# Introduction to Bayesian Inference: Supplemental Topics

Tom Loredo
Cornell Center for Astrophysics and Planetary Science
http://www.astro.cornell.edu/staff/loredo/bayes/

CASt Summer School — 7–8 June 2017

# Supplemental Topics

**1** **Parametric bootstrapping vs. posterior sampling**

**2** **Estimation and model comparison for binary outcomes**

**3** **Basic inference with normal errors**

**4** **Poisson distribution; the on/off problem**

**5** **Assigning priors**

# Supplemental Topics

# Likelihood-Based Parametric Bootstrapping

Likelihood $\mathcal{L}(\theta) \equiv p(D_{\mathsf{obs}}|\theta)$.
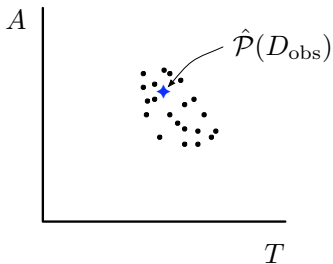Log-likelihood $L(\theta) = \ln \mathcal{L}(\theta)$.

For the Gaussian example,

$$
\begin{aligned}
\mathcal{L}(\mu) &= \prod_i \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{(x_i - \mu)^2}{2\sigma^2}\right] \\
&\propto \prod_i \exp\left[-\frac{(x_i - \mu)^2}{2\sigma^2}\right] \\
L(\mu) &= -\frac{1}{2}\sum_i \frac{(x_i - \mu)^2}{\sigma^2} + \mathsf{Const} \\
&= -\frac{\chi^2(\mu)}{2} + \mathsf{Const}
\end{aligned}
$$

# Incorrect Parametric Bootstrapping

$$\mathcal{P} = (A, T)$$



Histograms/contours of best-fit estimates from $D \sim p(D|\hat{\theta}(D_{\mathrm{obs}}))$ provide *poor* confidence regions—no better (possibly worse) than using a least-squares/$\chi^2$ covariance matrix.

What's wrong with the population of $\hat{\theta}$ points for this purpose?

# Incorrect Parametric Bootstrapping
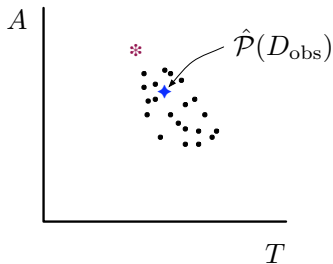
$$\mathcal{P} = (A, T)$$



Histograms/contours of best-fit estimates from $D \sim p(D|\hat{\theta}(D_{\mathrm{obs}}))$ provide *poor* confidence regions—no better (possibly worse) than using a least-squares/$\chi^2$ covariance matrix.

What's wrong with the population of $\hat{\theta}$ points for this purpose?

The estimates are skewed down and to the right, indicating the truth must be up and to the left. *Do not mistake variability of the estimator with the uncertainty of the estimate!*

Key idea: Use likelihood *ratios* to define confidence regions.
I.e., use $L$ or $\chi^2$ *differences* to define regions.

Estimate parameter values via *maximum likelihood* (min $\chi^2$)
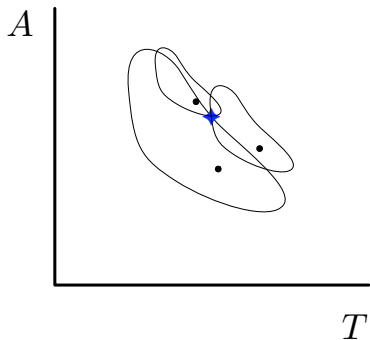$\rightarrow L_{\mathrm{max}}$.
Pick a constant $\Delta L$. Then

$$\Delta(D) = \{\theta : L(\theta) > L_{\mathrm{max}} - \Delta L\}$$
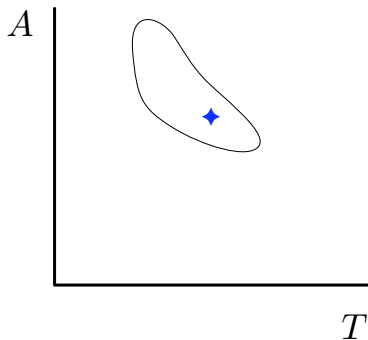
Coverage calculation:

1. Fix $\theta_0 = \hat{\theta}(D_{\mathrm{obs}})$ (plug-in approx'n)
2. Simulate a dataset from $p(D|\theta_0) \rightarrow L_D(\theta)$
3. Find maximum likelihood estimate $\hat{\theta}(D)$
4. Calculate $\Delta L = L_D(\hat{\theta}_D) - L_D(\theta_0)$
5. Goto (2) for $N$ total iterations
6. Histogram the $\Delta L$ values to find coverage vs. $\Delta L$
   (fraction of sim'ns with smaller $\Delta L$)

Report $\Delta(D_{\mathrm{obs}})$ with $\Delta L$ chosen for desired approximate CL.

Δ*L* Calibration        Reported Region

The CL is approximate due to:

- Monte Carlo error in calibrating Δ*L*
- The plug-in approximation

# Credible Region Via Posterior Sampling

Monte Carlo algorithm for finding credible regions:

1. Create a RNG that can sample $\theta$ from $p(\theta|D_{\text{obs}})$
2. Draw $N$ samples; record $\theta_i$ and $q_i = \pi(\theta_i)\mathcal{L}(\mu_i)$
3. Sort the samples by the $q_i$ values
4. An HPD region of probability $P$ is the $\theta$ region spanned by the $100P\%$ of samples with highest $q_i$

Note that no dataset other than $D_{\text{obs}}$ is ever considered.
$P$ is a property of the *particular interval* reported.

This complication is the rule rather than the exception!

Simple example: Estimate the mean *and standard deviation* of a normal distribution ($\mu = 5$, $\sigma = 1$, $N = 5$; 200 samples):

# Supplemental Topics

**1** Parametric bootstrapping vs. posterior sampling

**2** Estimation and model comparison for binary outcomes

**3** Basic inference with normal errors

**4** Poisson distribution; the on/off problem

**5** Assigning priors

# Binary Outcomes:
# Parameter Estimation

$M =$ Existence of two outcomes, $S$ and $F$; for each case or trial, the probability for $S$ is $\alpha$; for $F$ it is $(1 - \alpha)$

$H_i =$ Statements about $\alpha$, the probability for success on the next trial $\rightarrow$ seek $p(\alpha|D, M)$

$D =$ Sequence of results from $N$ observed trials:

FFSSSSFSSSFS ($n = 8$ successes in $N = 12$ trials)

*Likelihood:*

$$
\begin{aligned}
p(D|\alpha, M) &= p(\text{failure}|\alpha, M) \times p(\text{failure}|\alpha, M) \times \cdots \\
&= \alpha^n (1 - \alpha)^{N-n} \\
&= \mathcal{L}(\alpha)
\end{aligned}
$$

*Prior*

Starting with no information about $\alpha$ beyond its definition, use as an "uninformative" prior $p(\alpha|M) = 1$. Justifications:

- Intuition: Don't prefer any $\alpha$ interval to any other of same size
- Bayes's justification: "Ignorance" means that before doing the $N$ trials, we have no preference for how many will be successes:

$$P(n\,\text{success}|M) = \frac{1}{N+1} \qquad \rightarrow \qquad p(\alpha|M) = 1$$

Consider this a *convention*—an assumption added to $M$ to make the problem well posed.

