

# Bayesian Analysis (BA)

(A name I prefer to the usual Bayesian Inference)

- 1) Preamble: BA is path to truth quantified  
a scientific method quantified. (p. 2)
- 2) Unquantified or Qualitative Bayesian Analysis  
Priors & Posteriors (p. 3)
- 3) Bayes Theorem & a bit of History (p. 6)
- 4) Proof of Bayesian Analysis in the ideal limit — which can be approached arbitrarily closely (p. 11)  
→ Truth (5) Bayes Odds Ratio & Bayes Factor (p. 21)
- 7) Multinomial Theorem & Multinomial Probability (p. 23)
- 8) An Example of Bayesian Analysis  
A toy example: The Die Problem
- 9) Marginalization & Occam's Razor  
shortcuts to Bayesian analysis
- 10) BIC = Bayesian Information Criterion  
AIC = Akaike Information Criterion

17002  
(7.1)

# Preamble

Note in this lecture I give a proof of Bayesian analysis. Not an explanation of how to do it in practice — immense amount of tricks, procedures, formulas

I argue that Bayesian Analysis (BA = AKA Bayesian inference) on Bayesian statistics

but then our computers packages

is a true theory of binding truth or i.e., a true theory in an ideal limit that can be approached arbitrarily closely ideally

many ~~important~~ important theories are like this.

at least improvement. The Scientific Method (Qualitative and theory) proof of the Sci. Method

In fact virtually all, Except maybe TOE (theory of everything or 2nd law of Thermodynamics) But ~~they~~ are

E.g., Newtonian Physics,  $\hookrightarrow$  It is believed to be the macroscopic, <sup>limit</sup> of QM low velocity <sup>limit</sup> of Special relativity weak gravity <sup>limit</sup> of General relativity

Another ~~perspective~~ perspective is ~~it~~ <sup>Newtonian physics</sup> an approximate theory

absolutely exact. But maybe there exactness is quibbling

But you're truly (following a head) would call it a true emergent theory  $\rightarrow$  exactly true in the classical limit which can be approached arbitrarily closely and is approached very closely

in everyday life.

Maybe exactly mid? Well the system must be big enough?

Other examples 2<sup>nd</sup> law of thermodynamics  
(a law in qualitative sense too by reason)

Natural selection evolution

↳ unlike Newtonian Physics these can be proven by mathematical logic (and low? Well for my toy universe you imagine)

— and so can Bayesian Analysis

↳ aside from pure rigor or pure philosophical quibbling in my opinion.

## 17.2) Qualitative Bayesian Analysis

⇒ We do it all the time and so all other conscious beings to one degree or another.

All of life experience gives vague probabilities about how

what <sup>things</sup> will happen in ~~many~~ in a given situation

e.g., impossible, virtually impossible, highly unlikely, unlikely, possibly, likely, very likely, virtually certain, certain

170041

These are your prior probabilities  
or priors

Based on vague approximate frequency

$\frac{N_{\text{event}}}{N_{\text{trial}}}$   
often with a lot of mistakes and fluctuations

Then new experience happens in that situation

and you vaguely update

~~your~~ priors

to your posteriors

(i.e. posterior probabilities)

(which sounds awful and often  $\Gamma$ s)

Your updating is also qualitative and of variable accuracy,

but it works well enough mostly in everyday life ~~to help~~ to your advantage

Nature via Natural selection evolution has its own way of improving success rate relative to local conditions.

Of course, things like major impactors can change the rules of success suddenly  $\rightarrow$  like what happened to non-avian dinosaurs

Ex. Job interviews

17005

↳ especially if you are new at the game — or the 'rules' have changed → each interview causes you to update your priors to posteriors often in an intuitive (but still useful) way.

To some degree 'qualitative Bayesian' analysis trial & error

But

improvement if just restricted to each new experience ~~of a unique situation~~ from has very restricted generality.

But as <sup>more</sup> general conclusions about updates occur  
↳ qualitative Bayesian analysis approaches the scientific method.

6  
177006

Bayesian analysis itself  
is the sci. Method Quantified  
(I argue and the proof  
therefore of the Scientific  
method)

17.3)

Bayes Theorem

{ At root of BA  
and general to  
all probability  
theory in fact

Root of Bayesian analysis is  
Bayes Theorem — which is  
really simple and simple to prove.  
Easier to prove than to remember.

