

Introduction to Bayesian Methods- 2

Edoardo Milotti

Università di Trieste and INFN-Sezione di Trieste



Sunrise and eclipse over Graz on Jan. 4th 2011 (Robert Pölzl)

What is the probability that the sun rises tomorrow? (Laplace)

$$0 \leq \theta \leq 1$$

is the probability that the sun rises tomorrow, and we assume this probability to be uniformly distributed between 0 and 1.

Then

$$\begin{aligned} P(S) &= \sum_k P(SH_k) \\ &= \sum_k P(S|H_k)P(H_k) \\ &\rightarrow \int_0^1 P(S|\theta)p(\theta)d\theta \end{aligned}$$

Probability of the event “the sun rises tomorrow”

Probability of the event “the sun rises tomorrow” and that hypothesis H_k is true

hypothesis H_k corresponds to a certain value of p

$$P(S|\theta, N) = \theta$$

Probability that “the sun rises tomorrow”, if it already did so N times

$$\begin{aligned}
 p(\theta|N) &= \frac{P(N|\theta)}{P(N)} p(\theta) = \frac{P(N|\theta)}{\int_0^1 P(N|\theta)p(\theta)d\theta} p(\theta) \\
 &= \frac{\theta^N}{\int_0^1 \theta^N d\theta} = (N+1)\theta^N
 \end{aligned}$$



$$P(S|N) = \int_0^1 P(S|\theta, N) p(\theta|N) d\theta = \int_0^1 \theta \cdot (N+1)\theta^N d\theta = \frac{N+1}{N+2}$$

$$P(S|N) = \frac{N+1}{N+2} = 1 - \frac{1}{N+2}$$

Then, for a 27-year old person (about 10000 days), having observed that the sun rises every day means that

$$P(S|N) \approx 1 - \frac{1}{10000} \approx 0.99999$$

or also, that the probability that the sun does not rise tomorrow is about $1/10000 = 0.0001$

Bayesian inference

$$\begin{aligned} P(A_k | B) &= \frac{P(B | A_k) \cdot P(A_k)}{\sum_{k=1}^N P(B | A_k) \cdot P(A_k)} \\ &= \frac{P(B | A_k)}{\sum_{k=1}^N P(B | A_k) \cdot P(A_k)} \cdot P(A_k) \end{aligned}$$

$$P(H_k | D, I) = \frac{P(D | H_k, I)}{\sum_{k=1}^N P(D | H_k, I) \cdot P(H_k | I)} \cdot P(H_k | I)$$

$$P(H_k | D, I) = \frac{P(D | H_k, I)}{\sum_{k=1}^N P(D | H_k, I) \cdot P(H_k | I)} \cdot P(H_k | I)$$

(Posterior probability that k -th hypothesis is true, when we observe data D , with prior information I)

=

(Probability of observing data D , given the k -th hypothesis) / Normalization

.

(Prior probability that k -th hypothesis is true)

$$\begin{aligned}
P(H_k | D, I) &= \frac{P(D | H_k, I)}{P(D | I)} \cdot P(H_k | I) \\
&= \frac{P(D | H_k, I)}{\sum_{k=1}^N P(D | H_k, I) \cdot P(H_k | I)}
\end{aligned}$$

prior distribution $P(H_k, I)$

posterior distribution $P(H_k | D, I)$

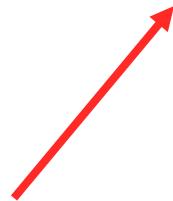
likelihood or sampling distribution $P(D | H_k, I)$

evidence (normalizing factor) $P(D | I) = \sum_{k=1}^N P(D | H_k, I) \cdot P(H_k | I)$

Testing hypotheses

$$P(H_k | D, I) = \frac{P(D | H_k, I)}{P(D | I)} \cdot P(H_k | I)$$

$$\frac{P(H_k | D, I)}{P(H_n | D, I)} = \left(\frac{P(D | H_k, I)}{P(D | H_n, I)} \right) \cdot \left(\frac{P(H_k | I)}{P(H_n | I)} \right)$$



Bayes' factor

When prior probabilities are the same (equally probable hypotheses), the posterior probability ratio depends only on the Bayes' factor:

$$\frac{P(H_k | D, I)}{P(H_n | D, I)} = \left(\frac{P(D | H_k, I)}{P(D | H_n, I)} \right)$$

Recursivity in a parameter estimate

$$P_2(\theta | \{d_1, d_2\}, I) = \frac{P(\{d_1, d_2\} | \theta, I)}{P(\{d_1, d_2\} | I)} \cdot P_0(\theta | I)$$

$$\begin{aligned} P(\{d_1, d_2\} | \theta, I) &= P(d_2 | d_1, \theta, I) \cdot P(d_1 | \theta, I) \\ &= P(d_2 | \theta, I) \cdot P(d_1 | \theta, I) \end{aligned}$$

here we assume that successive data
are independent

the same trick works
for the evidence

$$P(\{d_1, d_2\} | I) = P(d_2, I) \cdot P(d_1 | I)$$

$$\begin{aligned} P_2(\theta | \{d_1, d_2\}, I) &= \left(\frac{P(d_2 | \theta, I)}{P(d_2, I)} \right) \cdot \left(\frac{P(d_1 | \theta, I)}{P(d_1 | I)} \cdot P_0(\theta | I) \right) \\ &= \frac{P(d_2 | \theta, I)}{P(d_2, I)} \cdot P_1(\theta | I) \end{aligned}$$

A simple exercise (Skilling 1998)

Let T be the temperature of a liquid which can be either water or ethanol.

1. **We suppose first that the liquid is water:** then we take a uniform prior distribution for T , between 0°C and 100°C
2. The experimental apparatus and the measurement process is defined by the likelihood function $P(D|T, \text{water}, I)$. We assume that measurements are uniformly distributed within a range $\pm 5^\circ$. Therefore $P(D|T, \text{water}, I) = 0.1 (\text{ }^\circ\text{C})^{-1}$ in the interval $[T-5^\circ\text{C}, T+5^\circ\text{C}]$, and zero elsewhere.
3. We take a single measurement $D = -3^\circ\text{C}$.