*Prior Predictive*

$$
\begin{aligned}
p(D|M) &= \int d\alpha \; \alpha^n (1-\alpha)^{N-n} \\
&= B(n+1, N-n+1) = \frac{n!(N-n)!}{(N+1)!}
\end{aligned}
$$

A *Beta integral*, $B(a,b) \equiv \int dx \, x^{a-1}(1-x)^{b-1} = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$.

## Posterior

$$p(\alpha|D, M) \quad = \quad \frac{(N+1)!}{n!(N-n)!}\alpha^n(1-\alpha)^{N-n}$$

A *Beta distribution*. Summaries:

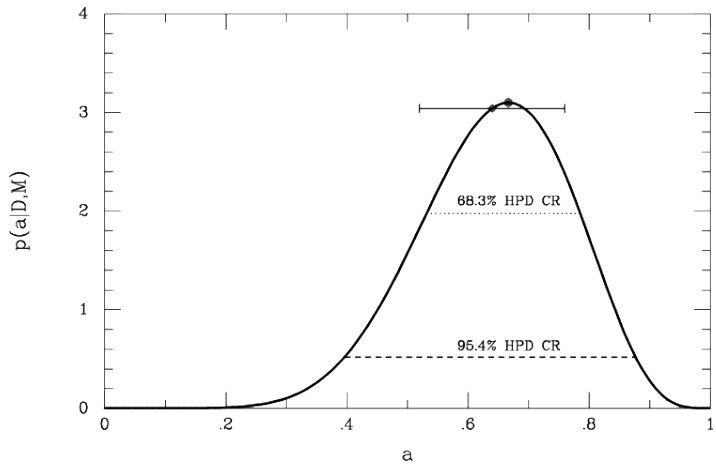- Best-fit: $\hat{\alpha} = \frac{n}{N} = 2/3$; $\langle\alpha\rangle = \frac{n+1}{N+2} \approx 0.64$

- Uncertainty: $\sigma_\alpha = \sqrt{\frac{(n+1)(N-n+1)}{(N+2)^2(N+3)}} \approx 0.12$
  Find credible regions numerically, or with incomplete beta function

Note that the posterior depends on the data only through $n$, not the $N$ binary numbers describing the sequence.

$n$ is a (minimal) *sufficient statistic*.

# Binary Outcomes: Model Comparison
## *Equal Probabilities?*

$M_1$: $\alpha = 1/2$
$M_2$: $\alpha \in [0, 1]$ with flat prior.

*Maximum Likelihoods*

$$M_1 : \qquad p(D|M_1) = \frac{1}{2^N} = 2.44 \times 10^{-4}$$

$$M_2 : \quad \mathcal{L}(\hat{\alpha}) \;=\; \left(\frac{2}{3}\right)^n \left(\frac{1}{3}\right)^{N-n} = 4.82 \times 10^{-4}$$

$$\frac{p(D|M_1)}{p(D|\hat{\alpha}, M_2)} = 0.51$$

Maximum likelihoods favor $M_2$ (failures more probable).

*Bayes Factor (ratio of model likelihoods)*

$$p(D|M_1) = \frac{1}{2^N}; \qquad \text{and} \qquad p(D|M_2) = \frac{n!(N-n)!}{(N+1)!}$$

$$\rightarrow B_{12} \equiv \frac{p(D|M_1)}{p(D|M_2)} = \frac{(N+1)!}{n!(N-n)!2^N}$$
$$= 1.57$$

Bayes factor (odds) favors $M_1$ (equiprobable).

Note that for $n = 6$, $B_{12} = 2.93$; for this small amount of data, we can never be very sure results are equiprobable.

If $n = 0$, $B_{12} \approx 1/315$; if $n = 2$, $B_{12} \approx 1/4.8$; for extreme data, 12 flips *can* be enough to lead us to strongly suspect outcomes have different probabilities.

(Frequentist significance tests can reject null for any sample size.)

# Binary Outcomes: Binomial Distribution

Suppose $D = n$ (number of heads in $N$ trials), rather than the actual sequence. What is $p(\alpha|n, M)$?

*Likelihood*

Let $\mathcal{S}$ = a sequence of flips with $n$ heads.

$$
\begin{aligned}
p(n|\alpha, M) &= \sum_{\mathcal{S}} \underbrace{p(S|\alpha, M)}_{} \; \underbrace{p(n|\mathcal{S}, \alpha, M)}_{} \\
&= \alpha^n(1-\alpha)^{N-n} C_{n,N}
\end{aligned}
$$

$\overset{\frown}{} \alpha^n (1-\alpha)^{N-n}$

$\llbracket \# \text{ successes} = n \rrbracket$

$C_{n,N} = \#$ of sequences of length $N$ with $n$ heads.

$$
\rightarrow p(n|\alpha, M) = \frac{N!}{n!(N-n)!}\alpha^n(1-\alpha)^{N-n}
$$

The *binomial distribution* for $n$ given $\alpha$, $N$.

*Posterior*

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

$$
\begin{aligned}
p(n|M) &= \frac{N!}{n!(N-n)!} \int d\alpha\; \alpha^n(1-\alpha)^{N-n} \\
&= \frac{1}{N+1}
\end{aligned}
$$

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

*Same result* as when data specified the actual sequence.

# Another Variation: Negative Binomial

Suppose $D = N$, the number of trials it took to obtain a predifined number of successes, $n = 8$. What is $p(\alpha|N, M)$?

*Likelihood*

$p(N|\alpha, M)$ is probability for $n - 1$ successes in $N - 1$ trials, times probability that the final trial is a success:

$$
\begin{aligned}
p(N|\alpha, M) &= \frac{(N-1)!}{(n-1)!(N-n)!}\alpha^{n-1}(1-\alpha)^{N-n}\alpha \\
&= \frac{(N-1)!}{(n-1)!(N-n)!}\alpha^{n}(1-\alpha)^{N-n}
\end{aligned}
$$

The *negative binomial distribution* for $N$ given $\alpha$, $n$.

*Posterior*

$$p(\alpha|D, M) = C'_{n,N} \frac{\alpha^n (1-\alpha)^{N-n}}{p(D|M)}$$

$$p(D|M) = C'_{n,N} \int d\alpha \, \alpha^n (1-\alpha)^{N-n}$$

$$\rightarrow p(\alpha|D, M) = \frac{(N+1)!}{n!(N-n)!} \alpha^n (1-\alpha)^{N-n}$$

*Same result* as other cases.

# Final Variation: "Meteorological Stopping"

Suppose $D = (N, n)$, the number of samples and number of successes in an observing run whose total number was determined by the weather at the telescope. What is $p(\alpha|D, M')$?

($M'$ adds info about weather to $M$.)

