

**STSCI 4780/5780**  
**Assigning direct probabilities:**  
**Sampling distributions and priors**

Tom Lored, CCAPS & SDS, Cornell University

© 2022-04-26

# The rules

## AND/OR/NOT

$$\begin{aligned}\text{'AND' (product rule): } P(A \wedge B|\mathcal{P}) &= P(A|\mathcal{P}) P(B|A \wedge \mathcal{P}) \\ &= P(B|\mathcal{P}) P(A|B \wedge \mathcal{P})\end{aligned}$$

$$\begin{aligned}\text{'OR' (sum rule): } P(A \vee B|\mathcal{P}) &= P(A|\mathcal{P}) + P(B|\mathcal{P}) \\ &\quad - P(A \wedge B|\mathcal{P})\end{aligned}$$

$$\text{'NOT': } P(\bar{A}|\mathcal{P}) = 1 - P(A|\mathcal{P})$$

## Key theorems: BT, LTP

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

$$\begin{aligned}\sum_i P(A, B_i|\mathcal{C}) &= \sum_i P(B_i|A, \mathcal{C})P(A|\mathcal{C}) = P(A|\mathcal{C}) \\ &= \sum_i P(B_i|\mathcal{C})P(A|B_i, \mathcal{C})\end{aligned}$$

## Well-posed problems

The rules (BT, LTP, . . . ) express desired probabilities in terms of other probabilities—they comprise a kind of *grammar* for inference

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

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

An inference problem is *well posed* only if all the needed direct probabilities are assignable. 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 convenient/defeasible assumptions (“robustness”)

## Lec04 recap: Essential contextual information

We can only be uncertain about a proposition,  $A$ , if there are alternatives (at least  $\bar{A}$ !); 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, e.g., for LTP)

**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)—the connection between data and hypotheses

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

*Where do predictive models (sampling dist'ns) come from? Seek patterns & approaches that may inform how we assign priors.*

# Directly assigned sampling distributions

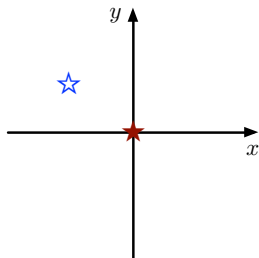
Some examples of reasoning leading to sampling distributions:

- Binomial distribution:
  - ▶ Ansatz: Probability for a Bernoulli trial,  $\alpha$
  - ▶ LTP  $\Rightarrow$  binomial for  $n$  successes in  $N$  trials via counting
- Poisson distribution:
  - ▶ Ansatz:  $P(\text{event in } dt|\lambda) \propto \lambda dt$ ; probabilities for events in disjoint intervals independent; no “clones”
  - ▶ Product & sum rules  $\Rightarrow$  Poisson for  $n$  in  $T$

- Gaussian distribution:

- ▶ Leading-order Taylor series for a smooth log PDF at its peak (cf. the Laplace approximation)
- ▶ 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. . .

# Herschel-Maxwell derivation of 2-D normal



- Assume for simplicity tht knowledge of  $x$  tells us nothing about  $y$ :

$$\rho(x, y) dx dy = f(x) dx \times h(y) dy$$

- Adopt same distribution in  $x$  and  $y$  (a kind of  $x, y$  equivalence/similarity):

$$\rho(x, y) dx dy = f(x) dx \times f(y) dy$$

- Express in polar coordinates,  $x = r \cos \theta$ ,  $y = r \sin \theta$ :

$$\rho(x, y) dx dy = f(r \cos \theta) f(r \sin \theta) r dr d\theta$$

- Require distribution to be independent of angle (rotational symmetry):

$$\begin{aligned} f(r \cos \theta) f(r \sin \theta) &= g(r) \\ \Rightarrow f(x) f(y) &= g\left(\sqrt{x^2 + y^2}\right) \end{aligned}$$

Solve this *functional equation*:

$$f(x)f(y) = g\left(\sqrt{x^2 + y^2}\right)$$

For  $y = 0$ ,  $f(x)f(0) = g(x)$ ; replacing  $g(\cdot)$  gives

$$f(x)f(y) = f\left(\sqrt{x^2 + y^2}\right) f(0)$$

$$\ln \left[ \frac{f(x)}{f(0)} \right] + \ln \left[ \frac{f(y)}{f(0)} \right] = \ln \left[ \frac{f\left(\sqrt{x^2 + y^2}\right)}{f(0)} \right]$$

Requires a function of  $x$  plus a function  $y$  that's a function only of  $x^2 + y^2$ :

$$\ln \left[ \frac{f(x)}{f(0)} \right] = ax^2$$



Normalization requires  $a < 0$  and determines the normalization constant  $f(0)$ :

$$f(x) = \sqrt{\frac{\alpha}{\pi}} e^{-\alpha x^2}, \quad \alpha > 0,$$
$$\rho(x, y) = \frac{\alpha}{\pi} e^{-\alpha(x^2 + y^2)}$$