Discovered by Thomas Bayes (c. 1701-1761)  
and published posthumously in 1763. (W.K)

Pierre-Simon Laplace (1749-1827) independently  
discovered it (and published it in 1774)  
(W.K)

- intro 267  
p. 3  
tells me  
W.K  
- all of  
- conditioning  
- information  
- about  
in  
BA  
proof

Consider 3 events A, B, K.  
K is background knowledge that  
is introduced as extra term because  
I need it later. It isn't needed for <sup>the</sup> proof  
Bayesian theorem itself.

I use ABK as joint event as an  
shorthand for union =  $\cup$  symbol  
in math  
which is too klutzy for me in this  
context

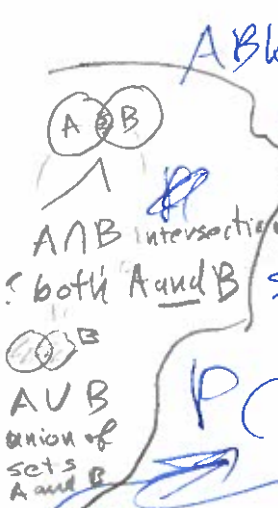
order  
has no  
meaning.

2025 May 01

7  
17007

# Conditional Probability, and Probability, Product Rule

Consider Events  $A, B, K$   $\leftarrow$  Need  $K$  for proof of Bayesian Analysis (BA)



$ABK$  is ~~joint event~~ <sup>joint event intersection</sup>  
 all 3 stone  $ABK$   
 $ABK = \overbrace{A \cap B \cap K}^{\text{intersection}}$

Background knowledge that in BA you are often called conditioning knowledge on context. (Lavello 2024 p. 3)  
 a shorthand and looks ugly in general

so NOT a product

$$P(ABK) = \frac{N_{ABK}}{N} = \frac{N_{ABK}}{N_K} \frac{N_K}{N}$$

Probability of all 3 at once  
 - order has no meaning

Frequentist definition  
 $N_{ABK}$  cases of  $N$  trials taken to infinity

Nothing forbids this factorization

$$P(ABK) = P(A|BK)P(B|K)P(K)$$

By symmetry

$BK$  is also given  $BK$  implicitly

$$P(ABK) = P(B|AK)P(A|K)P(K)$$

$\therefore$  Bayes theorem in symmetrical form

$$P(A|BK)P(B|K) = P(B|AK)P(A|K) \text{ GED}$$

Asymmetric form

$$P(A|BK) = \frac{P(B|AK)P(A|K)}{P(B|K)}, \quad P(B|AK) = \frac{P(A|BK)P(B|K)}{P(A|K)}$$

easier to prove than to remember.

Usually show with  $K$  implicit, but for proof of BA, we need  $K$  explicit.

Some assert this is an axiom or <sup>(Truth 2017 p. 6)</sup> can be so taken.

But I believe probability only has meaning in a frequentist sense even if  $N_{ABK}/N$  is only vaguely known or hypothesized imaginarily.  
 In fact, BA means you can improve from vaguely known to "truth".

$N_{ABK}$  &  $N$  exist in some ideal limit as Galileo would argue.

17008

That's all: a profound and profoundly simple logical rule of reality

# Root of Bayesian analysis

Bayesian analysis was greatly elaborated in 1939 ~~OK~~ (Scientific Inference 3rd ed 1973) by Harold Jeffreys (1891-1989) (no relation - but he manages to look like my father anyway)

## What is Bayesian analysis?

Actually, a actual practice. The ideal steps beyond step 1 are not done. Everyone seems to start step 1 with the hope that step 1's will gradually get better

It's a path to true theories by using Bayes theorem to update prior probabilities for posterior probabilities for theories until one true theory is left.

In toy cases this is easy. In hard cases, <sup>over</sup>hard, because there is a lot of slop in assigning priors. [Not to theories in almost cases but the priors on theory parameters.]

A theory with free parameters in one sense is a continuum of theories

They've not done in single papers. Now we add this date, now this and take away that

No one much used Bayesian analysis before 1970s. Since then it has had growing vogue.

Why was it unused before?

17009

Well, ~~it~~ <sup>statistical</sup> power is in dealing with statistical theories (i.e., theories that only give probabilities) and interesting cases in cosmology, epidemiology, social sciences

only come with vast data sets and vast computing power to manipulate them. But since 1990s ~~we~~ on, we have enough of both.

So nowadays, Bayesian analysis is an ubiquitous tool (though less wonderful than one can hope).

