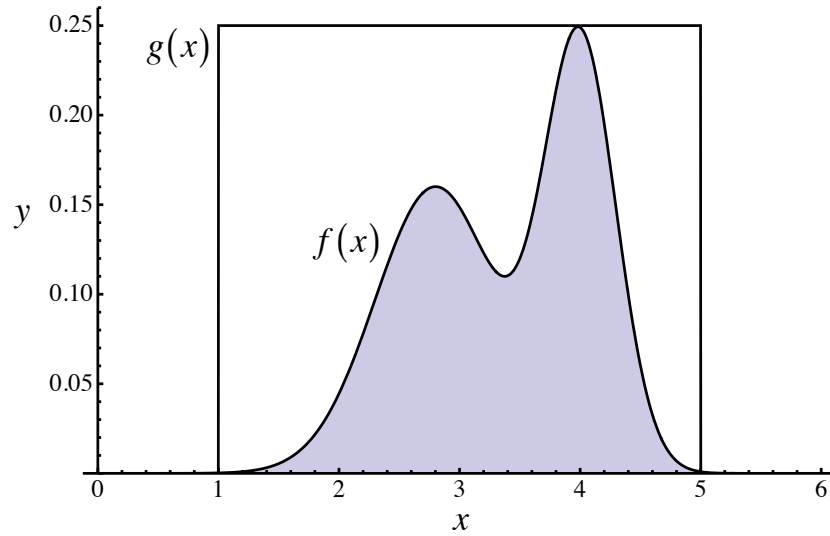
Our next important topic: Bayesian estimates often require complex numerical integrals. How do we confront this problem?



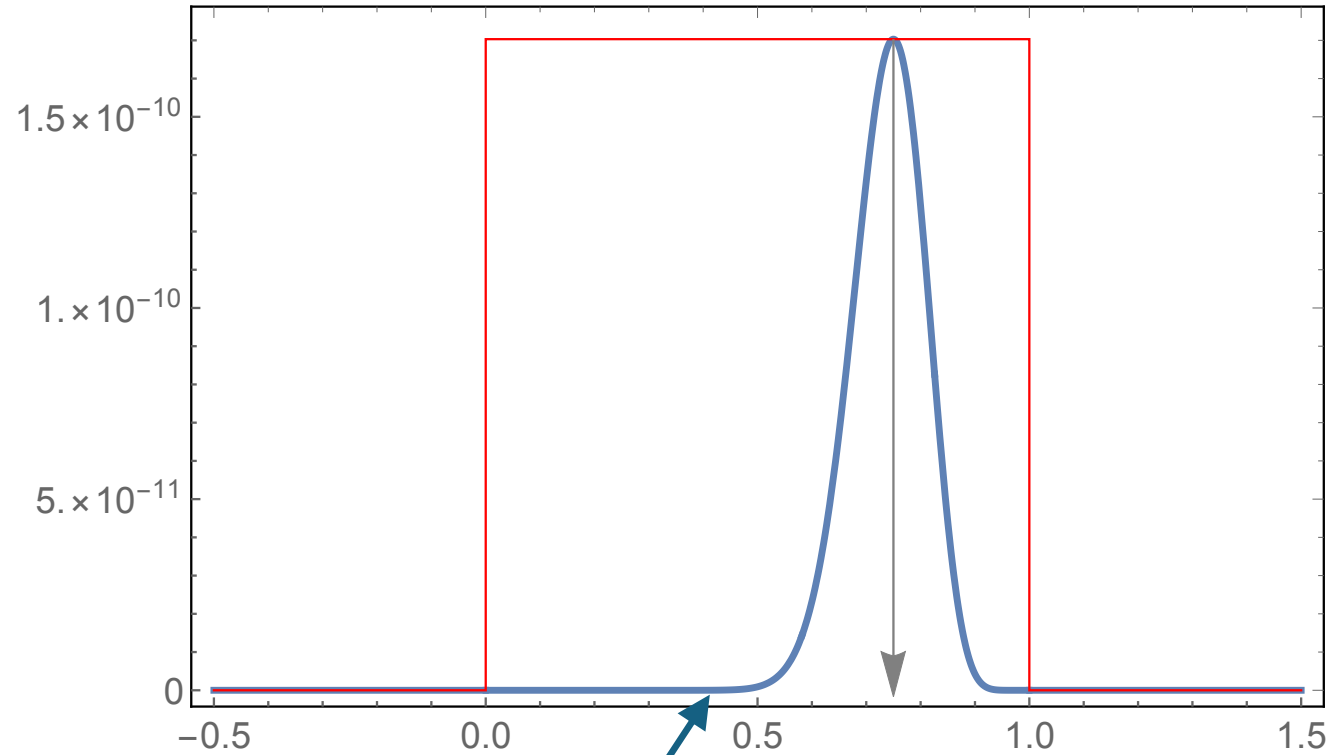
enter the Monte Carlo methods!

1. acceptance-rejection sampling
2. importance sampling
3. statistical bootstrap
4. Bayesian methods in a sampling-resampling perspective
5. Introduction to Markov chains and to Random Walks (RW)
6. Simulated annealing
7. The Metropolis algorithm
8. Markov Chain Monte Carlo (MCMC)
9. The Gibbs sampler
10. The efficiency of MCMC algorithms
11. Affine-invariant MCMC algorithms (EMCEE)

1. The acceptance rejection method



Example: generation of beta-distributed random numbers



$$p(x) = \frac{x^a(1-x)^b}{B(a+1, b+1)}$$

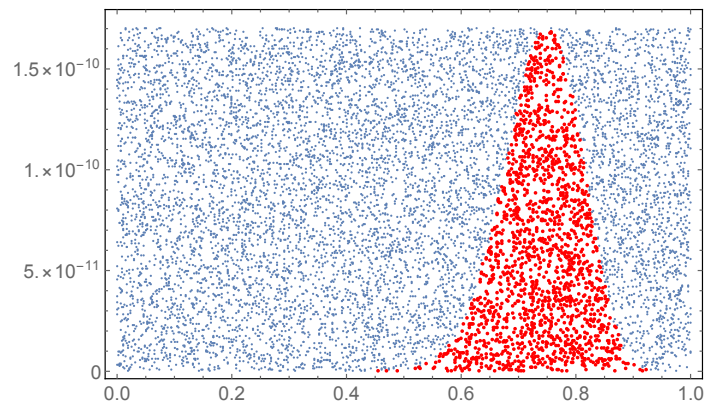
normalized distribution

$$p_0(x) = x^a(1-x)^b$$

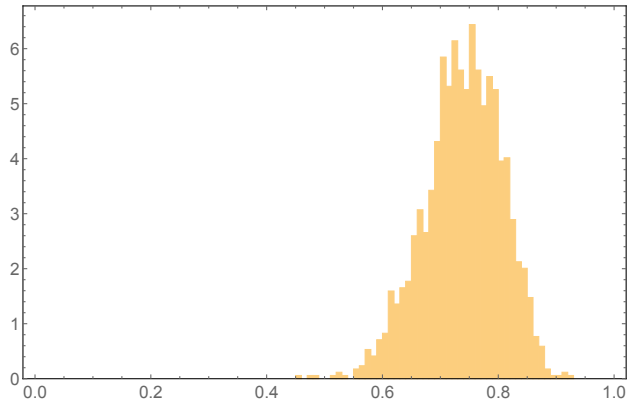
unnormalized distribution

$$x_{\max} = \frac{a}{a+b}$$

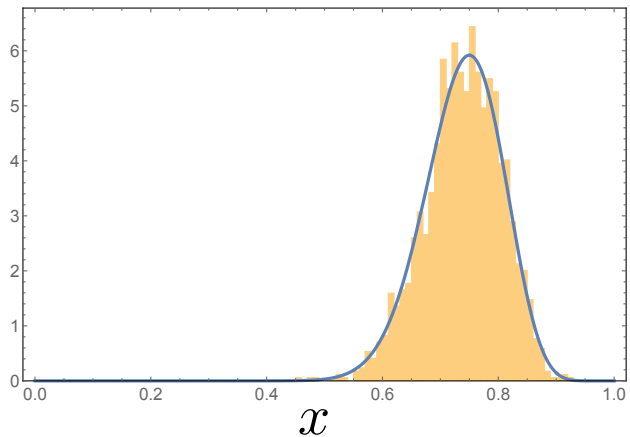
modal value



generated pairs
(red = accepted pairs)