4. The evidence $P(D)$ is

$$\begin{aligned} P(D|water, I) &= \int_T P(D|T, water, I)P(T)dT \\ &= \int_{0^\circ C}^{2^\circ C} \frac{(\circ C)^{-1}}{10} \cdot \frac{(\circ C)^{-1}}{100} dT (\circ C) = 0.002 (\circ C)^{-1} \end{aligned}$$

5. Using Bayes' theorem we find

$$\begin{aligned} P(T|D, water, I) &= \frac{P(D|T, water, I)}{P(D, water, I)} P(T) = \frac{0.1 (\circ C)^{-1}}{0.002 (\circ C)^{-1}} 0.01 (\circ C)^{-1} \\ &= 0.5 (\circ C)^{-1} \quad (0^\circ C < T < 2^\circ C) \end{aligned}$$

Now suppose that the liquid is ethanol, so that the temperature range is $-80^{\circ}\text{C} < T < 80^{\circ}\text{C}$

1. $P(T) = (160^{\circ}\text{C})^{-1}$ in $-80^{\circ}\text{C} < T < 80^{\circ}\text{C}$.
2. $P(D|T, \text{ethanol}, I) = 0.1 (\text{ }^{\circ}\text{C})^{-1}$ in $[T-5^{\circ}\text{C}, T+5^{\circ}\text{C}]$, and zero elsewhere.
3. We take a single measurement $D = -3^{\circ}\text{C}$.
4. The evidence $P(D, \text{ethanol}, I)$ is

$$P(D, \text{ethanol}, I) = \int_T P(D|T, \text{ethanol}, I) P(T) dT = \int_{-8^{\circ}\text{C}}^{2^{\circ}\text{C}} \frac{(\text{ }^{\circ}\text{C})^{-1}}{10} \cdot \frac{(\text{ }^{\circ}\text{C})^{-1}}{160} dT (\text{ }^{\circ}\text{C}) = 0.00625 (\text{ }^{\circ}\text{C})^{-1}$$

5. Using Bayes' theorem we find

$$P(T|D, \text{ethanol}, I) = \frac{P(D|T, \text{ethanol}, I)}{P(D, \text{ethanol}, I)} P(T) = \frac{0.1 (\text{ }^{\circ}\text{C})^{-1}}{0.00625 (\text{ }^{\circ}\text{C})^{-1}} \frac{1}{160} (\text{ }^{\circ}\text{C})^{-1} = 0.1 (\text{ }^{\circ}\text{C})^{-1}$$
$$(-8^{\circ}\text{C} < T < 2^{\circ}\text{C})$$

Assuming a prior for the water-ethanol choice, we can discriminate between water and ethanol

$$P_{\text{water}} = P_{\text{ethanol}} = 0.5$$

Indeed,

$$\begin{aligned} P(\text{water}|D,I) &= \frac{P(D|\text{water},I)}{P(D|\text{water},I)P(\text{water},I) + P(D|\text{ethanol},I)P(\text{ethanol},I)} P(\text{water},I) \\ &= \frac{P(D|\text{water},I)}{P(D|\text{water},I) + P(D|\text{ethanol},I)} \end{aligned}$$

and therefore the ratio of the posteriors is given by the Bayes' factor

$$\frac{P(\text{water}|D,I)}{P(\text{ethanol}|D,I)} = \frac{P(D|\text{water},I)}{P(D|\text{ethanol},I)}$$

We have found earlier that

$$P(D|water) = 0.002(\text{ }^{\circ}\text{C})^{-1}$$

$$P(D|ethanol) = 0.00625(\text{ }^{\circ}\text{C})^{-1}$$

therefore

$$\frac{P(ethanol|D,I)}{P(water|D,I)} = \frac{P(D|ethanol,I)}{P(D|water,I)} = 3.125$$

and we conclude that the observation favors the hypothesis of liquid ethanol.

$\log_{10}(B)$	B	Evidence support
0 to $1/2$	1 to 3.2	Not worth more than a bare mention
$1/2$ to 1	3.2 to 10	Substantial
1 to 2	10 to 100	Strong
> 2	> 100	Decisive

Interpretation of the Bayes factor B as evidence support according to Jeffreys.

In the case of the water-ethanol problem, and according to Jeffreys' categories, the preference for ethanol is “not worth more than a bare mention”, although it happens to be in the upper part of the range.

From discrete sets of hypothesis to the continuum. The Bayes' theorem in the context of parameter estimation.

$$P(H_k | D, I) = \frac{P(D | H_k, I)}{P(D | I)} \cdot P(H_k | I) = \frac{P(D | H_k, I)}{\sum_{k=1}^N P(D | H_k, I) \cdot P(H_k | I)}$$



$$dP(\theta | D, I) = \frac{P(D | \theta, I)}{\int_{\Theta} P(D | \theta, I) dP(\theta | I)} \cdot dP(\theta | I)$$

$$\frac{dP(\theta | D, I)}{d\theta} = \frac{P(D | \theta, I)}{\int_{\Theta} P(D | \theta, I) \cdot \frac{dP(\theta | I)}{d\theta} d\theta} \cdot \frac{dP(\theta | I)}{d\theta}$$

1. Example of Bayesian inference: estimate of the (probability) parameter of the binomial distribution

$$P(n|\theta, N) = \binom{N}{n} (1-\theta)^{N-n} \theta^n$$

this is the parameter that we want to infer from data

$$p(\theta|n, N) = \frac{P(n|\theta, N)}{\int_0^1 P(n|\theta, N) \cdot p(\theta) d\theta} \cdot p(\theta) =$$

uniform distribution: the least informative prior

$$= \frac{\binom{N}{n} (1-\theta)^{N-n} \theta^n}{\int_0^1 \binom{N}{n} (1-\theta)^{N-n} \theta^n \cdot p(\theta) d\theta} \cdot p(\theta) = \frac{(1-\theta)^{N-n} \theta^n}{\int_0^1 (1-\theta)^{N-n} \theta^n d\theta}$$

final result is a beta distribution

$$p(\theta | n, N) = \frac{(1-\theta)^{N-n} \theta^n}{\int_0^1 \theta^n (1-\theta)^{N-n} d\theta} = \frac{(1-\theta)^{N-n} \theta^n}{B(n+1, N-n+1)}$$

$$B(m, n) = \int_0^1 t^{m-1} (1-t)^{n-1} dt$$

beta function

$$= \frac{\Gamma(m)\Gamma(n)}{\Gamma(m+n)}$$

$$p(\theta | n, N) = \frac{\Gamma(N+2)}{\Gamma(n+1)\Gamma(N-n+1)} (1-\theta)^{N-n} \theta^n$$

$$= \frac{(N+1)!}{n!(N-n)!} (1-\theta)^{N-n} \theta^n$$

Mathematical digression: relationship between gamma and beta function

$$\Gamma(m)\Gamma(n) = \int_0^\infty s^{m-1} e^{-s} ds \int_0^\infty t^{n-1} e^{-t} dt$$