You may ask how can a theory have a probability of being true, isn't it just true or false

17010

speaking in an absolute sense <sup>not ~~working~~ leaving aside partially</sup>

true theories

or having some truth and some falseness definitely in

less absolute sense <sup>not leaving aside partially true theories</sup>

But the probability of truth is NOT in the absolute sense.

It's probability to our knowledge

Example

I've a coin in my hand. Is it heads or tails?

It's one or the other in an absolute sense.

But to your knowledge it's 50% - 50%

See, Nothing in my hand   
 → So you should include the theory of deception.

2026jan24

17011

17.4) Proof that Bayesian Analysis in the Ideal Limit  
is the Path to True Theories But You May  
Need to Restrict their Realm of Validity

a) In your truly's view, Bayesian analysis (BA) is the scientific method quantified. Many probably share this view, but many probably disagree since philosophers of science are great quibblers.

The proof is NOT absolutely vigorous.  
NOR beyond quibbling.

Say you have background knowledge which in jargon of Loewdo & Wolpert (2024, p.3) is called context  $C_0$  or conditioning information which is all knowledge relevant to the aspect of reality whose true theory is the goal of your BA.  $C_0$  includes general theories, specific information, anything at all relevant.

From  $C_0$ , you deduce a set of theories  $\{T_i\}$  which could be considered already in  $C_0$  or adding to  $C_0$ . The set  $\{T_i\}$  does NOT have to be complete even to  $C_0$  as it is, but in practice

17012)

it should be all theories you think most probable by

qualitative Bayesian analysis

Completeness is a rather vague concept any way. To recapitulate, the power of BA is that it works

in the ideal limit no matter how much vagueness there is in

choosing  $\{T_i\}$ , assigning their initial probabilities, the multiplicative factor frequentist probability assumption, etc.

Note,  $\{T_i\}$  does NOT have to include the true theory

but if it does NOT, you will have to do one or more episodes of qualitative Bayesian analysis to determine a new set of most probable theories.

b) The procedure of ideal BA analysis requires a series of new data acquisitions;

$D_1, D_2, D_3 \dots D_e \dots$

which I think does NOT have to be infinite

Note, the theories are assumed to be a finite discrete set, They have no free parameters

A theory with free parameters is in a sense a continuum family of theories. Of course, many theories have free parameters. For those, you need marginalization. See p. 170

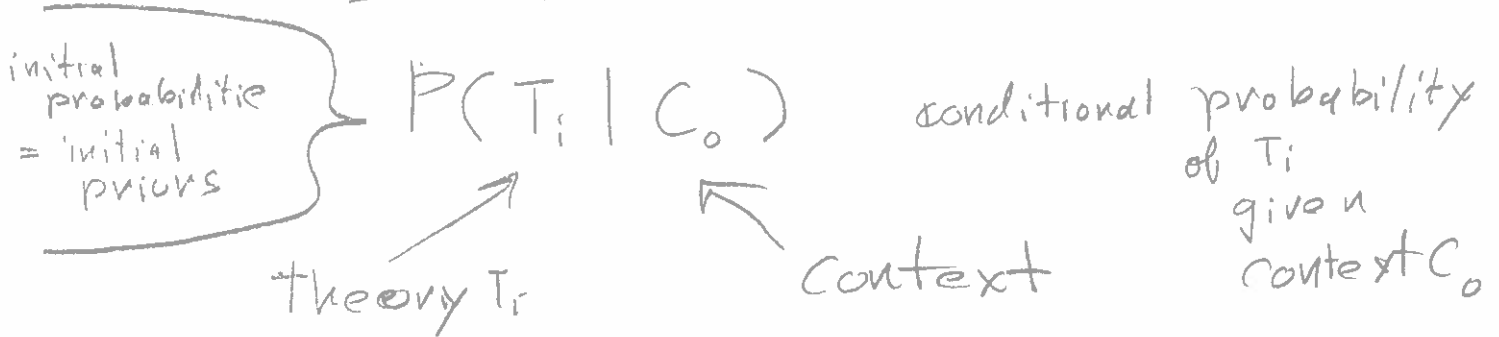
2026 Jun 29

17013

provided all knowledge relevant to the analyzed aspect of reality somehow.

Note  $D_e$  can be measurements, but also newly acquired knowledge or theories about some other aspects of reality

The 0th step in the procedure is initial probabilities to the set of theories  $\{T_i\}$ .



How do you do this?

In principle, anyway you like.

The ideal BA works no matter how.