Maxwell extended argument to 3-D (velocities)  $\Rightarrow$  Maxwellian velocity distribution in statistical mechanics

**Theme/pattern:**

Qualitative requirements

(typically from symmetry, invariance)

$\Rightarrow$  constraints on a PDF, in the form of functional equations

$\Rightarrow$  a specific family of dist'n's from solving the eqn's

# Assigning priors

## *Sources of prior information*

- *Discovery chains*: Analysis of previous experimental or 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 (more radical *subjective Bayesians* assert an agent's priors need only express an opinion that agrees with bets the agent is willing to make); this corresponds to  $\mathcal{C} \approx$  "What person  $X$  knows about the problem"
- *Population priors*: When it's meaningful to pool observations, we potentially can *learn* a shared prior—hierarchical/graphical/multilevel models do this (but then any unknown hyperparameters need priors)

## *“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
- *Objective Bayes (OBayes)*—Bayesian inference using priors that are largely or entirely determined from specification of the parameters and likelihood function in some mathematical/algorithmic way

Imagine a Stan++ package that uses the data, parameter, and model blocks to derive a prior when one isn't specified, purely from the problem specification—how might it work?

## Discrete uniform prior

Seek an *algorithm* that assigns a discrete PMF  $p_i = P(A_i|\mathcal{C})$  from symbolic expression of a problem

Consider the same problem, expressed in symbols in two different ways:

- Formulation 1,  $\mathcal{C}_1$  with a suite of  $N$  propositions denoted  $A_i$ , assigns

$$p_i = P(A_i|\mathcal{C}_1)$$

- Formulation 2,  $\mathcal{C}_2$  with the *same* suite of  $N$  propositions, but with different labels/symbols,  $B_i$ , assigns

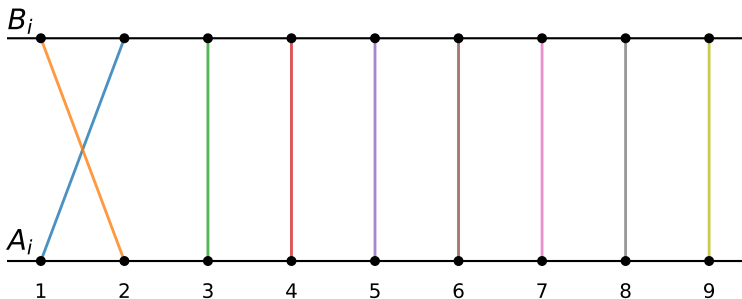
$$p'_i = P(B_i|\mathcal{C}_2)$$

These are two versions of the same problem—they differ in *form*, but not in *content*

Suppose *semantically* (in terms of meanings of symbols) we can identify equivalent propositions that happen to be represented by different symbols, e.g.,

$$B_1 \equiv A_2; \quad B_2 \equiv A_1$$

$$B_k \equiv A_k \text{ for } k = 3 \text{ to } N$$



*Transformation equations*: Equate probabilities that refer to equivalent propositions (reflecting *semantic* equivalence),

$$\Rightarrow p_1 = p'_2, p_2 = p'_1, p_k = p'_k \text{ for } k > 2$$

Suppose the contextual info doesn't distinguish among the  $N$  propositions. *Syntactically*, the two formulations will look equivalent (they have the same *form*; Stan++ would recognize them as differing only in choices of symbols—like a refactoring)

*Symmetry equations*: A rule that assigns a probability  $p_i$  to  $p(A_i|C_1)$  based solely on the pattern of symbols in the problem definition must assign the *same* value to  $p(B_i|C_2)$

$$P(A_i|C_1) = P(B_i|C_2) \\ \Rightarrow p_1 = p'_1, p_2 = p'_2, \dots$$

Reflects a *consistency requirement*—a formal rule should assign the same probabilities to problems with equivalent symbolic expressions (*syntactic* equivalence)

Combine transformation & symmetry:

$$p_1 = p_2$$

Consider further problems that differ from  $\mathcal{C}_1$  by different permutations of the labels for the propositions  $\Rightarrow$

$$p_i = p_j \text{ for all } i, j$$

Imposing normalization:

$$p_i = \frac{1}{N}$$

for a suite of  $N$  propositions with no information in  $\mathcal{C}$  distinguishing between them

Known as the *principle of insufficient reason* (adopted without name by Bernoulli, Laplace), aka *principle of indifference* (Keynes, Jaynes)

# Continuous uniform prior can't be universal

*“Method of inverse probability”*—Bayes's theorem with uniform/flat prior PDFs for *continuous* parameters (adopted for expedience)