$$s = x^2; \quad t = y^2; \quad \Rightarrow$$

$$\Gamma(m)\Gamma(n) = 4 \int_0^\infty x^{2m-1} e^{-x^2} dx \int_0^\infty y^{2n-1} e^{-y^2} dy$$

$$x = r \cos \theta; \quad y = r \sin \theta; \quad \Rightarrow$$

$$\Gamma(m)\Gamma(n) = 4 \int_0^\infty r^{2m+2n-1} e^{-r^2} dr \int_0^{\pi/2} \cos^{2m-1} \theta \sin^{2n-1} \theta d\theta$$

$$= \Gamma(m+n) \left(2 \int_0^{\pi/2} \cos^{2m-1} \theta \sin^{2n-1} \theta d\theta \right) \quad (t = \cos^2 \theta; \quad dt = -2 \cos \theta \sin \theta d\theta)$$

$$= \Gamma(m+n) \int_0^1 t^{m-1} (1-t)^{n-1} dt$$

$$= \Gamma(m+n) B(m,n)$$

$$\Rightarrow B(m,n) = \frac{\Gamma(m)\Gamma(n)}{\Gamma(m+n)} \quad \Rightarrow \quad B(m+1,n+1) = \frac{m!n!}{(m+n+1)!}$$

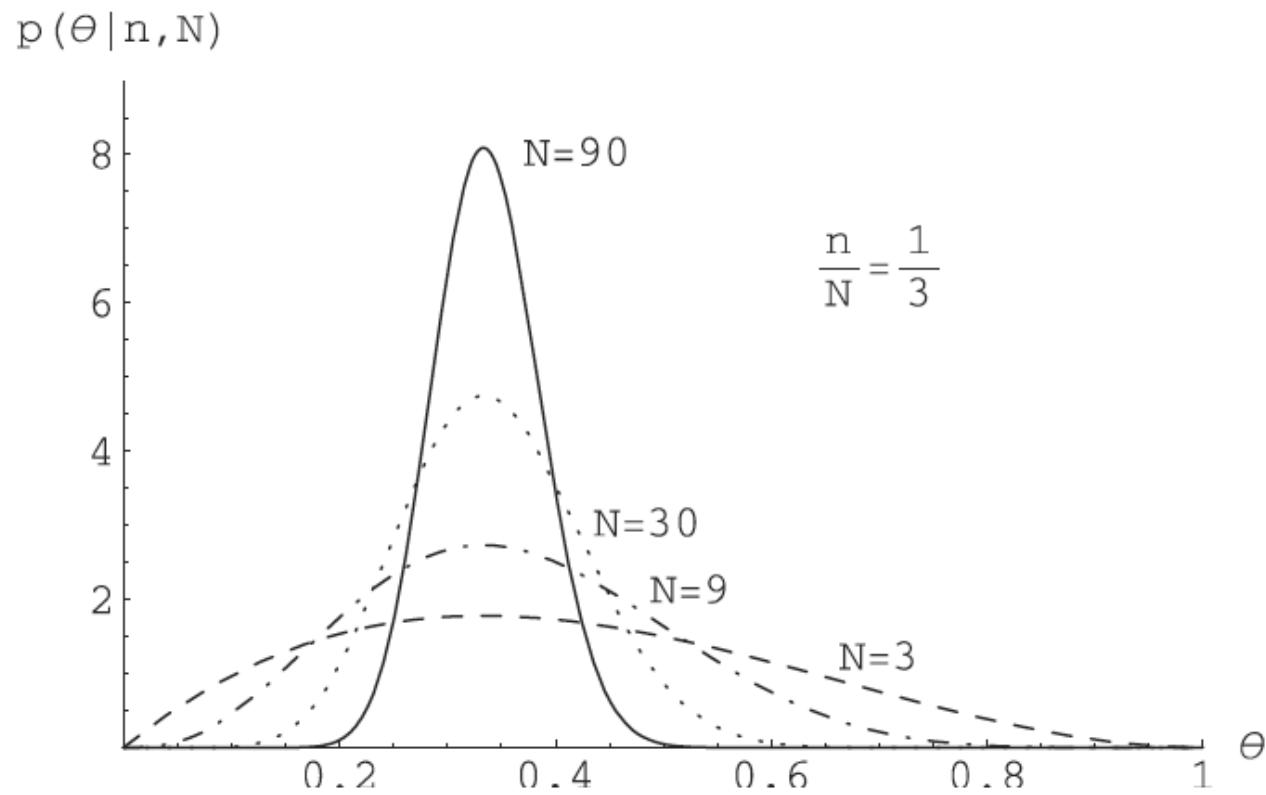


Figure 1. Posterior probability density function of the binomial parameter θ , having observed n successes in N trials.

From the knowledge of the posterior pdf we obtain all the momenta of the distribution

$$p(\theta | n, N) = \frac{(N+1)!}{n!(N-n)!} (1-\theta)^{N-n} \theta^n$$



$$\begin{aligned}\langle \theta \rangle &= \int_0^1 p(\theta | n, N) \theta d\theta = \frac{(N+1)!}{n!(N-n)!} \int_0^1 (1-\theta)^{N-n} \theta^{n+1} d\theta \\ &= \frac{(N+1)!}{n!(N-n)!} B(n+2, N-n+1) \\ &= \frac{(N+1)!}{n!(N-n)!} \cdot \frac{(n+1)!(N-n)!}{(N+2)!} \\ &= \frac{n+1}{N+2} \rightarrow \frac{n}{N}\end{aligned}$$

biased, asymptotically unbiased,
estimator

$$\begin{aligned}
\langle \theta^2 \rangle &= \int_0^1 p(\theta | n, N) \theta^2 d\theta = \frac{(N+1)!}{n!(N-n)!} \int_0^1 (1-\theta)^{N-n} \theta^{n+2} d\theta \\
&= \frac{(N+1)!}{n!(N-n)!} B(n+3, N-n+1) \\
&= \frac{(N+1)!}{n!(N-n)!} \cdot \frac{(n+2)!(N-n)!}{(N+3)!} \\
&= \frac{(n+2)(n+1)}{(N+3)(N+2)}
\end{aligned}$$

$$\begin{aligned}
\text{var } \theta &= \langle \theta^2 \rangle - \langle \theta \rangle^2 = \frac{(n+2)(n+1)}{(N+3)(N+2)} - \left(\frac{n+1}{N+2} \right)^2 = \\
&= \frac{(N-n+1)(n+1)}{(N+3)(N+2)^3}
\end{aligned}$$

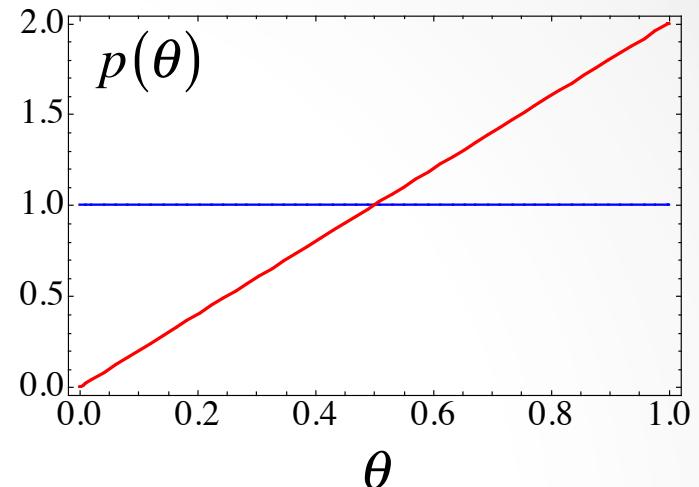
What happens if we try a different prior?