But since you thought  $\{T_i\}$  was the most probable set and may NOT have any reasonable quantitative ranking, you can use the principle of indifference (Wik).

Just set all  $P(T_i | C_0)$  equal.

You can normalize the probabilities if you like!

$$\sum_i P(T_i | C_0) = 1$$

17017)

but this isn't necessary since  
ideal BA normalizes updated  
theory probabilities automatically,  
(see p. 17017)

As you acquire the new sets  
of data  $D_1, D_2, D_3, \dots, D_e, \dots,$

the context updates

$$C_e = C_{e-1} D_e$$

Where as a shorthand  
I use 'product'

to mean union: i.e.

$$C_{e-1} D_e = C_{e-1} \cup D_e$$

Union of set S (Wiki)

I think this symbol klutz  
typed or handwritten.

You use  $D_e$  to update your  
prior probabilities at step  $e$   
to update your prior probabilities  
at step  $e$  (i.e.  $P(T_i | C_{e-1})$ )  
to your posterior probabilities at step  $e$   
(i.e.  $P(T_i | C_e)$ )



17016]

Actually,  $P(D_e | T_i, C_{e-1})$   
is NOT from a marginalization  
for our case of a discrete set  
of theories with no free parameters,  
but for consistency when do marginalization  
to get  $P(D_e | T_i, C_{e-1})$ . (see p. 17...)

we call it that here following the  
jargon of Kass & Raftery (1995, p. 776).

But what is denominator  $P(D_e | C_{e-1})$ ?

We can only do an estimate based  
on the set of theories  $\{T_i\}$  we are  
considering. We expand using factorization

$$P(D_e | C_{e-1}) = \sum_j \frac{N_{D_e}}{N_{T_j, C_{e-1}}} \frac{N_{T_j, C_{e-1}}}{N_{C_e}}$$

Really  
just a  
best  
estimate  
to your  
knowledge.

$$= \sum_j P(D_e | T_j, C_{e-1}) P(T_j | C_{e-1}) P(C_{e-1})$$

$$= \sum_j P(D_e | T_j, C_{e-1}) P(T_j | C_{e-1})$$

But  
this is 1  
since  
we  
know it.

knowledge. You do NOT know

$P(D_e | C_{e-1})$  in an absolute sense.

You could include theories NOT in your set  $\{T_i\}$ , but that  
seems pointless practically AND as a formalism

2026 Jan 24

17017

$$P(T_i | C_e) = \frac{P(D_e | T_i, C_{e-1}) P(T_i | C_{e-1})}{\sum_j P(D_e | T_j, C_{e-1}) P(T_j | C_{e-1})}$$

You can regard the coefficient as an updating factor from prior  $P(T_i | C_{e-1})$  to posterior  $P(T_i | C_e)$

Note the  $P(T_i | C_e)$  are automatically normalized as noted on p. 17018:  
 $\sum_i P(T_i | C_e) = 1.$

$$= \left( \frac{P(D_e | T_i, C_{e-1})}{\sum_j P(D_e | T_j, C_{e-1}) P(T_j | C_{e-1})} \right) P(T_i | C_{e-1})$$

$$= \left[ \frac{P(D_e | T_i, C_{e-1})}{\langle P(D_e | T_j, C_{e-1}) \rangle} \right] P(T_i | C_{e-1})$$

Weighted mean of the marginal likelihood

Except for format niceness

If  $P(D_e | T_i, C_{e-1})$

exceeds/subceeds the mean, the probability of theory  $T_i$  increases/decreases.

(i.e., knowing it exists and can be calculated), you do NOT need the denominator. You just update the relative probabilities.

Recall the denominator formula is actually just a best estimate to your knowledge.

d) As the BA Procedure continues } See p. 17016  
What happens

As data sets  $D_e$  continue, the probabilities of the theories will continue be updated.

17018

i) Quasi-ultimately in many cases,  
one  $T_{i^*}$  will grow in probability  
heading toward  $P(T_{i^*} | C_e) = 1$

and the other theories

will head toward  $P(T_i | C_e) = 0$ .  
Some may have reached  $P(T_i | C_e) = 0$

ii) However, unlucky sets of data  $D$   
may cause  $P(T_{i^*} | C_e)$   
to decrease sometimes  
for some steps.

during  
the  
procedure  
and  
can  
be  
dropped  
then,  
of  
course.

iii) Does  $T_{i^*}$  reach  $P(T_{i^*} | C_e) = 1$   
in finite or infinite steps.