normalized histogram of the
accepted x 's

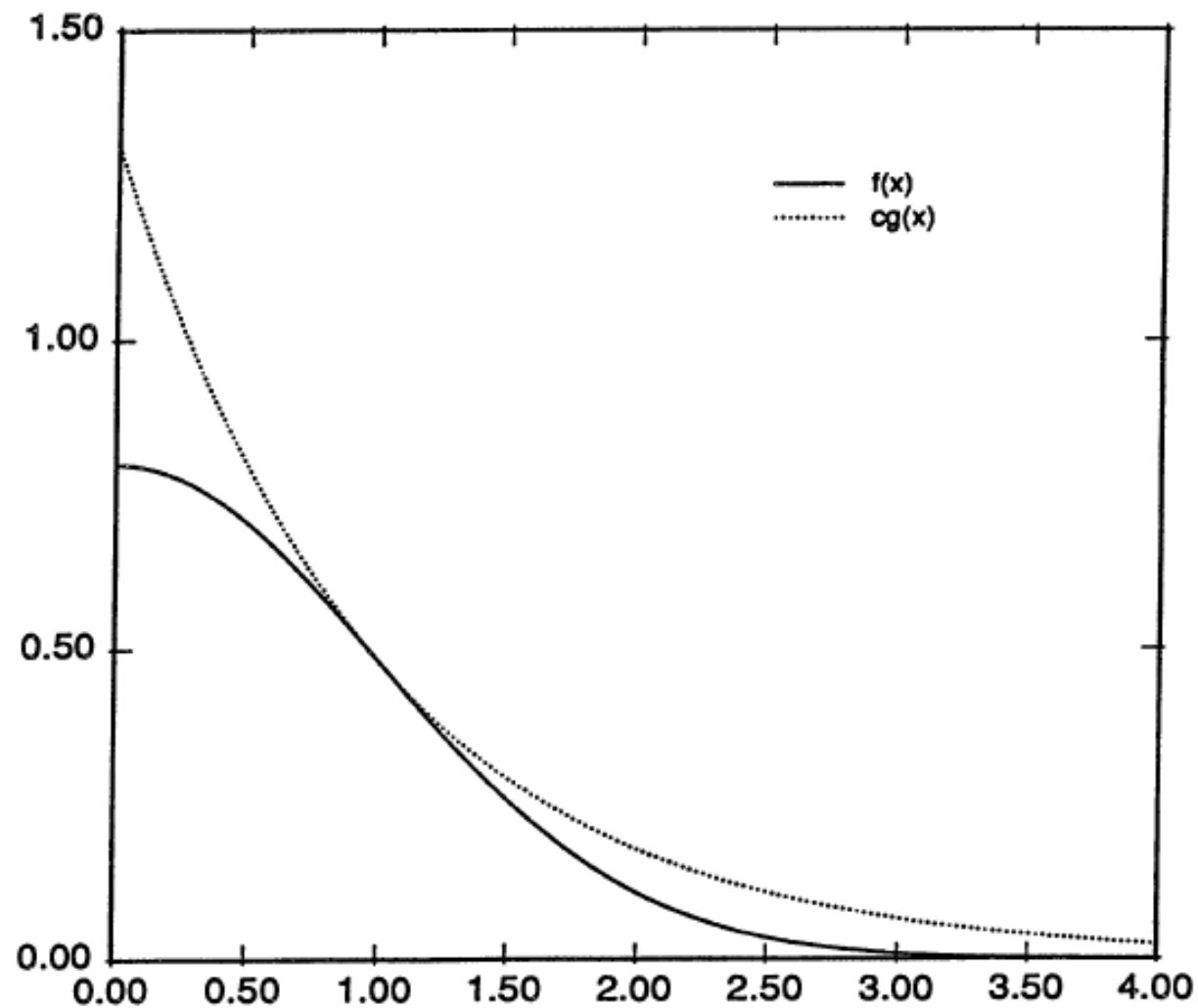


comparison with the plot of the
normalized beta distribution

Example: random numbers with semi-Gaussian distribution from exponentially distributed random numbers.

$$f(x) = \sqrt{\frac{2}{\pi}} \exp\left(-\frac{x^2}{2}\right) \quad x \geq 0$$

$$g(x) = \exp(-x)$$



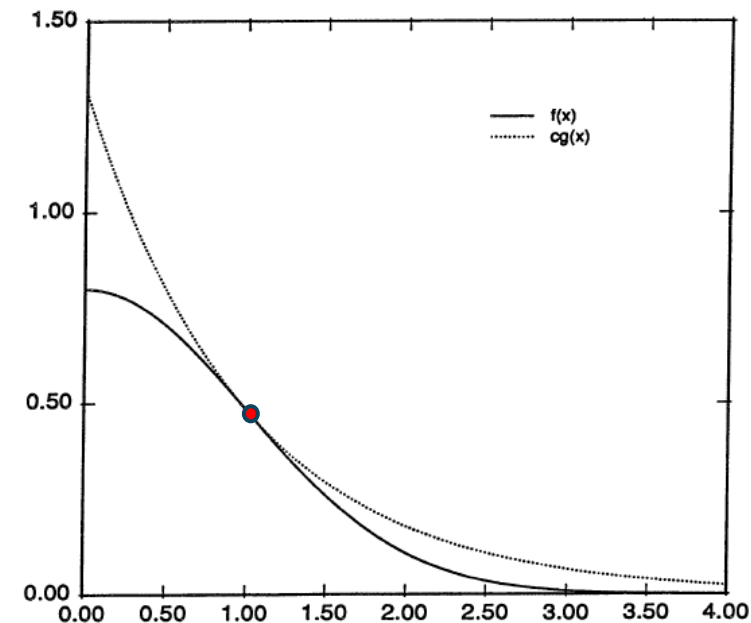
Definition of contact point (to maximize efficiency)