## Likelihood

$p(D|\alpha, M')$ is the binomial distribution times the probability that the weather allowed $N$ samples, $W(N)$:

$$p(D|\alpha, M') = W(N)\frac{N!}{n!(N-n)!}\alpha^n(1-\alpha)^{N-n}$$

Let $C_{n,N} = W(N)\binom{N}{n}$. We get the *same result* as before!

# Likelihood Principle

To define $\mathcal{L}(H_i) = p(D_{\text{obs}}|H_i, I)$, we must contemplate what other data we might have obtained. But the "real" sample space may be determined by many complicated, seemingly irrelevant factors; it may not be well-specified at all. Should this concern us?

*Likelihood principle*: The result of inferences depends only on how $p(D_{\text{obs}}|H_i, I)$ varies w.r.t. hypotheses. We can ignore aspects of the observing/sampling procedure that do not affect this dependence.

This happens because no sums of probabilities for hypothetical data appear in Bayesian results; Bayesian calculations *condition on $D_{obs}$*.

This is a sensible property that frequentist methods do not share. Frequentist probabilities are "long run" rates of performance, and depend on details of the sample space that are irrelevant in a Bayesian calculation.

# Goodness-of-fit Violates the Likelihood Principle

*Theory ($H_0$)*

> The number of "A" stars in a cluster should be 0.1 of the total.

*Observations*

> 5 A stars found out of 96 total stars observed.

*Theorist's analysis*

> Calculate $\chi^2$ using $\bar{n}_A = 9.6$ and $\bar{n}_X = 86.4$.
>
> Significance level is $p(> \chi^2 | H_0) = 0.12$ (or 0.07 using more rigorous binomial tail area). Theory is accepted.

### Observer's analysis

Actual observing plan was to keep observing until 5 A stars seen!

"Random" quantity is $N_{\mathrm{tot}}$, not $n_A$; it should follow the negative binomial dist'n. Expect $N_{\mathrm{tot}} = 50 \pm 21$.

$p(> \chi^2 | H_0) = 0.03$. Theory is *rejected*.

### Telescope technician's analysis

A storm was coming in, so the observations would have ended whether 5 A stars had been seen or not. The proper ensemble should take into account $p(\mathrm{storm})$ ...

### Bayesian analysis

The Bayes factor is the same for binomial or negative binomial likelihoods, and slightly favors $H_0$. Include $p(\mathrm{storm})$ if you want—it will drop out!

# Probability & frequency

Frequencies are relevant when modeling repeated trials, or repeated sampling from a population or ensemble.

*Frequencies are* observables

- When available, can be used to *infer* probabilities for next trial
- When unavailable, can be *predicted*

*Bayesian/Frequentist relationships*

- Relationships between probability and frequency
- Long-run performance of Bayesian procedures

# Probability & frequency in IID settings

*Frequency from probability*

Bernoulli's law of large numbers: In repeated i.i.d. trials, given $P(\text{success}|\ldots) = \alpha$, predict

$$\frac{n_{\text{success}}}{N_{\text{total}}} \to \alpha \quad \text{as} \quad N_{\text{total}} \to \infty$$

If $p(x)$ does not change from sample to sample, it may be interpreted as a frequency distribution.

*Probability from frequency*

Bayes's "An Essay Towards Solving a Problem in the Doctrine of Chances" $\to$ First use of Bayes's theorem:

Probability for success in next trial of i.i.d. sequence:

$$\mathsf{E}(\alpha) \to \frac{n_{\text{success}}}{N_{\text{total}}} \quad \text{as} \quad N_{\text{total}} \to \infty$$

If $p(x)$ does not change from sample to sample, it may be estimated from a frequency distribution.

# Supplemental Topics

1. Parametric bootstrapping vs. posterior sampling

2. Estimation and model comparison for binary outcomes

3. **Basic inference with normal errors**

4. Poisson distribution; the on/off problem

5. Assigning priors

# Inference With Normals/Gaussians

*Gaussian PDF*

$$p(x|\mu, \sigma) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad \text{over } [-\infty, \infty]$$

Common abbreviated notation: $x \sim N(\mu, \sigma^2)$

*Parameters*

$$\mu = \langle x \rangle \equiv \int dx \, x \, p(x|\mu, \sigma)$$

$$\sigma^2 = \langle (x-\mu)^2 \rangle \equiv \int dx \, (x-\mu)^2 \, p(x|\mu, \sigma)$$

## Gauss's Observation: Sufficiency

Suppose our data consist of $N$ measurements, $d_i = \mu + \epsilon_i$.
Suppose the noise contributions are independent, and
$\epsilon_i \sim \mathcal{N}(0, \sigma^2)$.

$$
\begin{aligned}
p(D|\mu, \sigma, M) &= \prod_i p(d_i|\mu, \sigma, M) \\
&= \prod_i p(\epsilon_i = d_i - \mu|\mu, \sigma, M) \\
&= \prod_i \frac{1}{\sigma\sqrt{2\pi}} \exp\left[-\frac{(d_i - \mu)^2}{2\sigma^2}\right] \\
&= \frac{1}{\sigma^N (2\pi)^{N/2}} e^{-Q(\mu)/2\sigma^2}
\end{aligned}
$$

Find dependence of $Q$ on $\mu$ by completing the square:

$$
\begin{aligned}
Q &= \sum_i (d_i - \mu)^2 \qquad\qquad [\text{Note: } Q/\sigma^2 = \chi^2(\mu)] \\
&= \sum_i d_i^2 + \sum_i \mu^2 - 2\sum_i d_i \mu \\
&= \left(\sum_i d_i^2\right) + N\mu^2 - 2N\mu\overline{d} \qquad \text{where } \overline{d} \equiv \frac{1}{N}\sum_i d_i \\
&= N(\mu - \overline{d})^2 + \left(\sum_i d_i^2\right) - N\overline{d}^2 \\
&= N(\mu - \overline{d})^2 + Nr^2 \quad \text{where } r^2 \equiv \frac{1}{N}\sum_i (d_i - \overline{d})^2
\end{aligned}
$$

Likelihood depends on $\{d_i\}$ only through $\overline{d}$ and $r$:

$$\mathcal{L}(\mu, \sigma) = \frac{1}{\sigma^N (2\pi)^{N/2}} \exp\left(-\frac{Nr^2}{2\sigma^2}\right) \exp\left(-\frac{N(\mu - \overline{d})^2}{2\sigma^2}\right)$$

The sample mean and variance are *sufficient statistics*.

This is a miraculous compression of information—the normal dist'n is highly *abnormal* in this respect!

# Estimating a Normal Mean

*Problem specification*

    Model: $d_i = \mu + \epsilon_i$, $\epsilon_i \sim N(0, \sigma^2)$, $\sigma$ is known $\rightarrow I = (\sigma, M)$.

    Parameter space: $\mu$; seek $p(\mu | D, \sigma, M)$

*Likelihood*

$$
\begin{aligned}
p(D | \mu, \sigma, M) &= \frac{1}{\sigma^N (2\pi)^{N/2}} \exp\left(-\frac{N r^2}{2\sigma^2}\right) \exp\left(-\frac{N(\mu - \overline{d})^2}{2\sigma^2}\right) \\
&\propto \exp\left(-\frac{N(\mu - \overline{d})^2}{2\sigma^2}\right)
\end{aligned}
$$

### *"Uninformative" prior*

Translation invariance $\Rightarrow p(\mu) \propto C$, a constant.