Let's try with a linear prior

$$p(\theta) = 2\theta$$

$$p(\theta | n, N) = \frac{P(n | \theta, N)}{\int_0^1 P(n | \theta, N) \cdot p(\theta) d\theta} \cdot p(\theta)$$

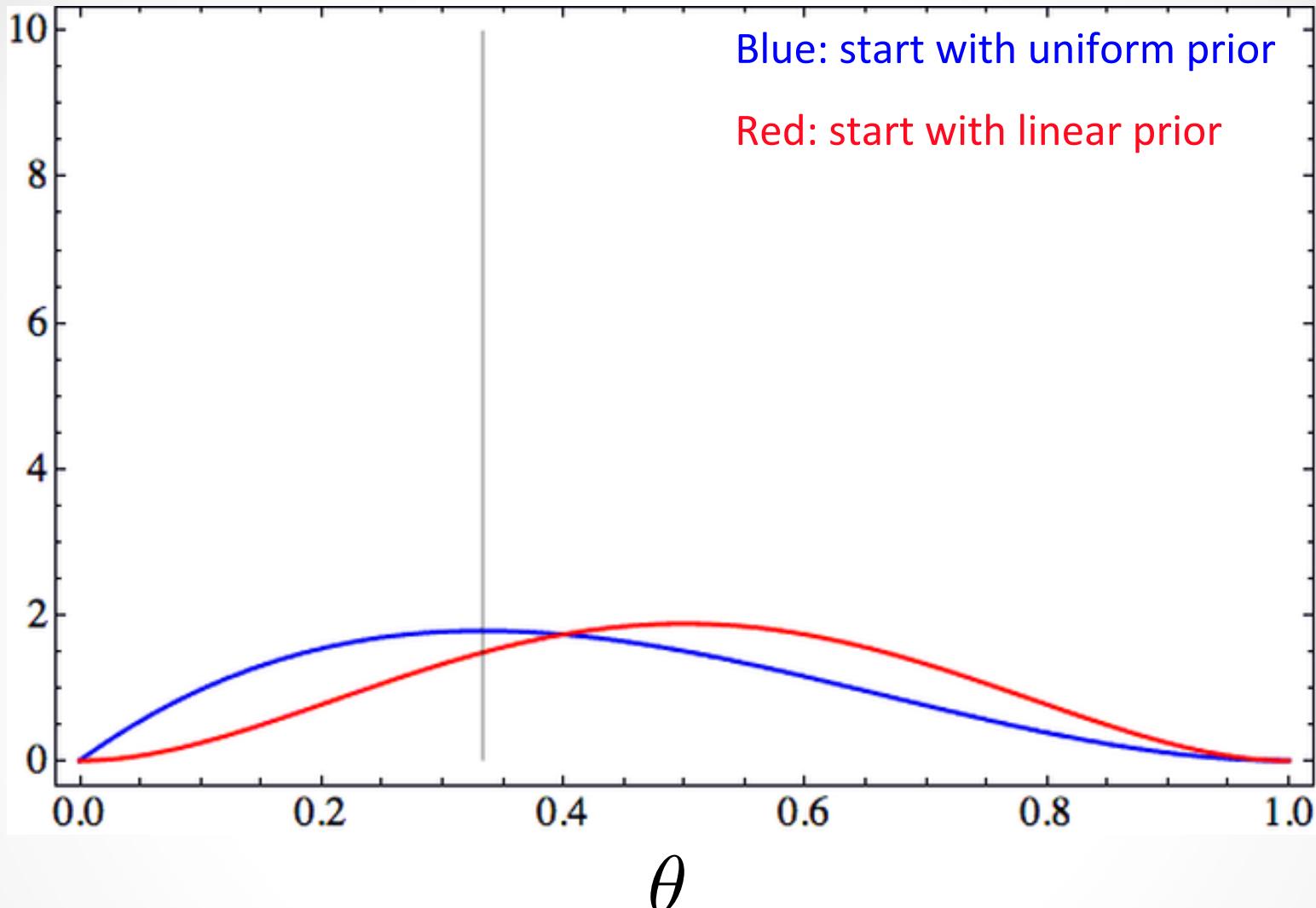
$$= \frac{\binom{N}{n} (1-\theta)^{N-n} \theta^n}{\int_0^1 \binom{N}{n} (1-\theta)^{N-n} \theta^n \cdot 2\theta d\theta} \cdot 2\theta = \frac{(1-\theta)^{N-n} \theta^{n+1}}{\int_0^1 (1-\theta)^{N-n} \theta^{n+1} d\theta}$$



$$p(\theta | n, N) = \frac{(N+2)!}{(n+1)!(N-n)!} (1-\theta)^{N-n} \theta^{n+1}$$



$$\begin{aligned}
 \langle \theta \rangle &= \int_0^1 p(\theta | n, N) \theta d\theta = \frac{(N+2)!}{(n+1)!(N-n)!} \int_0^1 (1-\theta)^{N-n} \theta^{n+2} d\theta \\
 &= \frac{(N+2)!}{(n+1)!(N-n)!} B(n+3, N-n+1) \\
 &= \frac{(N+2)!}{(n+1)!(N-n)!} \cdot \frac{(n+2)!(N-n)!}{(N+3)!} \\
 &= \frac{n+2}{N+3} \rightarrow \frac{n}{N}
 \end{aligned}$$

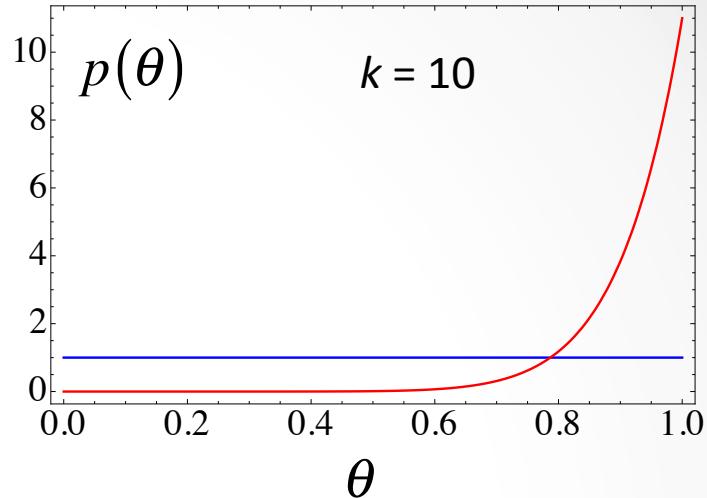


Taking few coin throws, the posterior from the linear prior is considerably biased. The bias disappears when the number of coin throws is large.