It can reach it in finite steps,

but if the theories  $\{T_i\}$

are statistical and you

gather statistical data  $D_e$

$P(T_{i^*} | C_e)$  may go 1

only as  $l \rightarrow \infty$ .

2026 Jun 24

17012

However, even for statistical theories,  
you may be able to access data  
with is not determined probabilistically  
and reach  $P(T_i | C_e) = 1$   
for finite  $l$ .

I give a toy example of the  
ideal BA procedure on p. 17  
in which this happens.

iv) What if at step  $l$  all  $P(T_i | C_e) = 0$ ?

Then you have to do qualitative Bayesian analysis  
and invent a new set of theories  $\{T_i\}$ .

Of course, you don't have wait  
for this catastrophe.

You can drop theories from the  
procedure when their  $P(T_i | C_e) \ll 1$   
and not wait for  $P(T_i | C_e) = 0$ .

And the growing  $C_e$  may lead  
you to new theories that seem  
probable along the way and you  
can add them with reasonable reassignment  
of probabilities.

17020

v) Does  $P(T_{i*} | C_e) = 1$

mean  $T_{i*}$  is the true theory?

No. It just means all the other theories have zero probability.

If you keep doing the procedure with just one theory at some data set  $D_e$ ,

$$P(T_{i*} | C_e) = 0.$$

Then you have to do qualitative Bayesian analysis and invent a new set of theories and start over.

vi) But what if  $C_e$  is enormous and the  $P(T_{i*} | C_{e+1}) = 0$ ,

In this case, you may well want to retain  $T_{i*}$ , but just limit its scope to  $C_e$ .

2024 Jan 24

17021

In one sense, this is lowering the bar.

But theories  $T_i$  that are adequate for a vast realm  $C_i$  usually shouldn't be dismissed.

Example Newtonian physics was once hypothesized to be absolutely fundamental for all phenomena.

But even early people couldn't explain electromagnetism or chemistry by it though they may hoped that that would happen someday.

However, by later 19th century that hope was diminishing.

For example, classical electromagnetism which seemed an extremely robust theory was based on Maxwell's equations which were known NOT to be invariant under the Galilean transformations. Famously, Einstein was led to special relativity by, among other things, the belief that

17022

Maxwell's equations had to be more fundamental than Newtonian physics.

The advent of classical electromagnetism, special relativity, general relativity, quantum mechanics, and quantum field theory vastly limited the scope of Newtonian physics.

However, saying it's just an approximate theory seems inadequate considering its vast realm of applicability.

Following a hunch, I'd say it is an exact theory in the classical limit

where relative velocities  $\ll c$ ,

size scale  $\gg$  microscopic,

the local gravity field

is sufficiently uniform,

and one avoids some tricky electromagnetic effects.

It's a true or fundamental emergent theory in the classical limit.

Newtonian can never be proven false, it can only maybe become limited.

This is true of other profound theories too.

We believe general relativity fails in the quantum gravity realm even though we don't have an established quantum gravity theory.

2024 Jan 27

17023

vii) Does BA guarantee that you will get to a true theory even if only limited to vast  $C_x$ ?

I think only if you can keep interrogating reality until you know everything relevant to the aspect of reality of interest.

To give an example of case where practically you may never know the true theory.

Say you just measure a series of 64-bit floating point numbers (15 to 17 digits) in the range  $(0, 1)$  (and so excluding endpoints 0 and 1).

You measure a long range and all ordinary statistical tests find them to be randomly distributed over range  $(0, 1)$ .

But are the numbers fundamentally random (as set by atmospheric noise; i.e., atmospheric radio noise; Wik) or computer generated random by a deterministic algorithm?

(17024)

e.g., the Mersenne-Twister (MT)

which has repeat cycle =  $2^{19937} - 1$

(WIK: Mersenne-Twister: characteristics)

Say you had an exaflop computer

and it takes 100 flops (just guessing) to compute

one MT random number, how

long until you complete a cycle?

$$t = \frac{(100 \text{ flops}) * (2^{19937} - 1)}{10^{18} \text{ flops/s}}$$

$$\approx \frac{10^2 * 10^{[0.3 * (20000)]} \text{ flops}}{10^{18} \text{ flops/s}}$$

$$= \frac{10^{6000}}{10^{18}} \text{ s} \approx 10^{6000} \text{ s} * \left( \frac{1 \text{ year}}{\pi * 10^7 \text{ s}} \right)$$

$$\approx 10^{6000} \text{ years} \approx 10^{6000} \text{ Gyr}$$

Practically, you could never tell as long as MT numbers were as good as claimed.