This prior is *improper* unless bounded.

### *Prior predictive/normalization*

$$
\begin{aligned}
p(D|\sigma, M) &= \int d\mu \; C \exp\left(-\frac{N(\mu - \overline{d})^2}{2\sigma^2}\right) \\
&= C(\sigma/\sqrt{N})\sqrt{2\pi}
\end{aligned}
$$

... minus a tiny bit from tails, using a proper prior.

*Posterior*

$$p(\mu|D, \sigma, M) = \frac{1}{(\sigma/\sqrt{N})\sqrt{2\pi}} \exp\left(-\frac{N(\mu - \overline{d})^2}{2\sigma^2}\right)$$

Posterior is $N(\overline{d}, w^2)$, with standard deviation $w = \sigma/\sqrt{N}$.

68.3% HPD credible region for $\mu$ is $\overline{d} \pm \sigma/\sqrt{N}$.

Note that $C$ drops out $\rightarrow$ limit of infinite prior range is well behaved.

### Informative Conjugate Prior

Use a normal prior, $\mu \sim N(\mu_0, w_0^2)$.

*Conjugate* because the posterior turns out also to be normal.

### Posterior

Normal $N(\tilde{\mu}, \tilde{w}^2)$, but mean, std. deviation *"shrink"* towards prior.

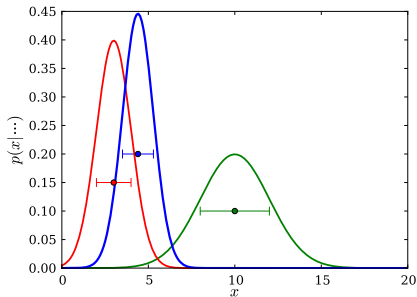Define $B = \frac{w^2}{w^2 + w_0^2}$, so $B < 1$ and $B = 0$ when $w_0$ is large. Then

$$
\begin{aligned}
\widetilde{\mu} &= \overline{d} + B \cdot (\mu_0 - \overline{d}) \\
\widetilde{w} &= w \cdot \sqrt{1 - B}
\end{aligned}
$$

*"Principle of stable estimation"* — The prior affects estimates only when data are not informative relative to prior.

Conjugate normal examples:

- Data have $\overline{d} = 3$, $\sigma/\sqrt{N} = 1$

- Priors at $\mu_0 = 10$, with $w = \{5, 2\}$

# Estimating a Normal Mean: Unknown $\sigma$

*Problem specification*

　Model: $d_i = \mu + \epsilon_i$, $\epsilon_i \sim N(0, \sigma^2)$, $\sigma$ is *unknown*

　Parameter space: $(\mu, \sigma)$; seek $p(\mu|D, M)$

*Likelihood*

$$
\begin{aligned}
p(D|\mu, \sigma, M) &= \frac{1}{\sigma^N (2\pi)^{N/2}} \exp\left(-\frac{Nr^2}{2\sigma^2}\right) \exp\left(-\frac{N(\mu - \overline{d})^2}{2\sigma^2}\right) \\
&\propto \frac{1}{\sigma^N} e^{-Q/2\sigma^2} \\
&\text{where} \quad Q = N\left[r^2 + (\mu - \overline{d})^2\right]
\end{aligned}
$$

*Uninformative Priors*

Assume priors for $\mu$ and $\sigma$ are independent.

Translation invariance $\Rightarrow p(\mu) \propto C$, a constant.

Scale invariance $\Rightarrow p(\sigma) \propto 1/\sigma$ (flat in $\log \sigma$).

*Joint Posterior for $\mu$, $\sigma$*

$$p(\mu, \sigma | D, M) \propto \frac{1}{\sigma^{N+1}} e^{-Q(\mu)/2\sigma^2}$$

*Marginal Posterior*

$$p(\mu|D, M) \propto \int d\sigma \; \frac{1}{\sigma^{N+1}} e^{-Q/2\sigma^2}$$

Let $\tau = \frac{Q}{2\sigma^2}$ so $\sigma = \sqrt{\frac{Q}{2\tau}}$ and $|d\sigma| = \tau^{-3/2}\sqrt{\frac{Q}{2}}\, d\tau$

$$\Rightarrow p(\mu|D, M) \;\propto\; 2^{N/2} Q^{-N/2} \int d\tau \; \tau^{\frac{N}{2}-1} e^{-\tau}$$
$$\propto\; Q^{-N/2}$$

Write $Q = Nr^2 \left[ 1 + \left( \frac{\mu - \overline{d}}{r} \right)^2 \right]$ and normalize:

$$p(\mu|D, M) = \frac{(\frac{N}{2} - 1)!}{(\frac{N}{2} - \frac{3}{2})! \sqrt{\pi}} \frac{1}{r} \left[ 1 + \frac{1}{N} \left( \frac{\mu - \overline{d}}{r/\sqrt{N}} \right)^2 \right]^{-N/2}$$

"Student's $t$ distribution," with $t = \frac{(\mu - \overline{d})}{r/\sqrt{N}}$

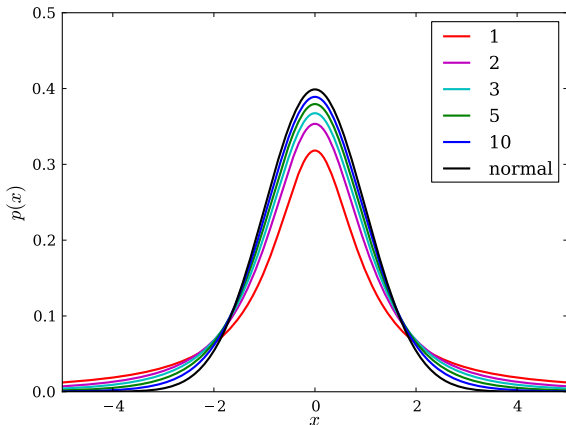A "bell curve," but with power-law tails

Large $N$:

$$p(\mu|D, M) \sim e^{-N(\mu - \overline{d})^2/2r^2}$$

This is the rigorous way to "adjust $\sigma$ so $\chi^2/\text{dof} = 1$."

It doesn't just plug in a best $\sigma$; it slightly broadens posterior to account for $\sigma$ uncertainty.

Student $t$ examples:

- $p(x) \propto \dfrac{1}{\left(1+\frac{x^2}{n}\right)^{\frac{n+1}{2}}}$

- Location $= 0$, scale $= 1$

- Degrees of freedom $= \{1, 2, 3, 5, 10, \infty\}$

# Gaussian Background Subtraction

Measure background rate $b = \hat{b} \pm \sigma_b$ with source off.

Measure total rate $r = \hat{r} \pm \sigma_r$ with source on.

Infer signal source strength $s$, where $r = s + b$.