Now we try with a very non-uniform prior

We take

$$p(\theta) = (k+1)\theta^k; \quad k \gg 1$$

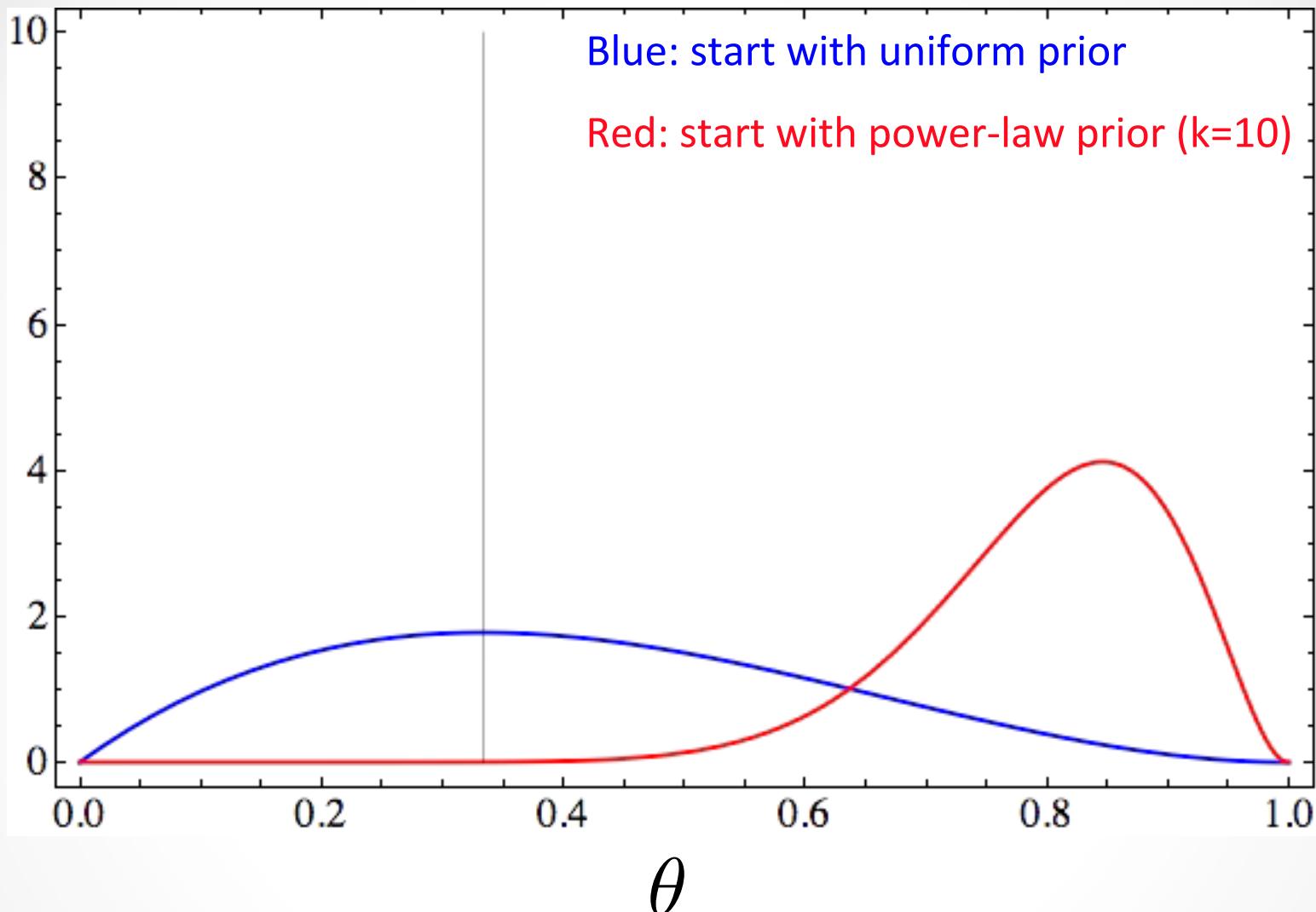


$$\begin{aligned}
 p(\theta | n, N) &= \frac{p(n | \theta, N)}{\int_0^1 P(n | \theta, N) \cdot p(\theta) d\theta} \cdot p(\theta) \\
 &= \frac{\binom{N}{n} (1-\theta)^{N-n} \theta^n}{\int_0^1 \binom{N}{n} (1-\theta)^{N-n} \theta^n \cdot (k+1)\theta^k d\theta} \cdot (k+1)\theta^k = \frac{(1-\theta)^{N-n} \theta^{n+k}}{\int_0^1 (1-\theta)^{N-n} \theta^{n+k} d\theta}
 \end{aligned}$$

$$p(\theta | n, N) = \frac{(N+k+1)!}{(n+k)!(N-n)!} (1-\theta)^{N-n} \theta^{n+k}$$



$$\begin{aligned}
 \langle \theta \rangle &= \int_0^1 p(\theta | n, N) \theta d\theta = \frac{(N+k+1)!}{(n+k)!(N-n)!} \int_0^1 (1-\theta)^{N-n} \theta^{n+k+1} d\theta \\
 &= \frac{(N+k+1)!}{(n+k)!(N-n)!} B(n+k+2, N-n+1) \\
 &= \frac{(N+k+1)!}{(n+k)!(N-n)!} \cdot \frac{(n+k+1)!(N-n)!}{(N+k+2)!} \\
 &= \frac{n+k+1}{N+k+2} \rightarrow \frac{n}{N}
 \end{aligned}$$



In this case, initial bias due to the prior is very large.