Inverse probability was heavily used by Laplace and subsequently dominated statistical practice until late 19th/early 20th century

Criticism (Boole, Venn, others):

- Investigator 1 analyzes data using a model with parameter  $\theta$ , with prior PDF  $\pi_1(\theta) = C_1$
- Investigator 2 analyzes the *same* data using the *same* model, but parameterized in terms of  $\phi = \Phi(\theta)$ ; assigns the prior PDF  $\pi_2(\phi) = C_2$
- To an interval  $d\phi$  corresponding to models with  $\theta \in d\theta$ , investigator 2 assigns probability

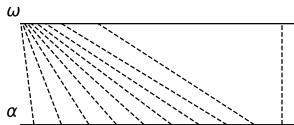
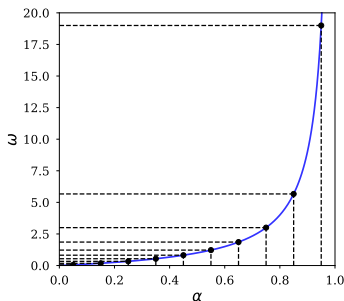
$$\pi_2(\phi)d\phi = C_2d\phi = C_2\Phi'(\theta)d\theta$$

This won't be constant wrt  $\theta$  in general  $\Rightarrow$  investigators adopting uniform priors with different parameterizations may disagree in inferences with the same model and data



Consider two parameterizations of a binomial sampling dist'n:

- Success probability,  $\alpha$ :  $p(n|\alpha) \propto \alpha^n(1 - \alpha)^{N-1}$
- Odds,  $\omega \equiv \alpha/(1 - \alpha)$ :  $p(n|\alpha) \propto \left(\frac{\omega}{1+\omega}\right)^n \left(\frac{1}{1+\omega}\right)^{N-1}$



Uniform over  $\alpha$  is inconsistent with uniform over  $\omega$ !

Is there a “natural” way to pick intervals we should consider equally probable a priori?

## Non-informative priors from transformation groups

The principle of indifference could not be justified in this case because the two versions of the problem have different symbolic structure (parameterizations)  $\rightarrow$  there is no obvious counterpart to the (syntactic) *symmetry equations*

*Sometimes* a Bayesian inference problem with continuous parameters *does* have a symmetry that identifies symmetry equations—particular transformations that make different parameterizations of the problem symbolically equivalent

The likelihood function plays a key role in identifying such symmetries; i.e., we aren't assigning a prior  $\pi(\theta)$  in a vacuum, but using information about  $\theta$ 's meaning encoded in the likelihood function

## *Poisson rate (for discrete data)*

$\mathcal{C}_1$ : Estimate a Poisson rate, using a parameter  $r$  in counts/sec, from  $n$  counts observed in time  $\delta$  sec, based on

$$p(n|r, \mathcal{C}_1) = \frac{(r\delta)^n}{n!} e^{-r\delta}, \quad p(r|\mathcal{C}_1) = f(r)$$

$\mathcal{C}_2$ : Estimate the same a Poisson rate, using a parameter  $R$  in counts/hr, from  $n$  counts observed in time  $\Delta$  hr, based on

$$p(n|R, \mathcal{C}_2) = \frac{(R\Delta)^n}{n!} e^{-R\Delta}, \quad p(R|\mathcal{C}_2) = g(R)$$

Since they are describing the same situation,  $R$  and  $r$ , and  $\Delta$  and  $\delta$  must be related; for  $\alpha \equiv 1/3600$ :

$$\Delta = \alpha\delta, \quad R = r/\alpha$$

which ensures that  $r\delta = R\Delta$ , so they assign the same probabilities to  $n$  in equivalent situations

Repeating:

$$\Delta = \alpha\delta, \quad R = r/\alpha$$

Since they are describing the same situation, the priors must be related via change-of-variables (transformation eq'ns):

$$g(R) dR = f(r) dr \quad \Rightarrow \quad g(R) = f(\alpha R) \times \alpha$$

This particular scale transformation makes the two formulations of the problem look equivalent (to Stan++, only the choice of variable names is different).

If we have no prior information distinguishing the formulations—we have no information favoring a particular time scale for the phenomenon—we should assign priors of the same form (symmetry eq'ns):

$$g(\cdot) = f(\cdot)$$

Repeating:

$$g(R) = f(\alpha R) \times \alpha; \quad g(\cdot) = f(\cdot)$$

$$\rightarrow f(R) = \alpha f(\alpha R)$$

This holds for any  $(R, \alpha)$ . Substitute  $R = 1$ :

$$f(1) = \alpha f(\alpha)$$

$$f(\alpha) = \frac{C}{\alpha}$$

For Poisson rate inference,  $\pi(r) \propto 1/r$  expresses ignorance of the time scale for the phenomenon