With flat priors,

$$p(s, b | D, M) \propto \exp\left[-\frac{(b - \hat{b})^2}{2\sigma_b^2}\right] \times \exp\left[-\frac{(s + b - \hat{r})^2}{2\sigma_r^2}\right]$$

Marginalize $b$ to summarize the results for $s$ (complete the square to isolate $b$ dependence; then do a simple Gaussian integral over $b$):

$$p(s|D, M) \propto \exp\left[-\frac{(s - \hat{s})^2}{2\sigma_s^2}\right] \qquad \begin{aligned} \hat{s} &= \hat{r} - \hat{b} \\ \sigma_s^2 &= \sigma_r^2 + \sigma_b^2 \end{aligned}$$

$\Rightarrow$ Background *subtraction* is a special case of background *marginalization*; i.e., marginalization "told us" to subtract a background estimate.

Recall the standard derivation of background uncertainty via "propagation of errors" based on Taylor expansion (statistician's *Delta-method*).

Marginalization provides a generalization of error propagation—without approximation!

# Supplemental Topics

1. Parametric bootstrapping vs. posterior sampling

2. Estimation and model comparison for binary outcomes

3. Basic inference with normal errors

4. **Poisson distribution; the on/off problem**

5. Assigning priors

# Poisson Dist'n: Infer a Rate from Counts

*Problem:*

Observe $n$ counts in $T$; infer rate, $r$

*Likelihood*

$$\mathcal{L}(r) \equiv p(n|r, M) = p(n|r, M) = \frac{(rT)^n}{n!} e^{-rT}$$

*Prior*

Two simple standard choices (or conjugate gamma dist'n):

- $r$ known to be nonzero; it is a scale parameter:
$$p(r|M) = \frac{1}{\ln(r_u/r_l)} \frac{1}{r}$$

- $r$ may vanish; require $p(n|M) \sim$ Const:
$$p(r|M) = \frac{1}{r_u}$$

## Prior predictive

$$
\begin{aligned}
p(n|M) &= \frac{1}{r_u}\frac{1}{n!}\int_0^{r_u} dr (rT)^n e^{-rT} \\
&= \frac{1}{r_u T}\frac{1}{n!}\int_0^{r_u T} d(rT)(rT)^n e^{-rT} \\
&\approx \frac{1}{r_u T} \quad \text{for} \quad r_u \gg \frac{n}{T}
\end{aligned}
$$

## Posterior

A gamma distribution:

$$
p(r|n, M) = \frac{T(rT)^n}{n!} e^{-rT}
$$

# Gamma Distributions

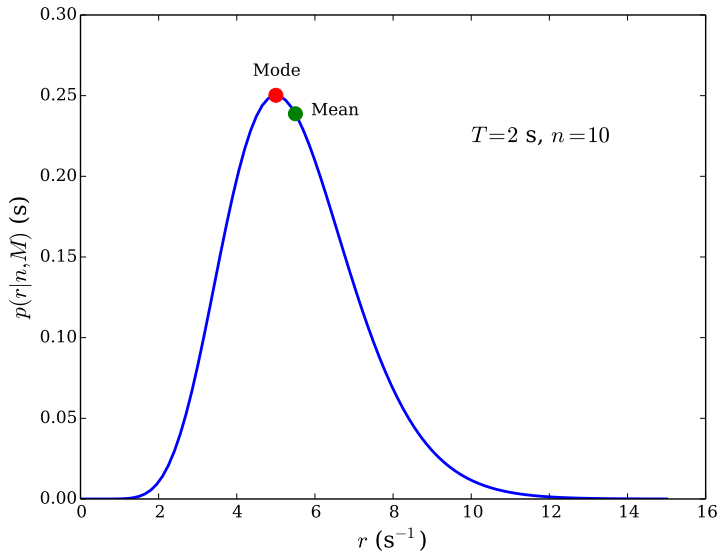A 2-parameter family of distributions over nonnegative $x$, with shape parameter $\alpha$ and scale parameter $s$:

$$p_\Gamma(x|\alpha, s) = \frac{1}{s\Gamma(\alpha)} \left(\frac{x}{s}\right)^{\alpha-1} e^{-x/s}$$

Moments:

$$\mathsf{E}(x) = s\alpha \qquad \mathsf{Var}(x) = s^2\alpha$$

Our posterior corresponds to $\alpha = n + 1$, $s = 1/T$.

- Mode $\hat{r} = \frac{n}{T}$; mean $\langle r \rangle = \frac{n+1}{T}$ (shift down 1 with $1/r$ prior)
- Std. dev'n $\sigma_r = \frac{\sqrt{n+1}}{T}$; credible regions found by integrating (can use incomplete gamma function)

## Conjugate prior

Note that a gamma distribution prior is the conjugate prior for the Poisson sampling distribution:

$$
\begin{aligned}
p(r|n, M') &\propto \mathrm{Gamma}(r|\alpha, s) \times \mathrm{Pois}(n|rT) \\
&\propto r^{\alpha-1}e^{-r/s} \times r^n e^{-rT} \\
&\propto r^{\alpha+n-1}\exp[-r(T + 1/s)]
\end{aligned}
$$

For $\alpha = 1$, $s \to \infty$, the gamma prior becomes an "uninformative" flat prior; posterior is proper even for $n = 0$

## Useful conventions

- Use a flat prior for a rate that may be zero

- Use a log-flat prior ($\propto 1/r$) for a nonzero scale parameter

- Use proper (normalized, bounded) priors

- Plot posterior with abscissa that makes prior flat (use $\log r$ abscissa for scale parameter case)

# Infer a Signal in a Known Background

*Problem:*

As before, but $r = s + b$ with $b$ known; infer $s$

$$p(s|n, b, M) = C \frac{T\left[(s+b)T\right]^n}{n!} e^{-(s+b)T}$$

$$
\begin{aligned}
C^{-1} &= \frac{e^{-bT}}{n!} \int_0^\infty d(sT)\,(s+b)^n T^n e^{-sT} \\
&= \sum_{i=0}^n \frac{(bT)^i}{i!} e^{-bT}
\end{aligned}
$$

A sum of Poisson probabilities for background events; it can be evaluated using the incomplete gamma function. (Helene 1983)

# The On/Off Problem

*Basic problem*

- Look off-source; unknown background rate $b$
  Count $N_{\mathrm{off}}$ photons in interval $T_{\mathrm{off}}$

- Look on-source; rate is $r = s + b$ with unknown signal $s$
  Count $N_{\mathrm{on}}$ photons in interval $T_{\mathrm{on}}$

- Infer $s$

*Conventional solution*

$$\hat{b} = N_{\mathrm{off}}/T_{\mathrm{off}}; \quad \sigma_b = \sqrt{N_{\mathrm{off}}}/T_{\mathrm{off}}$$
$$\hat{r} = N_{\mathrm{on}}/T_{\mathrm{on}}; \quad \sigma_r = \sqrt{N_{\mathrm{on}}}/T_{\mathrm{on}}$$
$$\hat{s} = \hat{r} - \hat{b}; \quad \sigma_s = \sqrt{\sigma_r^2 + \sigma_b^2}$$

But $\hat{s}$ can be *negative*!