Note on posterior distributions:

the relationship between binomial distribution and beta function is quite important and common, and leads to the formal definition of the Beta distribution:

$$B(\theta|a,b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \theta^{a-1} (1-\theta)^{b-1}$$

There are other important dualities between distributions. This topic is discussed in depth in

J. M. Bernardo: “Reference Posterior Distributions for Bayesian Inference”, J. R. Statist. Soc. B **41** (1979), 113

Maximum a posteriori (MAP) estimate

Consider the case with a uniform prior: from the posterior distribution

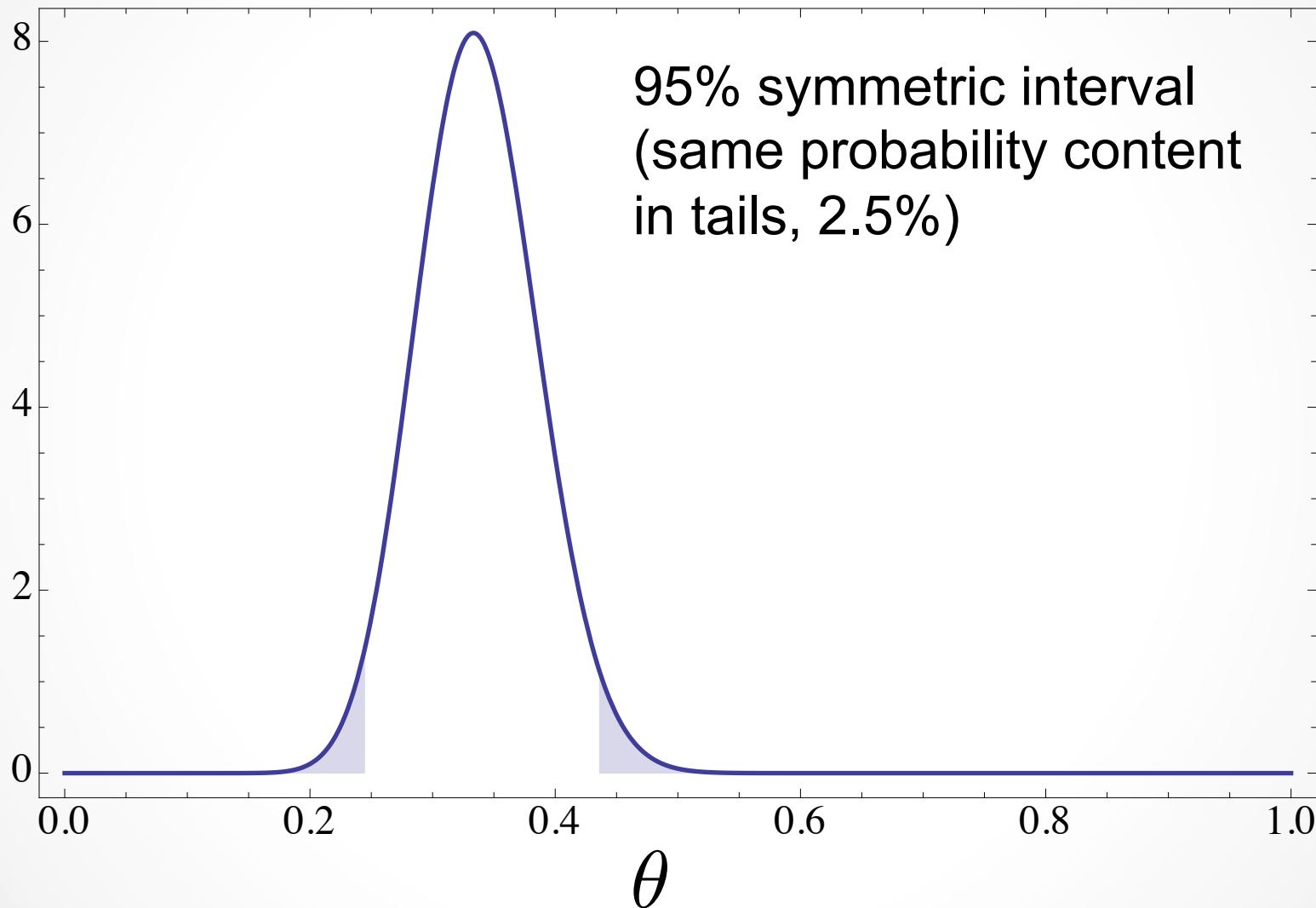
$$p(\theta | n, N) = \frac{(N+1)!}{n!(N-n)!} (1-\theta)^{N-n} \theta^n$$

we easily find that the posterior pdf is maximized by the parameter value

$$\theta = n/N$$

which is the unbiased estimate of the parameter (unlike the mean value ...)

Confidence intervals (case of initial uniform prior)



2. Example of Bayesian inference: estimate of a Poissonian rate

$$P(n | \lambda, \Delta t) = \frac{(\lambda \Delta t)^n e^{-\lambda \Delta t}}{n!}$$

uniform distribution (step function, not normalizable, improper prior!)

$$\begin{aligned} p(\lambda | n, \Delta t) &= \frac{P(n | \lambda, \Delta t)}{\int_{\lambda} P(n | \lambda, \Delta t) \cdot p(\lambda) d\lambda} \cdot p(\lambda) = \\ &= \frac{(\lambda \Delta t)^n e^{-\lambda \Delta t}}{\int_0^{\infty} (\lambda \Delta t)^n e^{-\lambda \Delta t} d(\lambda \Delta t)} = \frac{(\lambda \Delta t)^n e^{-\lambda \Delta t}}{n!} \end{aligned}$$