$$f(x) = \sqrt{\frac{2}{\pi}} \exp\left(-\frac{x^2}{2}\right) \quad x \geq 0$$

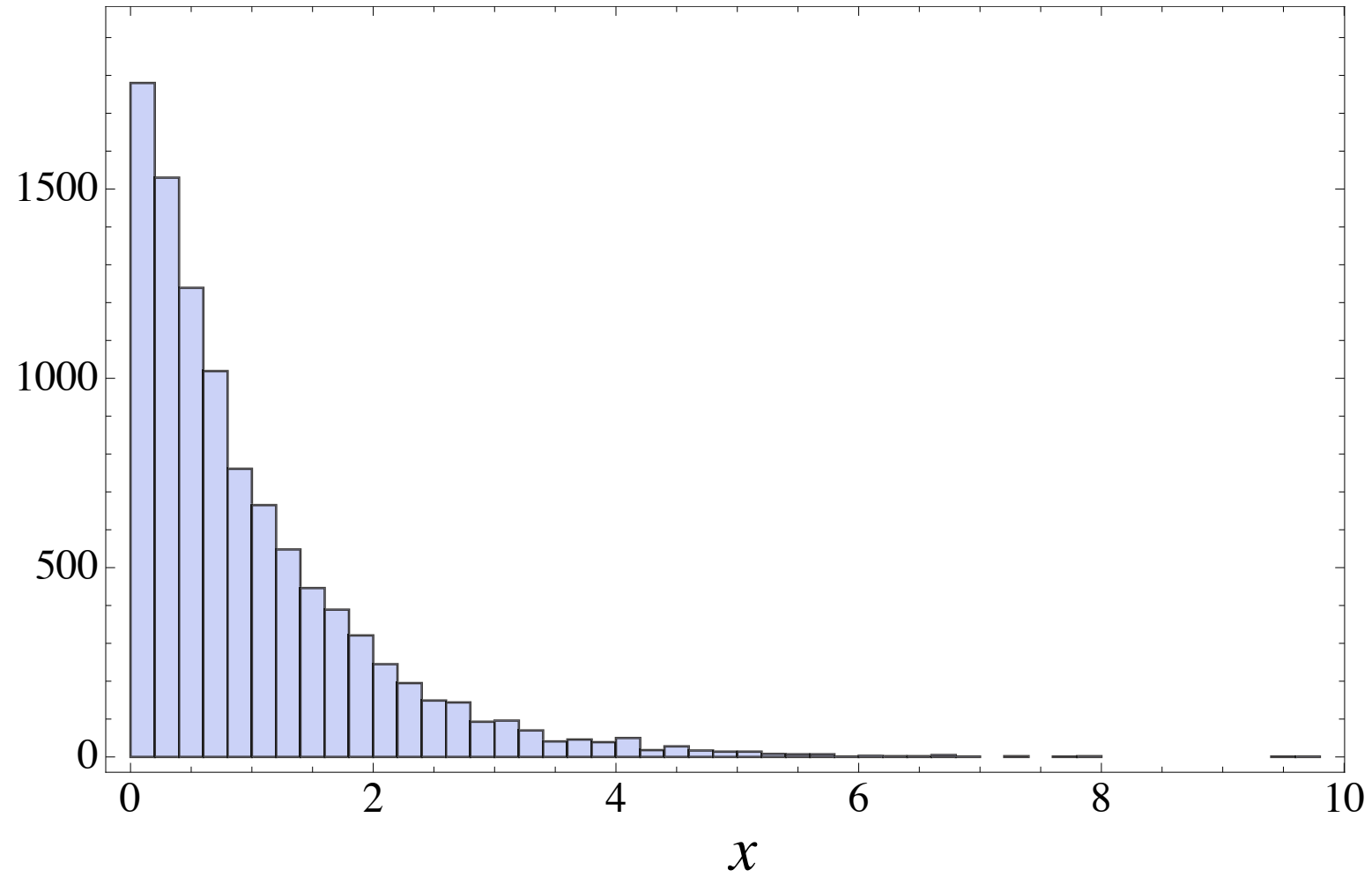
$$g(x) = \exp(-x)$$

$$\Rightarrow \begin{cases} f(x) = cg(x) \\ f'(x) = cg'(x) \end{cases} \Rightarrow \begin{cases} \sqrt{\frac{2}{\pi}} \exp\left(-\frac{x^2}{2}\right) = c \exp(-x) \\ x \sqrt{\frac{2}{\pi}} \exp\left(-\frac{x^2}{2}\right) = c \exp(-x) \end{cases}$$

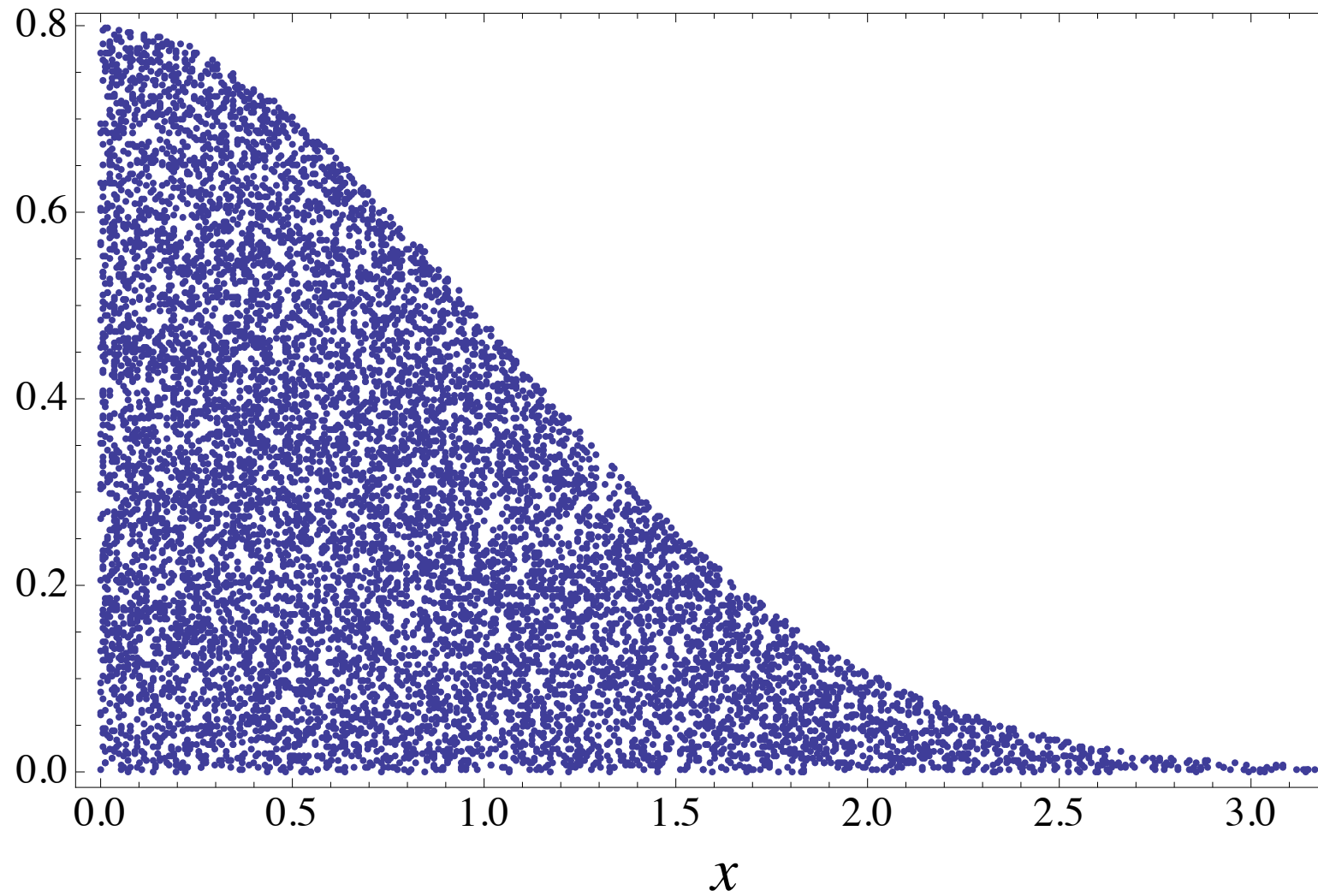
$$\Rightarrow x = 1; \quad c = \sqrt{\frac{2}{\pi}} \exp\left(-\frac{x^2}{2} + x\right) \approx 1.31549$$