## Bayesian Solution to On/Off Problem

First consider off-source data; use it to estimate $b$:

$$p(b|N_{\text{off}}, I_{\text{off}}) = \frac{T_{\text{off}}(bT_{\text{off}})^{N_{\text{off}}} e^{-bT_{\text{off}}}}{N_{\text{off}}!}$$

Use this as a prior for $b$ to analyze on-source data. For on-source analysis $I_{\text{all}} = (I_{\text{on}}, N_{\text{off}}, I_{\text{off}})$:

$$p(s, b|N_{\text{on}}) \quad \propto \quad p(s)p(b)[(s+b)T_{\text{on}}]^{N_{\text{on}}} e^{-(s+b)T_{\text{on}}} \qquad || \ I_{\text{all}}$$

$p(s|I_{\text{all}})$ is flat, but $p(b|I_{\text{all}}) = p(b|N_{\text{off}}, I_{\text{off}})$, so

$$p(s, b|N_{\text{on}}, I_{\text{all}}) \quad \propto \quad (s+b)^{N_{\text{on}}} b^{N_{\text{off}}} e^{-sT_{\text{on}}} e^{-b(T_{\text{on}}+T_{\text{off}})}$$

Now marginalize over $b$;

$$
\begin{aligned}
p(s|N_{\mathrm{on}}, I_{\mathrm{all}}) &= \int db \; p(s, b \mid N_{\mathrm{on}}, I_{\mathrm{all}}) \\
&\propto \int db \; (s+b)^{N_{\mathrm{on}}} b^{N_{\mathrm{off}}} e^{-sT_{\mathrm{on}}} e^{-b(T_{\mathrm{on}}+T_{\mathrm{off}})}
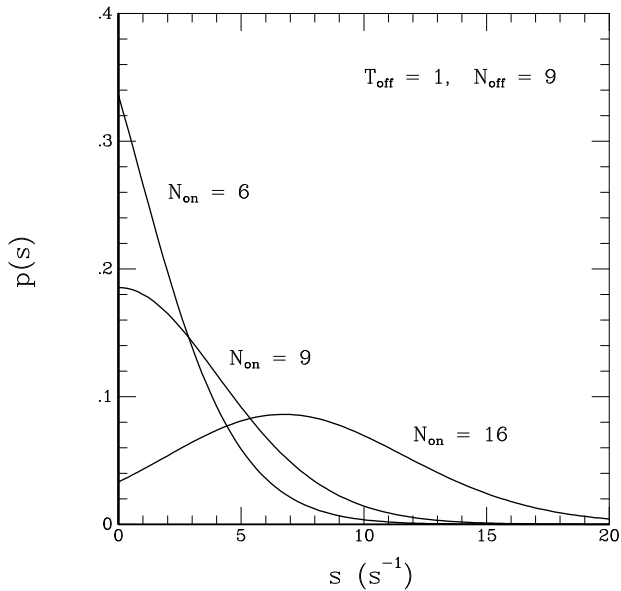\end{aligned}
$$

Expand $(s+b)^{N_{\mathrm{on}}}$ and do the resulting $\Gamma$ integrals:

$$
\begin{aligned}
p(s|N_{\mathrm{on}}, I_{\mathrm{all}}) &= \sum_{i=0}^{N_{\mathrm{on}}} C_i \frac{T_{\mathrm{on}}(sT_{\mathrm{on}})^i e^{-sT_{\mathrm{on}}}}{i!} \\
C_i &\propto \left(1 + \frac{T_{\mathrm{off}}}{T_{\mathrm{on}}}\right)^i \frac{(N_{\mathrm{on}}+N_{\mathrm{off}}-i)!}{(N_{\mathrm{on}}-i)!}
\end{aligned}
$$

Posterior is a weighted sum of Gamma distributions, each assigning a different number of on-source counts to the source. (Evaluate via recursive algorithm or confluent hypergeometric function.)
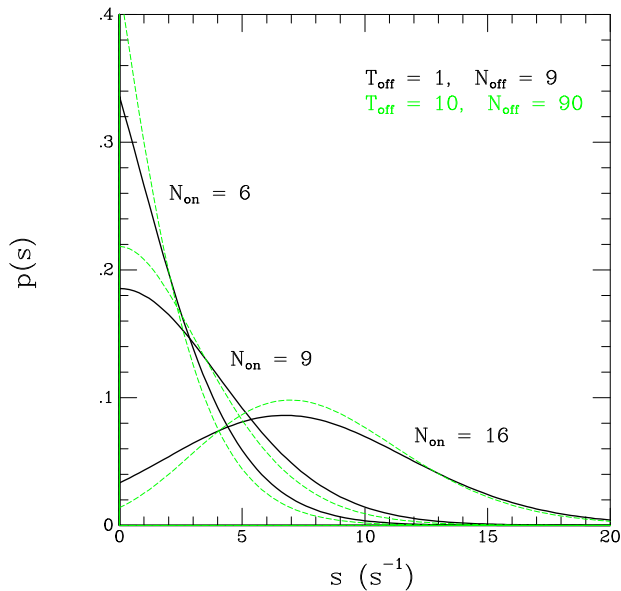
Example On/Off Posteriors—Short Integrations

Example On/Off Posteriors—Long Background Integrations

## Second Solution of the On/Off Problem

Consider all the data at once; the likelihood is a product of Poisson distributions for the on- and off-source counts:

$$
\begin{aligned}
\mathcal{L}(s, b) &\equiv p(N_{\mathrm{on}}, N_{\mathrm{off}} | s, b, I) \\
&\propto [(s + b) T_{\mathrm{on}}]^{N_{\mathrm{on}}} e^{-(s+b)T_{\mathrm{on}}} \times (b T_{\mathrm{off}})^{N_{\mathrm{off}}} e^{-b T_{\mathrm{off}}}
\end{aligned}
$$

Take joint prior to be flat; find the joint posterior and marginalize over $b$;

$$
\begin{aligned}
p(s | N_{\mathrm{on}}, I_{\mathrm{on}}) &= \int db \, p(s, b | I) \, \mathcal{L}(s, b) \\
&\propto \int db \, (s + b)^{N_{\mathrm{on}}} b^{N_{\mathrm{off}}} e^{-s T_{\mathrm{on}}} e^{-b(T_{\mathrm{on}} + T_{\mathrm{off}})}
\end{aligned}
$$

$\rightarrow$ same result as before.

## Third Solution: Data Augmentation

Suppose we knew the number of on-source counts that are from the background, $N_b$. Then the on-source likelihood is simple:

$$p(N_{\rm on}|s, N_b, I_{\rm all}) = {\rm Pois}(N_{\rm on} - N_b; sT_{\rm on}) = \frac{(sT_{\rm on})^{N_{\rm on} - N_b}}{(N_{\rm on} - N_b)!} e^{-sT_{\rm on}}$$

*Data augmentation:* Pretend you have the "missing data," then marginalize to account for its uncertainty:

$$
\begin{aligned}
p(N_{\rm on}|s, I_{\rm all}) &= \sum_{N_b=0}^{N_{\rm on}} p(N_b|I_{\rm all})\, p(N_{\rm on}|s, N_b, I_{\rm all}) \\
&= \sum_{N_b} \text{Predictive for } N_b \times {\rm Pois}(N_{\rm on} - N_b; sT_{\rm on})
\end{aligned}
$$

$$
\begin{aligned}
p(N_b|I_{\rm all}) &= \int db\; p(b|N_{\rm off}, I_{\rm off})\, p(N_b|b, I_{\rm on}) \\
&= \int db\; {\rm Gamma}(b) \times {\rm Pois}(N_b; bT_{\rm on})
\end{aligned}
$$

$\rightarrow$ same result as before.