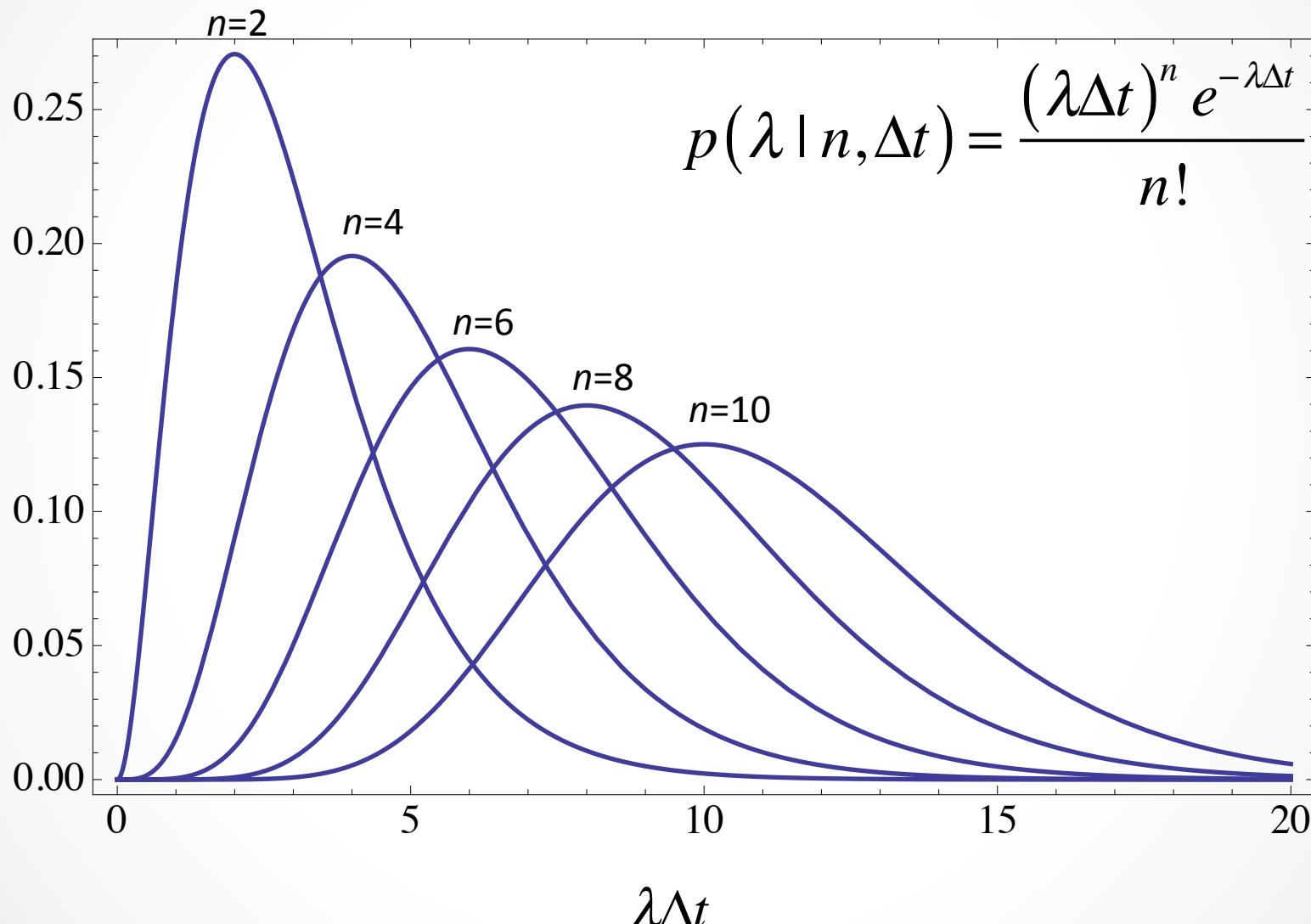
here the random variable is the parameter λ

$$\begin{aligned}\langle \lambda \Delta t \rangle &= \int_0^\infty p(\lambda | n, \Delta t) (\lambda \Delta t) d(\lambda \Delta t) = \int_0^\infty \frac{(\lambda \Delta t)^n e^{-\lambda \Delta t}}{n!} (\lambda \Delta t) d(\lambda \Delta t) \\ &= \frac{1}{n!} \int_0^\infty (\lambda \Delta t)^{n+1} e^{-\lambda \Delta t} d(\lambda \Delta t) = \frac{(n+1)!}{n!} = n+1\end{aligned}$$

$$\begin{aligned}\langle (\lambda \Delta t)^2 \rangle &= \int_0^\infty \frac{(\lambda \Delta t)^n e^{-\lambda \Delta t}}{n!} (\lambda \Delta t)^2 d(\lambda \Delta t) \\ &= \frac{1}{n!} \int_0^\infty (\lambda \Delta t)^{n+2} e^{-\lambda \Delta t} d(\lambda \Delta t) = \frac{(n+2)!}{n!} = (n+2)(n+1)\end{aligned}$$

$$\text{var}(\lambda \Delta t) = \langle (\lambda \Delta t)^2 \rangle - \langle \lambda \Delta t \rangle^2 = (n+2)(n+1) - (n+1)^2 = n+1$$

Posterior distributions of the adimensional parameter $\lambda\Delta t$ when n events are observed



3. Example of Bayesian inference: Gaussian model (estimate of the mean with given distribution variance)

$$p(d|\mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(d-\mu)^2}{2\sigma^2}\right]$$

1. we know nothing about the position of the mean

$$\begin{aligned} p(\mu|d, \sigma) &= \frac{P(d|\mu, \sigma)}{\int_{\mu} P(d|\mu, \sigma) \cdot p(\mu) d\mu} \cdot p(\mu) = \\ &= \frac{\exp\left[-(d-\mu)^2/2\sigma^2\right]}{\int_{-\infty}^{+\infty} \exp\left[-(d-\mu)^2/2\sigma^2\right] d\mu} = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(d-\mu)^2}{2\sigma^2}\right] \end{aligned}$$

uniform prior (improper!)

here the random variable is μ

2. positive mean (now we know that $\mu > 0$)

step function (another improper prior)

$$\begin{aligned} p(\mu | d, \sigma) &= \frac{P(d | \mu, \sigma)}{\int_{\mu} P(d | \mu, \sigma) \cdot p(\mu) d\mu} \cdot p(\mu) = \\ &= \frac{\exp[-(d - \mu)^2 / 2\sigma^2]}{\int_0^{+\infty} \exp[-(d - \mu)^2 / 2\sigma^2] d\mu} \cdot s(\mu) = N \exp\left[-\frac{(d - \mu)^2}{2\sigma^2}\right] \cdot s(\mu) \end{aligned}$$

normalization

NB: because of the step function, posteriors always vanish for $\mu < 0$

3. Gaussian prior

$$P(\mu | \mu_1, I) = \frac{1}{\sqrt{2\pi\sigma_1^2}} \exp\left[-\frac{(\mu - \mu_1)^2}{2\sigma_1^2}\right]$$



this is an informative prior, we know that the mean is centered about μ_1

$$\begin{aligned} p(\mu | d, \sigma) &= \frac{P(d | \mu, \sigma)}{\int_{-\infty}^{+\infty} P(d | \mu, \sigma) \cdot p(\mu) d\mu} \cdot p(\mu) = \\ &= \frac{\exp\left[-\frac{(d - \mu)^2}{2\sigma^2}\right]}{\int_{-\infty}^{+\infty} \exp\left[-\frac{(d - \mu)^2}{2\sigma^2}\right] \exp\left[-\frac{(\mu - \mu_1)^2}{2\sigma_1^2}\right] d\mu} \cdot \exp\left[-\frac{(\mu - \mu_1)^2}{2\sigma_1^2}\right] \end{aligned}$$