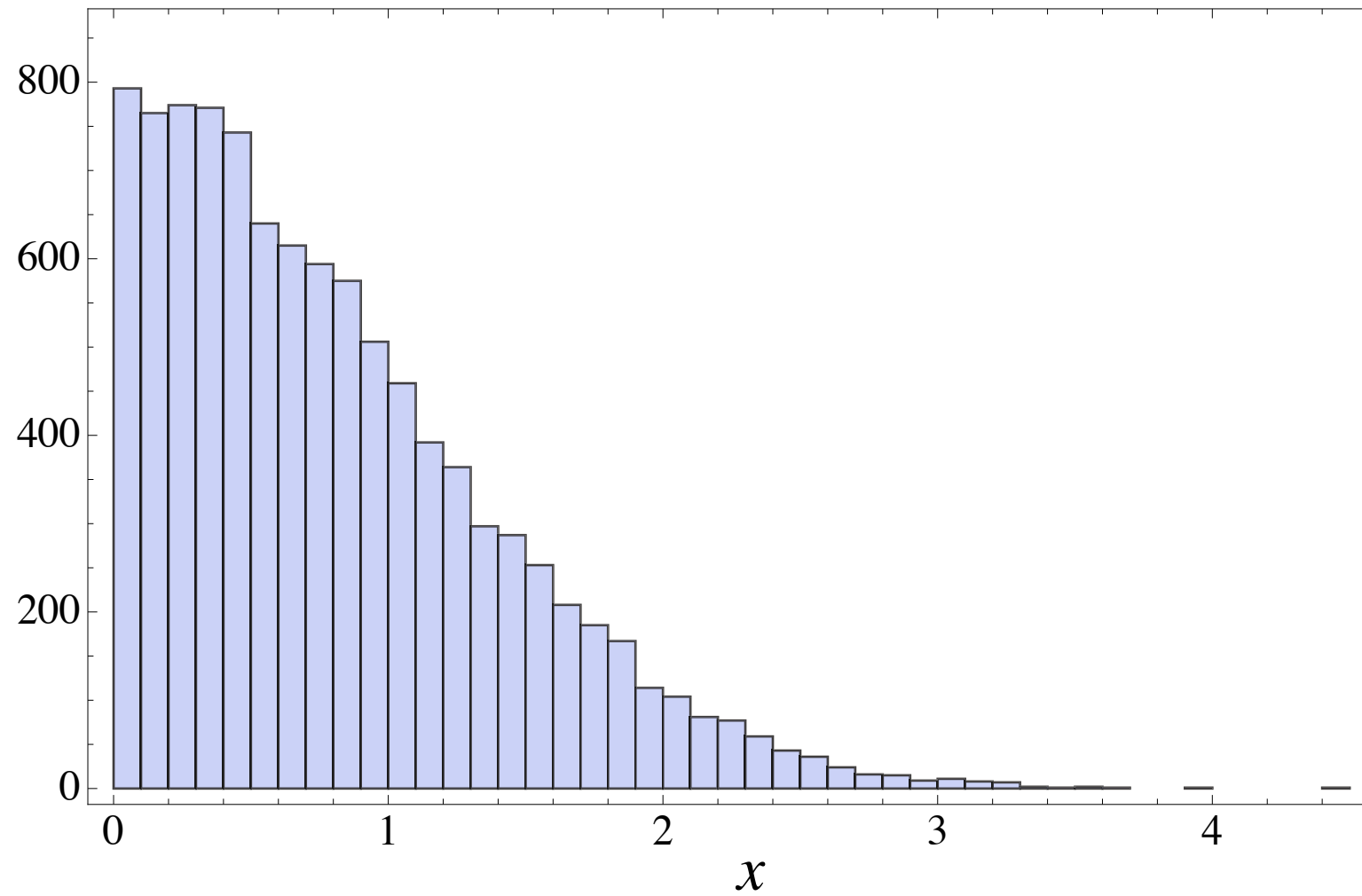
Exponentially distributed values



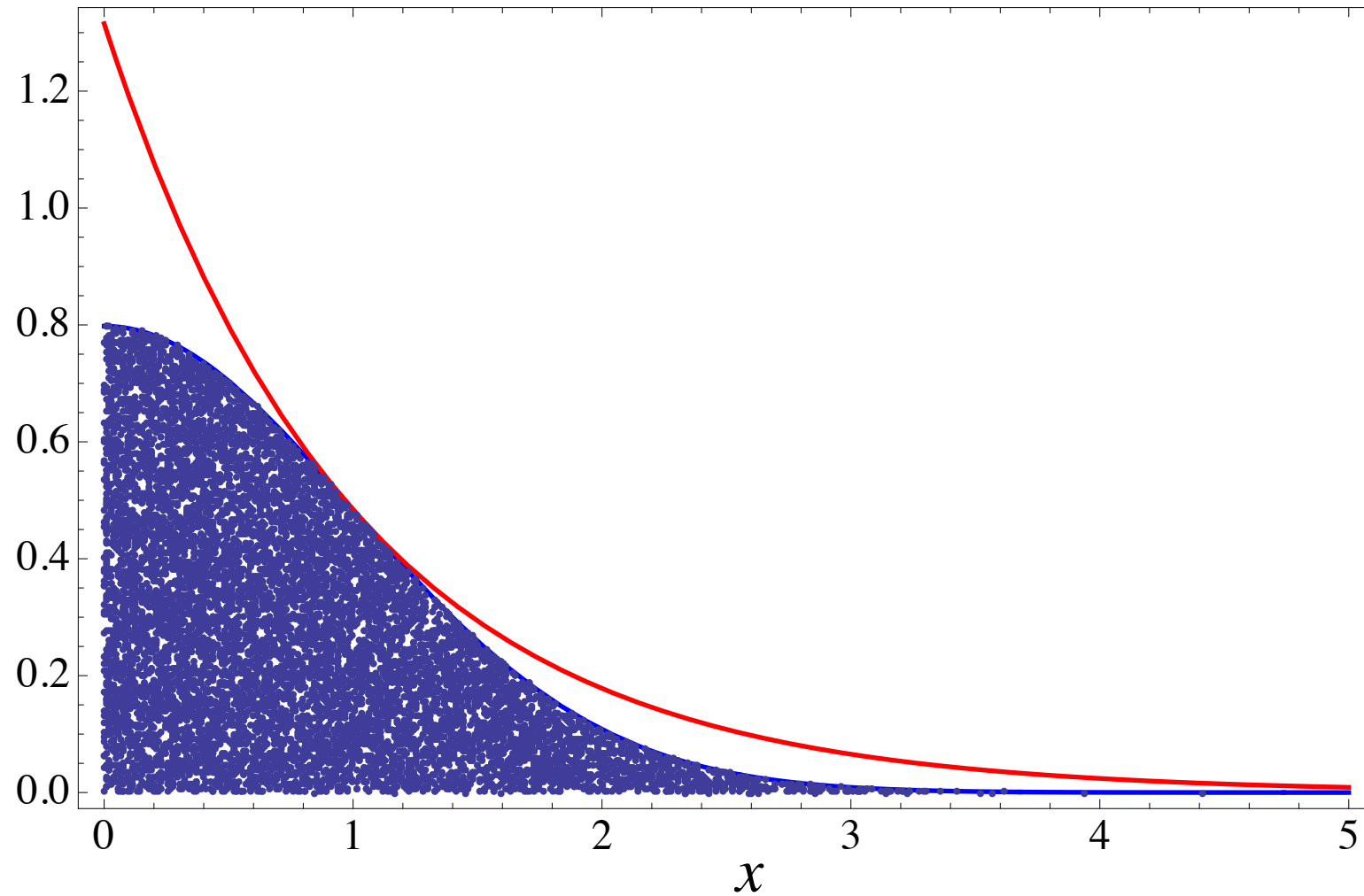
A/R accepted values (10000 accepted sample pairs)