# A profound consistency

We solved the on/off problem in multiple ways, always finding the same final results.

This reflects something fundamental about Bayesian inference.

R. T. Cox proposed two necessary conditions for a quantification of uncertainty:

- It should duplicate deductive logic when there is no uncertainty

- Different decompositions of arguments should produce the same final quantifications (internal consistency)

Great surprise: These conditions are *sufficient*; they lead to the probability axioms. E. T. Jaynes and others refined and simplified Cox's analysis.

# Multibin On/Off

The more typical on/off scenario:

Data = spectrum or image with counts in many bins

Model $M$ gives signal rate $s_k(\theta)$ in bin $k$, parameters $\theta$

To infer $\theta$, we need the likelihood:

$$\mathcal{L}(\theta) = \prod_k p(N_{\mathrm{on}\,k}, N_{\mathrm{off}\,k} | s_k(\theta), M)$$

For each $k$, we have an on/off problem as before, only we just need the marginal likelihood for $s_k$ (not the posterior). The same $C_i$ coefficients arise.

XSPEC and CIAO/Sherpa provide this as an option

van Dyk[+](2001) does the same thing via data augmentation (Monte Carlo)

# Supplemental Topics

**1** Parametric bootstrapping vs. posterior sampling

**2** Estimation and model comparison for binary outcomes

**3** Basic inference with normal errors

**4** Poisson distribution; the on/off problem

**5** Assigning priors

# Well-Posed Problems

The rules (BT, LTP, . . . ) express desired probabilities in terms of other probabilities

To get a numerical value *out*, at some point we have to put numerical values *in*

*Direct probabilities* are probabilities with numerical values determined directly by premises (via modeling assumptions, symmetry arguments, previous calculations, desperate presumption . . . )

An inference problem is *well posed* only if all the needed probabilities are assignable based on the context. We may need to add new assumptions as we see what needs to be assigned. We may not be entirely comfortable with what we need to assume! (Remember Euclid's fifth postulate!)

Should explore how results depend on uncomfortable assumptions ("robustness")

# Contextual/prior/background information

Bayes's theorem moves the data and hypothesis propositions wrt the solidus:

$$P(H_i|D_{\mathrm{obs}}, I) = P(H_i|I) \frac{P(D_{\mathrm{obs}}|H_i, I)}{P(D_{\mathrm{obs}}|I)}$$

It lets us *change the premises*

"Prior information" or "background information" or "context" = information that is **always** a premise (for the current calculation)

Notation: $P(\cdot|\cdot, I)$ or $P(\cdot|\cdot, \mathcal{C})$ or $P(\cdot|\cdot, M)$ or ...

The context can be a notational nuisance! "Skilling conditional":

$$P(H_i|D_{\mathrm{obs}}) = P(H_i) \frac{P(D_{\mathrm{obs}}|H_i)}{P(D_{\mathrm{obs}})} \qquad \| \, \mathcal{C}$$

# Essential contextual information

We can only be uncertain about *A* if there are alternatives; what they are will bear on our uncertainty. *We must explicitly specify relevant alternatives.*

**Hypothesis space:** The set of alternative hypotheses of interest (and auxiliary hypotheses needed to predict the data)

**Data/sample space:** The set of possible data we may have predicted before learning of the observed data

**Predictive model:** Information specifying the likelihood function (e.g., the conditional predictive dist'n/sampling dist'n)

**Other prior information:** Any further information available or necessary to assume to make the problem *well posed*

Bayesian literature often uses **model** to refer to *all* of the contextual information used to study a particular dataset and predictive model

# Directly assigned sampling distributions

Some examples of reasoning leading to sampling distributions:

- Binomial distribution:
    - ▶ Ansatz: Probability for a Bernoulli trial, $\alpha$
    - ▶ LTP $\Rightarrow$ binomial for $n$ successes in $N$ trials

- Poisson distribution:
    - ▶ Ansatz: $P(\text{event in } dt|\lambda) \propto \lambda dt$;
      probabilities for events in disjoint intervals independent
    - ▶ Product & sum rules $\Rightarrow$ Poisson for $n$ in $T$

- Gaussian distribution:
    - ▶ CLT: Probability theory for sum of many quantities with independent, finite-variance PDFs
    - ▶ Sufficiency (Gauss): Seek distribution with sample mean as sufficient statistic (also sample variance)
    - ▶ Asymptotic limits: large $n$ Binomial, Poisson
    - ▶ Others: Herschel's invariance argument (2-D), maximum entropy. . .

# Assigning priors

*Sources of prior information*

- Analysis of previous experimental/observational data (but begs the question of what prior to use for the first such analysis)

- *Subjective priors*: Elicit a prior from an expert in the problem domain, e.g., via ranges, moments, quantiles, histograms

- *Population priors*: When it's meaningful to pool observations, we potentially can *learn* a shared prior—multilevel modeling does this

*"Non-informative" priors*

- Seek a prior that in some sense (TBD!) expresses a lack of information prior to considering the data

- No universal solution—this notion must be problem-specific, e.g., exploiting symmetries

## Priors derived from the likelihood function

Few common problems beyond location/scale problems admit a transformation group argument $\rightarrow$ we need a more general approach to formal assignment of priors that express "ignorance" in some sense

There is no universal consensus on how to do this (yet? ever?)

A common underlying idea: The same $\mathcal{C}$ appears in the prior, $p(\theta|\mathcal{C})$, and the likelihood, $p(D|\theta, \mathcal{C})$—the prior "knows" about the likelihood function, although it doesn't know what data values will be plugged into it

**Jeffreys priors**: Uses Fisher information to define a (parameter-dependent) scale defining a prior; parameterization invariant, but strange behavior in many dimensions

**Reference priors**: Uses information theory to define a prior that (asymptotically) has the least effect on the posterior; complicated algorithm; gives good frequentist behavior to Bayesian inferences

# Jeffreys priors

Heuristic motivation:

- Dimensionally, $\pi(\theta) \propto 1/(\theta \text{ scale})$

- If we have data $D$, a natural inverse scale at $\theta$, from the likelihood function, is the square root of the *observed Fisher information* (recall Laplace approximation):

$$I_D(\theta) \equiv -\frac{\mathrm{d}^2 \log \mathcal{L}_D(\theta)}{\mathrm{d}\theta^2}$$