$$p(\mu | d, \sigma) = \frac{\exp\left[-\frac{(d-\mu)^2}{2\sigma^2} - \frac{(\mu-\mu_1)^2}{2\sigma_1^2}\right]}{\int_{-\infty}^{+\infty} \exp\left[-\frac{(d-\mu)^2}{2\sigma^2} - \frac{(\mu-\mu_1)^2}{2\sigma_1^2}\right] d\mu}$$

this is a Gaussian expression and we can rearrange the exponent

$$\begin{aligned} \frac{(d-\mu)^2}{\sigma^2} + \frac{(\mu-\mu_1)^2}{\sigma_1^2} &= \frac{\sigma_1^2(d-\mu)^2 + \sigma^2(\mu-\mu_1)^2}{\sigma^2\sigma_1^2} = \frac{(\sigma_1^2 + \sigma^2)\mu^2 - 2(\sigma_1^2 d + \sigma^2 \mu_1)\mu + (\sigma_1^2 d^2 + \sigma^2 \mu_1^2)}{\sigma^2\sigma_1^2} \\ &= \frac{1}{\sigma^2\sigma_1^2} \left[(\sigma_1^2 + \sigma^2)\mu^2 - 2 \frac{(\sigma_1^2 d + \sigma^2 \mu_1)}{\sqrt{\sigma_1^2 + \sigma^2}} \sqrt{\sigma_1^2 + \sigma^2} \mu + \frac{(\sigma_1^2 d + \sigma^2 \mu_1)^2}{(\sigma_1^2 + \sigma^2)} - \frac{(\sigma_1^2 d + \sigma^2 \mu_1)^2}{(\sigma_1^2 + \sigma^2)} + (\sigma_1^2 d^2 + \sigma^2 \mu_1^2) \right] \\ &= \frac{1}{\sigma^2\sigma_1^2} \left[(\sigma_1^2 + \sigma^2) \left(\mu - \frac{(\sigma_1^2 d + \sigma^2 \mu_1)}{(\sigma_1^2 + \sigma^2)} \right)^2 - \frac{(\sigma_1^2 d + \sigma^2 \mu_1)^2}{(\sigma_1^2 + \sigma^2)} + (\sigma_1^2 d^2 + \sigma^2 \mu_1^2) \right] \end{aligned}$$

$$\begin{aligned}
p(\mu | d, \sigma) &= \frac{\exp \left\{ -\frac{1}{2\sigma^2\sigma_1^2} \left[(\sigma_1^2 + \sigma^2) \left(\mu - \frac{(\sigma_1^2 d + \sigma^2 \mu_1)}{(\sigma_1^2 + \sigma^2)} \right)^2 - \frac{(\sigma_1^2 d + \sigma^2 \mu_1)^2}{(\sigma_1^2 + \sigma^2)} + (\sigma_1^2 d^2 + \sigma^2 \mu_1^2) \right] \right\}}{\int_{-\infty}^{+\infty} \exp \left\{ -\frac{1}{2\sigma^2\sigma_1^2} \left[(\sigma_1^2 + \sigma^2) \left(\mu - \frac{(\sigma_1^2 d + \sigma^2 \mu_1)}{(\sigma_1^2 + \sigma^2)} \right)^2 - \frac{(\sigma_1^2 d + \sigma^2 \mu_1)^2}{(\sigma_1^2 + \sigma^2)} + (\sigma_1^2 d^2 + \sigma^2 \mu_1^2) \right] \right\} d\mu} \\
&= \frac{1}{\sqrt{2\pi}} \sqrt{\frac{\sigma^2 \sigma_1^2}{(\sigma_1^2 + \sigma^2)}} \exp \left\{ -\frac{(\sigma_1^2 + \sigma^2)}{2\sigma^2\sigma_1^2} \left(\mu - \frac{(\sigma_1^2 d + \sigma^2 \mu_1)}{(\sigma_1^2 + \sigma^2)} \right)^2 \right\}
\end{aligned}$$

$$\boxed{
\begin{aligned}
\mu_2 &= \frac{\sigma_1^2 d + \sigma^2 \mu_1}{\sigma_1^2 + \sigma^2} \\
\sigma_2^2 &= \frac{\sigma^2 \sigma_1^2}{\sigma_1^2 + \sigma^2} \quad \left(\frac{1}{\sigma_2^2} = \frac{1}{\sigma^2} + \frac{1}{\sigma_1^2} \right)
\end{aligned}}$$

Recursivity of inferential process

$$\sigma_2^2 = \frac{\sigma^2 \sigma_1^2}{\sigma_1^2 + \sigma^2} \quad \left(\frac{1}{\sigma_2^2} = \frac{1}{\sigma^2} + \frac{1}{\sigma_1^2} \right)$$

$$\mu_2 = \frac{d/\sigma^2 + \mu_1/\sigma_1^2}{1/\sigma^2 + 1/\sigma_1^2}$$

weighted mean

$$\mu_2 = \frac{\sigma_1^2 d + \sigma^2 \mu_1}{\sigma_1^2 + \sigma^2} = \frac{\sigma_1^2 d - \sigma_1^2 \mu_1 + (\sigma^2 + \sigma_1^2) \mu_1}{\sigma_1^2 + \sigma^2} = \mu_1 + \frac{\sigma_1^2}{\sigma_1^2 + \sigma^2} (d - \mu_1)$$

same structure as a Kalman filter

Kalman filter (super-simple example)

$$\mu_2 = \mu_1 + K(d_2 - \mu_1)$$

$$K = \frac{\sigma_1^2}{\sigma_1^2 + \sigma^2} \quad (\text{filter gain})$$

+ linear variation (with Gaussian noise)

$$\mu(t_3) = \mu(t_2) + v(t_3 - t_2)$$

$$\sigma_\mu^2(t_3) = \sigma_\mu^2(t_2) + \sigma_v^2(t_3 - t_2)$$

$$\mu_3 = \mu_2 + K_3(d_3 - \mu_2)$$

$$K_3 = \frac{\sigma_\mu^2(t_3)}{\sigma_\mu^2(t_3) + \sigma_3^2}$$

References:

- G. D' Agostini, Rep. Prog. Phys. **66** (2003) 1383
- V. Dose, Rep. Prog. Phys. **66** (2003) 1421
- J. Skilling, J. of Microscopy **190** (1998) 28
- P. S. Maybeck, *Stochastic models, estimation, and control (vol. 1)*, Academic Press 1979