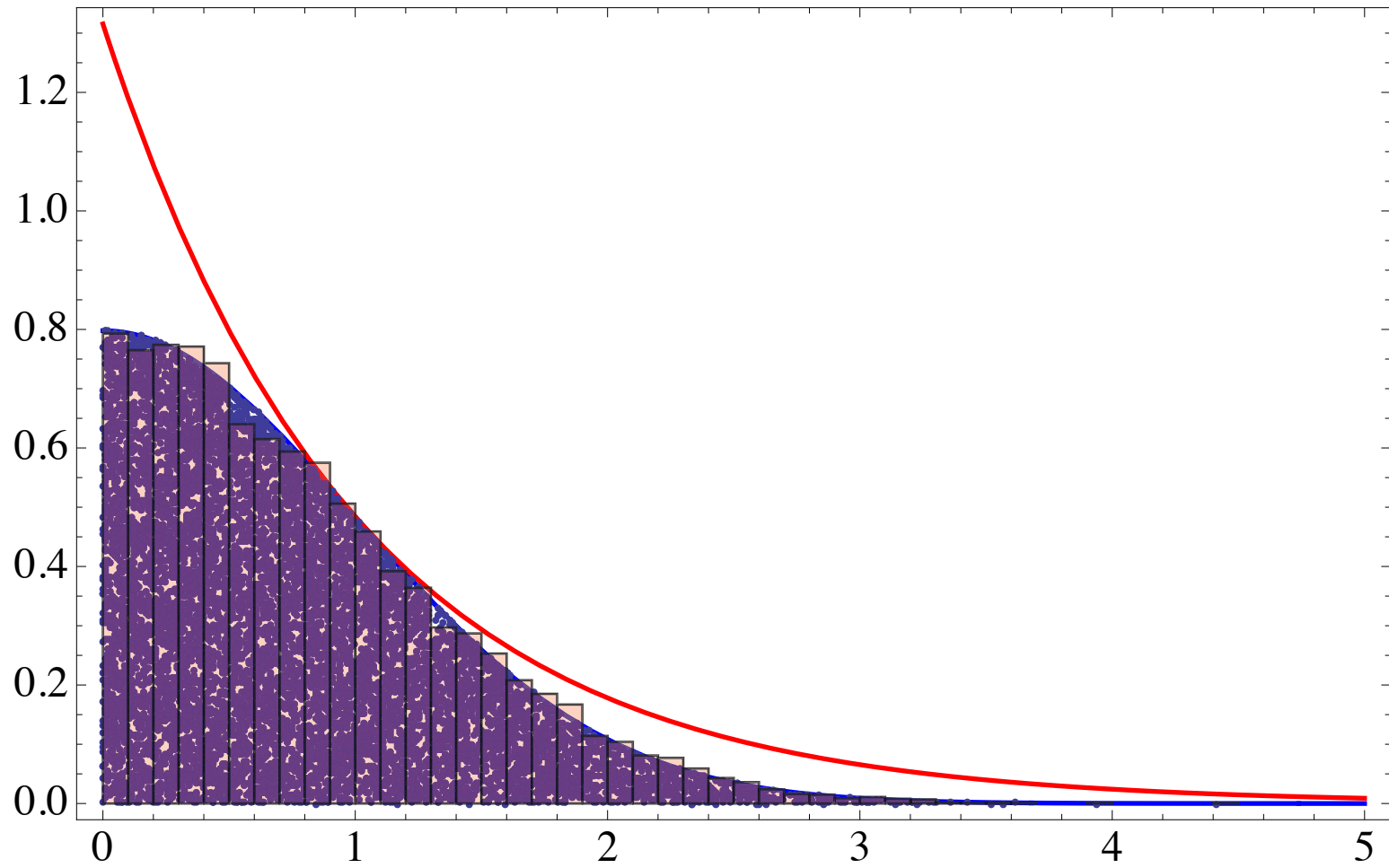


Histogram of accepted x values



Comparison with the original distributions





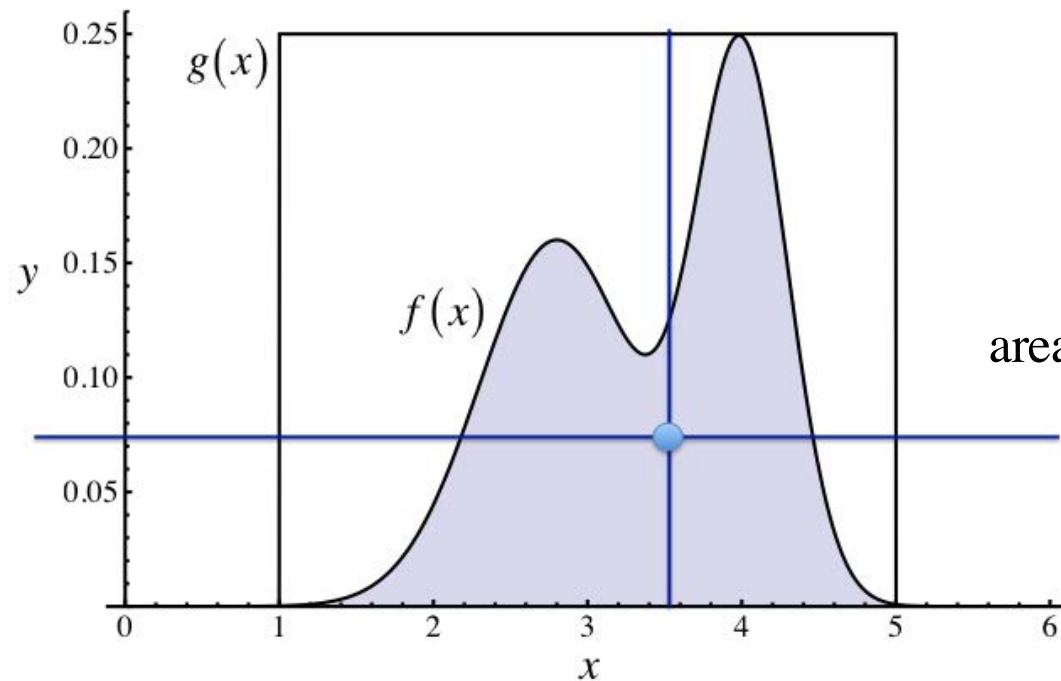
Short summary:

1. we create a data set by randomly sampling from the exponential distribution
2. we use the acceptance-rejection algorithm to resample the data set with the target distribution (the half-Gaussian)

This is a sampling – resampling technique (see later ...)

Notice that in this method we generate pairs of real numbers that are uniformly distributed between $f(x)$ and the x -axis, therefore we can use these pairs to estimate the total area under the curve

(here the reference area is the area of the enclosing rectangle which corresponds to a uniform distribution)



$$\text{area} = \frac{\# \text{ of accepted pairs}}{\# \text{ of pairs}} \text{reference area}$$

In general, if $h(x) = f(x)p(x)$, where p is a pdf

$$\int_a^b h(x) dx = \int_a^b f(x)p(x) dx = E_p[f(x)] \approx \frac{1}{N} \sum_{n=1}^N f(x_n)$$

here the x are i.i.d with pdf $p(x)$

and we find that the variance of this estimate of the integral is

$$\frac{1}{N} \left\{ \frac{1}{N-1} \sum_{n=1}^N [f(x_n) - E_p[f(x)]]^2 \right\}$$

We encounter a problem with this method when we must sample functions that have many narrow peaks.

2. Importance sampling

this pdf is troublesome ...

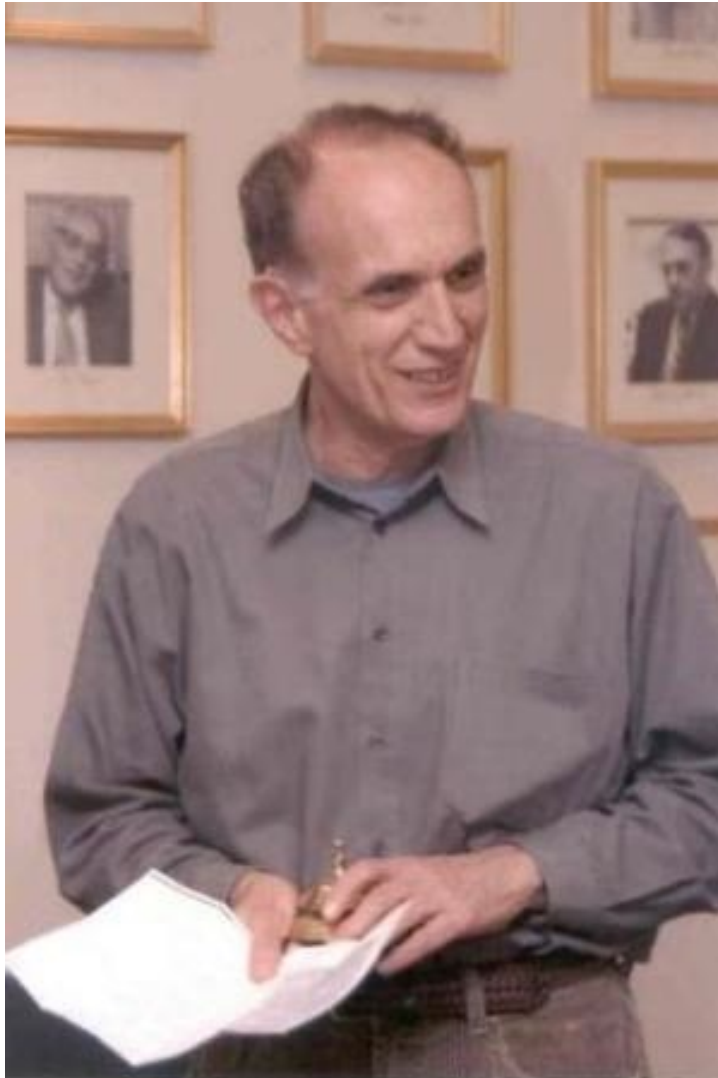
therefore, we use this ...