Note that this prior is *improper*—a generic feature of these types of priors when the parameter space is infinite

## *Location parameter (for continuous data)*

For a PDF  $p(x|\mu)$  for  $x$  with parameter  $\mu$ , if it is of the form  $p(x|\mu) = h(x - \mu)$ , then  $\mu$  is called a *location parameter*

$\mathcal{C}_1$ : Estimate a location parameter,  $\mu$ , using a sampling distribution for data  $x_i$  of the form

$$p(x|\mu, \mathcal{C}_1) = h(x - \mu);$$

Denote the  $\mathcal{C}_1$  prior by  $P(\mu \in d\mu | \mathcal{C}_1) = f(\mu) d\mu$

$\mathcal{C}_2$ : In a coordinate system shifted by  $\Delta$ , use data  $x'_i = x_i + \Delta$  to estimate  $\mu' = \mu + \Delta$ .

Since  $dx' = dx$ , the sampling distribution for  $x'$  has the same form,

$$p(x'|\mu', \mathcal{C}_2) = h(x' - \mu').$$

Denote the  $\mathcal{C}_2$  prior by  $P(\mu' \in d\mu' | \mathcal{C}_2) = g(\mu') d\mu'$

**Symmetry:** Note that  $h(x' - \mu') \sim h(x - \mu)$  (symbolic similarity); the two problems both look the same (same  $h(\cdot)$ ) and assign the same probability density to equivalent data

Provided there is no information in  $\mathcal{C}_1$  or  $\mathcal{C}_2$  identifying a special location, a formal rule should assign the same functions as priors:

$$f(u) = g(u)$$

**Transformation:** The shift in coordinates,  $\mu' = \mu + \Delta$ , implies that a consistent assignment of prior probabilities must obey

$$f(\mu)d\mu = f(\mu' - \Delta)d\mu' = g(\mu')d\mu'$$

Since symmetry implies  $f = g$ ,

$$f(u - \Delta) = f(u)$$

Only a *constant PDF*  $f(u) = C$  satisfies this functional eq'n

## Scale parameter

For a PDF  $p(x|\sigma)$  for  $x$  with parameter  $\sigma$ , if it is of the form  $p(x|\sigma) = h(x/\sigma)/\sigma$ , then  $\sigma$  is called a *scale parameter* (E.g., Poisson distribution above, with  $\sigma = 1/r$ )

$\mathcal{C}_1$ : Estimate a scale parameter,  $\sigma$ , using a sampling distribution for data  $x_i$  of the form

$$p(x|\sigma, \mathcal{C}_1) = \frac{1}{\sigma} h(x/\sigma);$$

Denote the  $\mathcal{C}_2$  prior by  $P(\sigma \in d\sigma | \mathcal{C}_1) = f(\sigma) d\sigma$

$\mathcal{C}_2$ : In a coordinate system rescaled by  $s$  (e.g., changing units), use data  $x'_i = sx_i$  to estimate  $\sigma' = s\sigma$ .

Since  $dx' = s dx$ ,  $p(x'|\cdots) = p(x|\cdots)/s$ , so the sampling distribution in terms of  $x'$  is,

$$p(x'|\sigma', \mathcal{C}_2) = \frac{1}{\sigma'} h(x'/\sigma');$$

Denote the  $\mathcal{C}_2$  prior by  $P(\sigma' \in d\sigma' | \mathcal{C}_2) = g(\sigma') d\sigma'$



**Symmetry:** Note that  $h(x'/\sigma')/\sigma' \sim h(x/\sigma)/\sigma$ ; the two problems both look the same (same  $h(\cdot)$ ) and assign the same probability density to equivalent data

Provided there is no information in  $\mathcal{C}_1$  or  $\mathcal{C}_2$  identifying a special scale, a formal rule should assign the same functions as priors:

$$f(u) = g(u)$$

**Transformation:** The shift in scale,  $\sigma' = s\sigma$ , implies that a consistent assignment of prior probabilities must obey

$$f(\sigma)d\sigma = f(\sigma'/s) d\sigma'/s = g(\sigma') d\sigma'$$

Since symmetry implies  $f = g$ ,

$$f\left(\frac{u}{s}\right) = sf(u)$$

Only  $f(u) = C/u$  satisfies this; this PDF is flat in  $\log(u)$

# Takeaways

- The axioms/theorems (BT, LTP) are only half of probability theory (grammar); we also need *rules that assign values/functions* in particular settings (vocabulary)
- *Symmetries* can play a key role, providing a kind of definition for “uninformative” — Probability assignment rules should not distinguish between problems equivalent in form
- Symmetries can lead to *functional equations* for PDFs (sampling distributions or priors)
- In inference problems, *the form of the likelihood function* can identify relevant symmetries — We aren't assigning PDFs to greek letters, but to symbols with meaning