- For a prior, we don't know $D$; for each $\theta$, average over $D$ predicted by the sampling distribution; this defines the *(expected) Fisher information*:
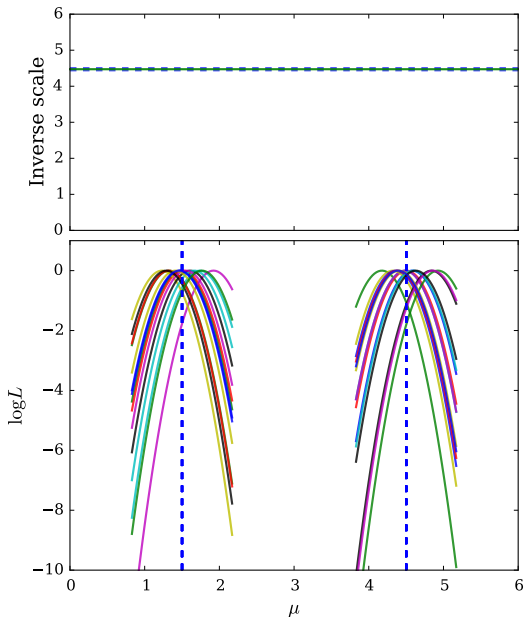
$$I(\theta) \equiv -\mathbb{E}_D \left[ \frac{\mathrm{d}^2 \log \mathcal{L}_D(\theta)}{\mathrm{d}\theta^2} \,\middle|\, \theta \right]$$

*Jeffreys' prior*:
$$\pi(\theta) \propto [I(\theta)]^{1/2}$$

- Note the proportionality—the prior scale depends on how much the likelihood function scale *changes* vs. $\theta$

- Puts more weight in regions of parameter space where the data are expected to be more informative

- Parameterization invariant, due to use of derivatives and vanishing expectation of the *score function*
  $$S_D(\theta) = \frac{d \log \mathcal{L}_D(\theta)}{d\theta}$$

- Typically improper when parameter space is non-compact

- Improves *frequentist* performance of posterior intervals w.r.t. intervals based on flat priors

- Only considered sound for a single parameter (or considering a single parameter at a time in some multiparameter problems)

# Jeffreys prior for normal mean



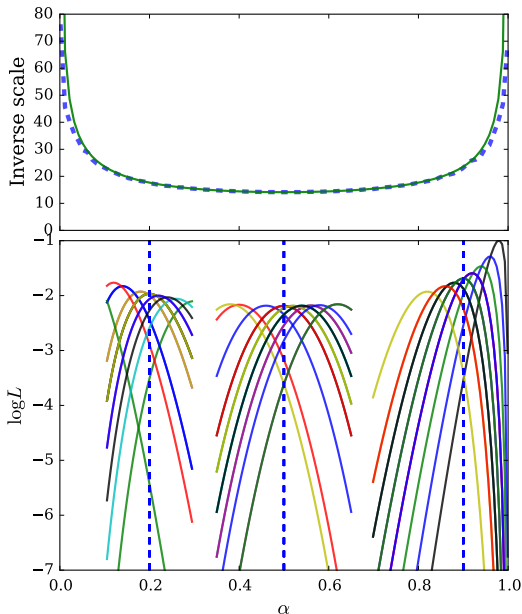$N = 20$ samples from normals with $\sigma = 1$

Likelihood width is independent of $\mu \Rightarrow$

$$\pi(\mu) = \text{Const}$$

Another justification of the uniform prior

Prior is improper without prior limits on the range

# Jeffreys prior for binomial probability



Binomial success counts $n$ from $N = 50$ trials

$$\pi(\mu) = \frac{1}{\pi \alpha^{1/2}(1-\alpha)^{1/2}}$$
$$= \text{Beta}(1/2, 1/2)$$

# Information Gain as Entropy Change

## *Entropy and uncertainty*

Shannon entropy = a scalar measure of the degree of uncertainty expressed by a probability distribution

$$
\begin{aligned}
\mathcal{S} &= \sum_i p_i \log \frac{1}{p_i} \qquad \text{``Average surprisal''} \\
&= -\sum_i p_i \log p_i
\end{aligned}
$$

## *Information gain*

Information gain upon learning $D$ = decrease in uncertainty:

$$
\begin{aligned}
\mathcal{I}(D) &= \mathcal{S}[\{p(H_i)\}] - \mathcal{S}[\{p(H_i|D)\}] \\
&= \sum_i p(H_i|D) \log p(H_i|D) - \sum_i p(H_i) \log p(H_i)
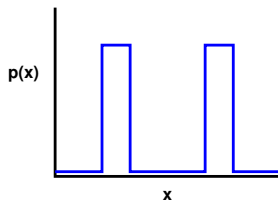\end{aligned}
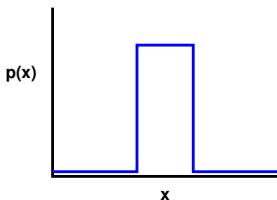$$

# A 'Bit' About Entropy

*Entropy of a Gaussian*

$$p(x) \propto e^{-(x-\mu)^2/2\sigma^2} \quad \rightarrow \quad \mathcal{I} \propto -\log(\sigma)$$

$$p(\vec{x}) \propto \exp\left[-\tfrac{1}{2}\vec{x} \cdot \mathbf{V}^{-1} \cdot \vec{x}\right] \quad \rightarrow \quad \mathcal{I} \propto -\log(\det \mathbf{V})$$

$\rightarrow$ Asymptotically like log Fisher matrix

*A log-measure of "volume" or "spread," not range*



These distributions have the same entropy/amount of information.

### Expected information gain

When the data are yet to be considered, the *expected* information gain averages over $D$; straightforward use of the product rule/Bayes's theorem gives:

$$\mathbb{E}\mathcal{I} = \int dD \, p(D) \, \mathcal{I}(D)$$
$$= \int dD \, p(D) \, \sum_i p(H_i|D) \log \left[ \frac{p(H_i|D)}{p(H_i)} \right]$$

For a continuous hypothesis space labeled by parameter(s) $\theta$,

$$\mathbb{E}\mathcal{I} = \int dD \, p(D) \, \int d\theta p(\theta|D) \log \left[ \frac{p(\theta|D)}{p(\theta)} \right]$$

This is the expectation value of the *Kullback-Leibler divergence* between the prior and posterior:

$$\mathcal{D} \equiv \int d\theta p(\theta|D) \log \left[ \frac{p(\theta|D)}{p(\theta)} \right]$$

# Reference priors

Bernardo (later joined by Berger & Sun) advocates *reference priors*, priors chosen to maximize the KLD between prior and posterior, as an "objective" expression of the idea of a "non-informative" prior: reference priors let the data most strongly dominate the prior (on average)

- Rigorous definition invokes asymptotics and delicate handling of non-compact parameter spaces to make sure posteriors are proper
- For 1-D problems, the reference prior is the Jeffreys prior
- In higher dimensions, the reference prior is *not* the Jeffreys prior; it behaves better
- The construction in higher dimensions is complicated and depends on separating interesting vs. nuisance parameters (but see Berger, Bernardo & Sun 2015, "Overall objective priors")
- Reference priors are typically improper on non-compact spaces
- They give Bayesian inferences good frequentist properties
- A constructive numerical algorithm exists