$$\int_a^b h(x) dx = \int_a^b f(x) p(x) dx = \int_a^b \left[f(x) \frac{p(x)}{q(x)} \right] q(x) dx$$
$$= E_q \left[f(x) \frac{p(x)}{q(x)} \right] \approx \frac{1}{N} \sum_{n=1}^N f(x_n) \frac{p(x_n)}{q(x_n)}$$

here the x are i.i.d with pdf $q(x)$

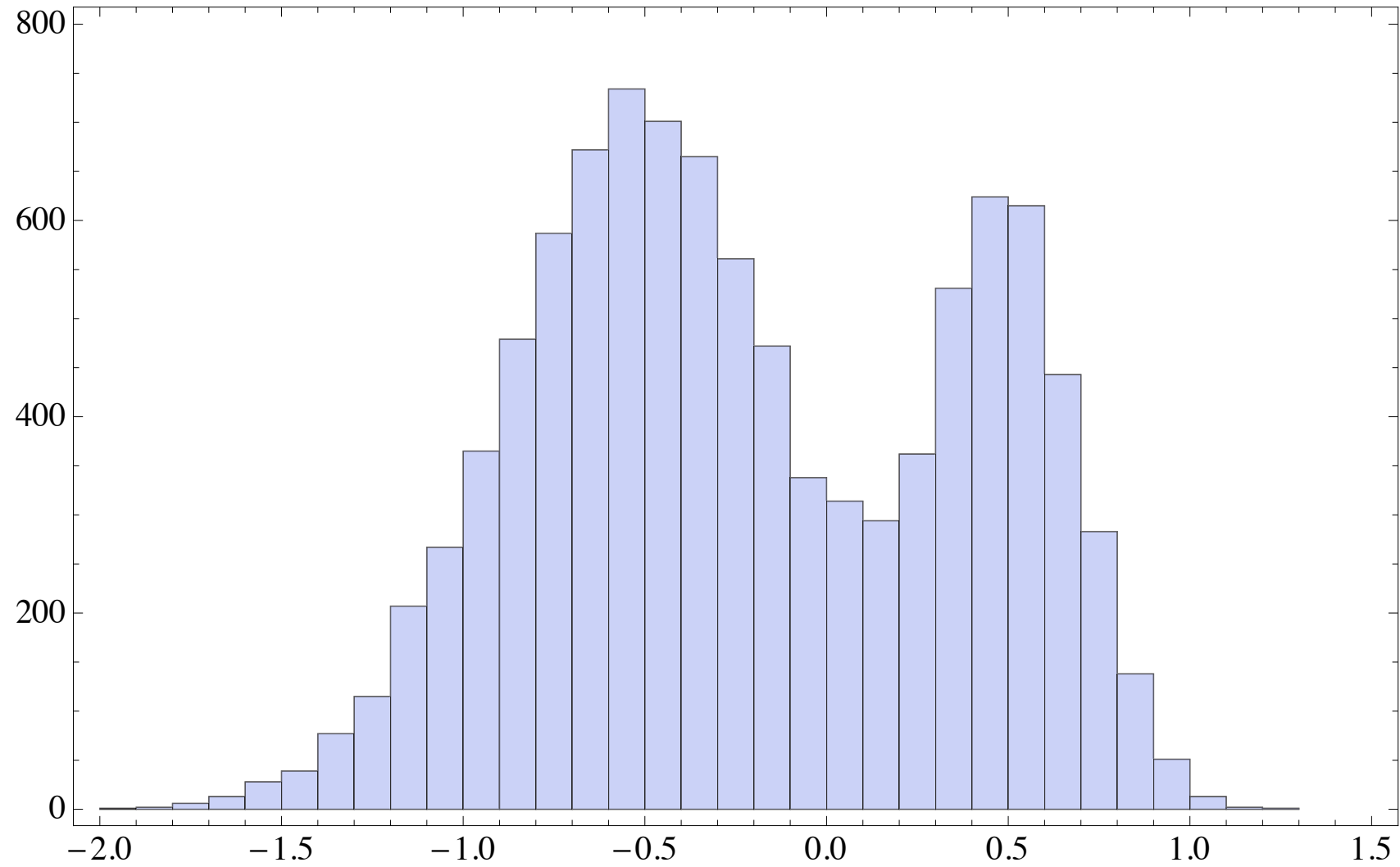
These methods are still not very efficient and there is a better alternative, the Markov Chain Monte Carlo method (see later)

3. *An important resampling technique: the Bootstrap method (B. Efron, 1977)*

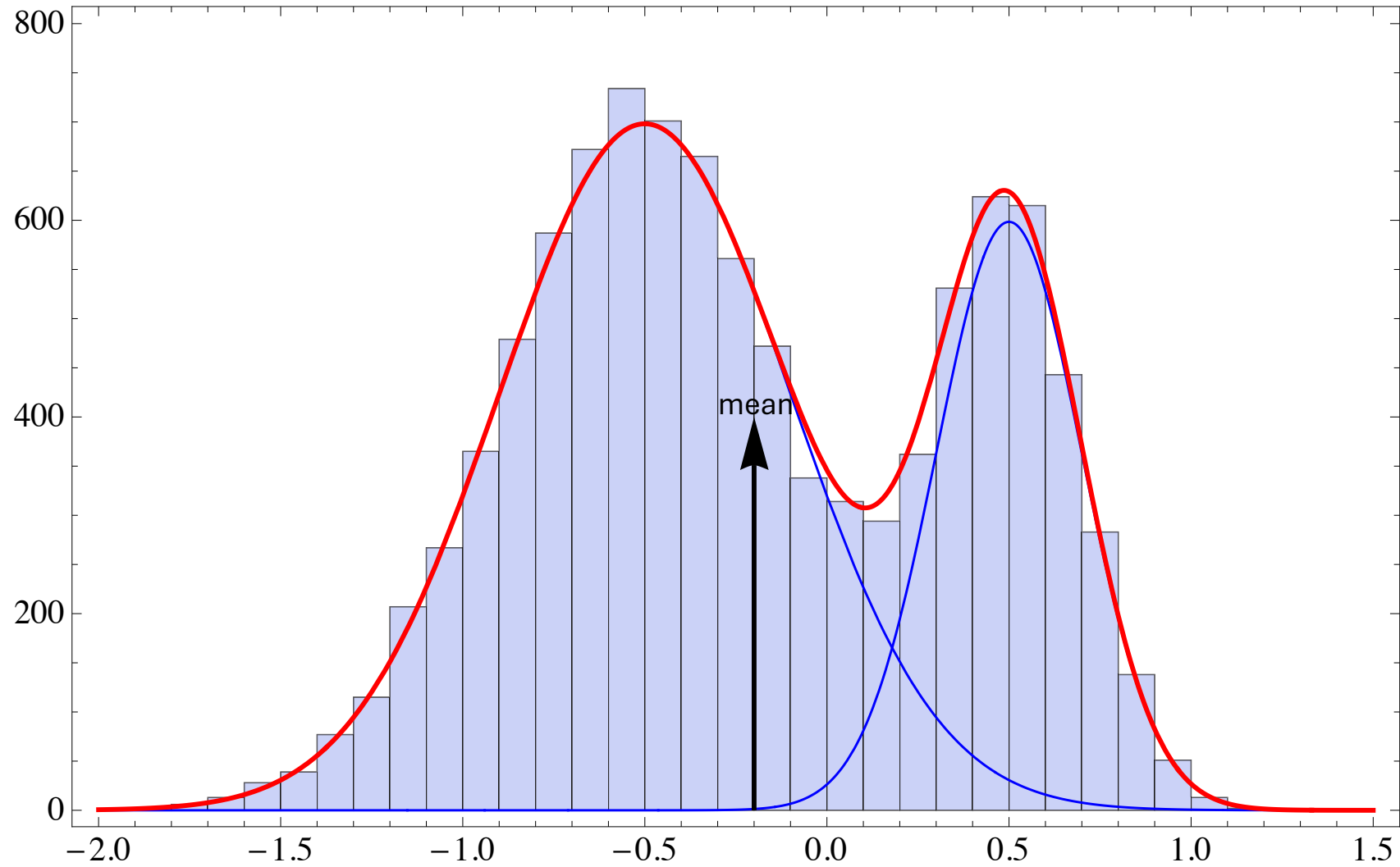


The bootstrap method is a resampling technique that helps calculating many statistical estimators

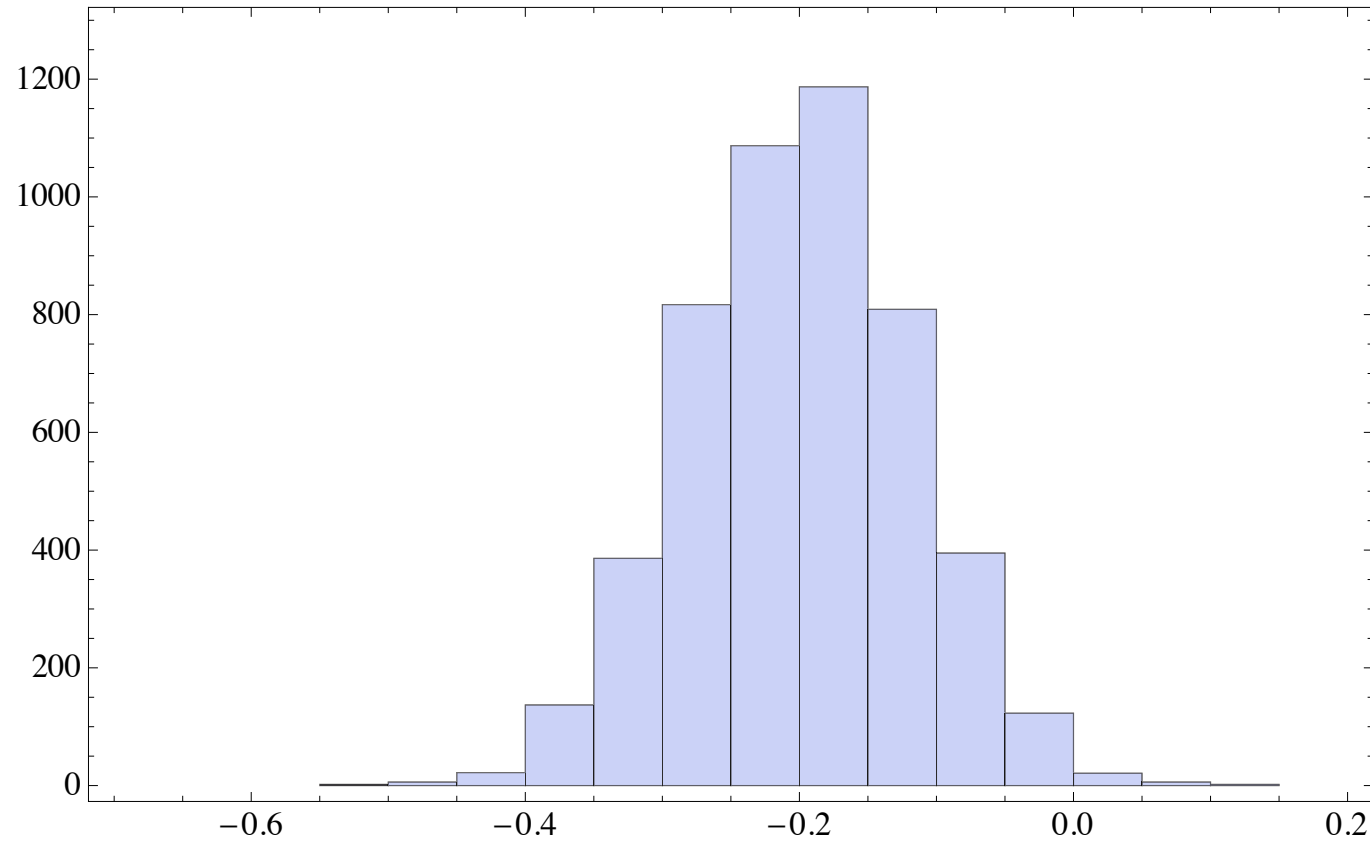
consider the distribution of a set of measurements



the distribution of data approximates the "true" underlying distribution (in this case a mixture model)

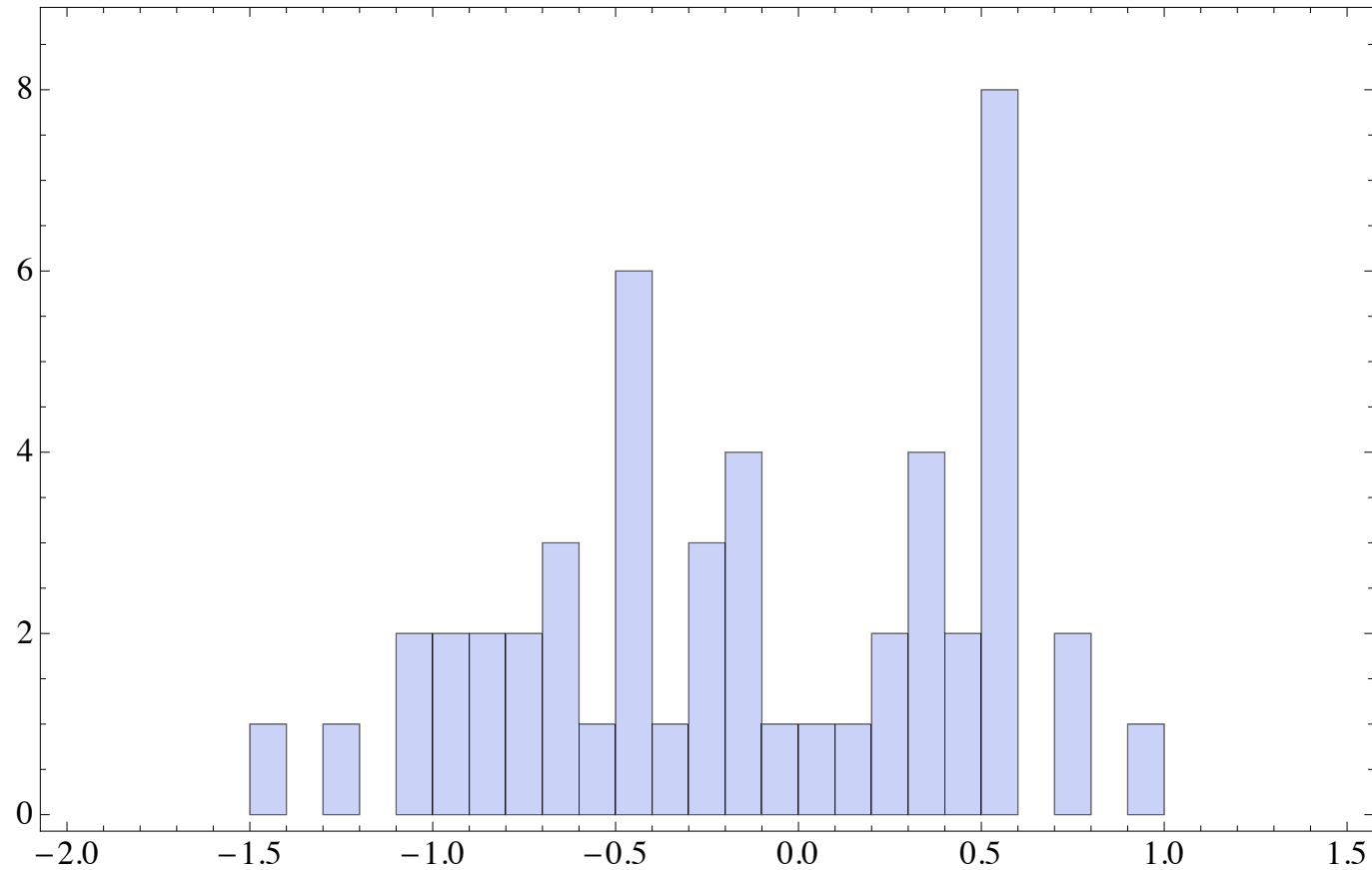


distribution of mean value obtained from 5000 sets of data (sample size = 50)



You can do this if you have large datasets ... but what if you have only a handful of measurements?

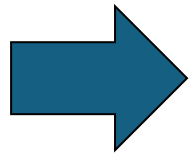
example: single dataset (same size as before, 50 measurements)



the discrete distribution is a rough representation of the underlying continuous distribution ... and yet it can be used just as before ...

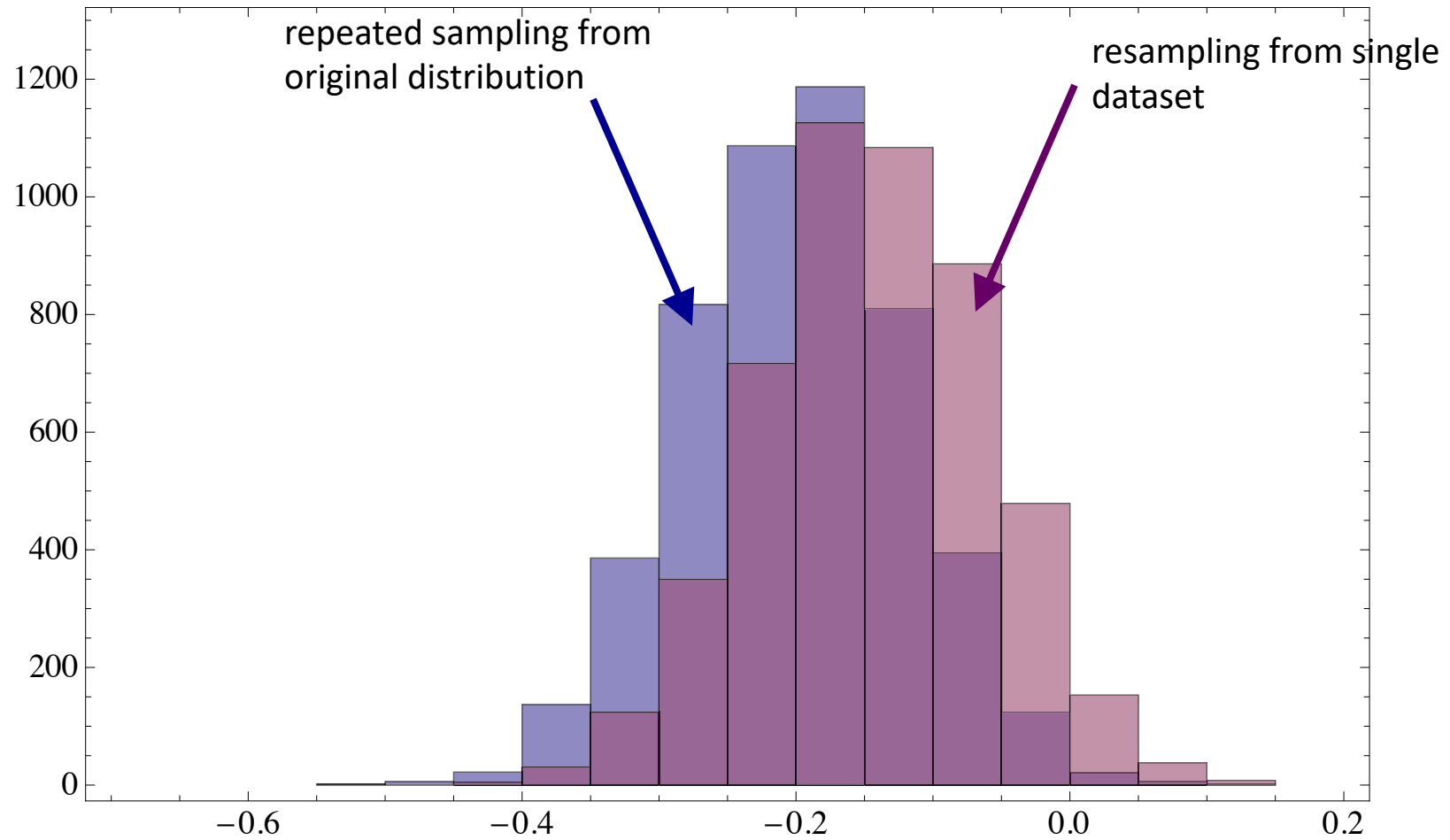
Bootstrap recipe:

if you need to find the distribution of the mean
(or any other statistical estimator) use the
dataset itself to generate new datasets



resample from dataset (with replacement)

distribution of mean value



true mean: -0.2

mean from repeated sampling (size = 250000): -0.200222 ± 0.0813632

mean from resampling dataset (size = 50): -0.142699 ± 0.0838678

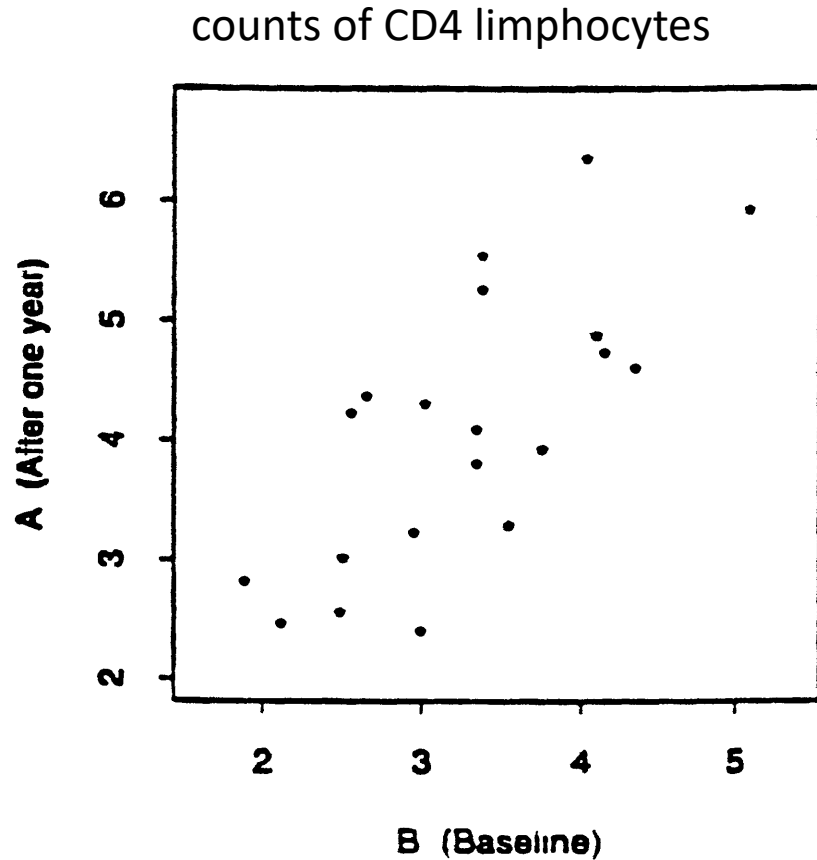


FIG. 1. *The cd4 data; cd4 counts in hundreds for 20 subjects, at baseline and after one year of treatment with an experimental anti-viral drug; numerical values appear in Table 1.*

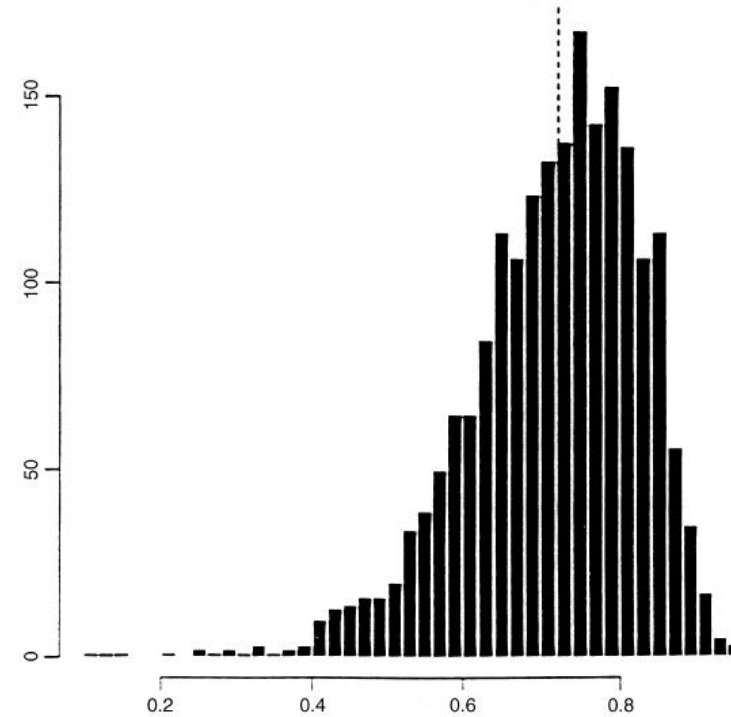


FIG. 3. *Histogram of 2,000 bootstrap correlation coefficients; bivariate normal sampling model.*

bootstrap estimate of correlation coefficient distribution

Example from Di Ciccio & Efron, *Statistics of Science* **11** (1996) 189 and Efron, *Statistics of Science* **13** (1998) 95

See Python example ...