But if you could see the MT numbers were just coming from an isolated computer, you'd know they were algorithm generated.

2026 Jan 24

This example is a pathological case, 117025  
but it does prove an interesting point,  
events can be completely  
deterministic as to source,  
but completely random  
for virtually all purposes  
as to receiver.

### viii) Speeding Up Bayesian Analysis: Ideal, Practical, or Qualitative

To speed any of these up,  
choose your data  $D_e$  acquisitions  
to be as decisive as possible.

To quote Ernest Rutherford (1871-1937)

"If you need statistics,  
you did the wrong experiment"  
(Trotter 2017, p. 4)

The pith of this aphorism is as above,  
choose your data acquisitions  
to be as decisive as possible.

17026

But Rutherford lived in simpler times where in physics at least, you need less statistics and, of course, Rutherford did use statistics as needed.

But in the modern age, in areas of cosmology, epidemiology, psychology, economics, social science, and AI statistics is what we've got.

But we have vast data sets and vast computing power make use of them. Much of Bayesian analysis was worked out in the 1950s by Harold Jeffreys (no relation), but only with increasing computing power over the decades has Bayesian analysis grown in importance.

I think I'd never or barely heard of it before year 2000 and only since teaching cosmology I have become interested and NOT for practice, but just to understand why it's the path to truth (i.e., true theories)

2025, Jan 27

1702Y

ix) A key Point about Ideal Bayesian Analysis

You may start out very bad guesses at probable theories and your data acquisitions may be far less than decisive, but you will still approach truth if you keep doing the procedure (without making and repeating mistakes).

People often Bayesian analysis (or Bayesian Inference) is theoretically well justified and hopefully this section nonrigorously is a valid justification.

People often say other statistical inference is NOT so well justified (but I have personal expertise on this point).

17028

17.5 Bayesian Analysis: What People Actually Do  
Including Bayes Odds Ratios, and  
Bayes Factor But Deferring  
Marginalization to p. 17053

So far as I know, no one in practical work does more than one explicit BA step of ideal BA in a paper. From paper to paper, it's just qualitative Bayesian analysis,

The ideal BA procedure described on p. 17011 - 17027, is just to show that ideal limit is true and you can approach that truth even just in one explicit Bayesian step.

In fact, what they do in that one-step BA is effectively to conflate multiple steps by expanding and contracting  $D_1$  by including and excluding specific data sets.

A big practical problem is systematic error in data sets. By definition, the size of systematic error is unknown or you would be able to correct for it.

Recent papers applying BA to cosmological data have spent a lot of effort to detect systematic error by finding inconsistency between data sets.

Weird things turn up: eg.,

i) Two data sets support the same model but with inconsistent free parameters. To some degree, this is the Hubble tension

ii) Two data sets support some of the same parameters, but for different models. I can't think of an example, but this may happen.

One thing people do NOT do is bother with the denominator  $P(D_e | C_{e-1})$

It is easy to calculate given that you know the  $P(T_i | C_{e-1})$ 's and

$$= \sum_j P(D_e | T_j, C_{e-1}) P(T_j | C_{e-1})$$

(see p. 17016)

$P(D_e | C_{e-1})$ 's, but you only need relative probabilities to judge the probability of truth  $\{T_i\}$ .  
re:  $C_{e-1}$ , adequacy of theories in your set  $\{T_i\}$ .

17030

This is true the ideal Bayesian analysis too. See p. 17017

What they do is compute the Bayesian posterior odds ratio on just the Bayes factor

(jargon of Kass & Raftery 1995, p. 776)

From p. 17017

$$OR_{postij} \equiv \frac{P(T_i | C_l)}{P(T_j | C_l)} = \frac{P(D_l | T_i, C_{l-1}) P(T_i | C_{l-1})}{P(D_l | T_j, C_{l-1}) P(T_j | C_{l-1})} = K_{ij}^{(l)} * OR_{preij}$$

$K_{ij}^{(l)}$  means on step  $l$ , not a power

Bayes  $K$  factor, but usage from Wik, NOT Jeffreys original definition (Kass, p. 776)

prior odds ratio

$$K_{ij}^{(l)} \equiv \frac{P(D_l | T_i, C_{l-1})}{P(D_l | T_j, C_{l-1})}$$

$$OR_{preij} \equiv \frac{P(T_i | C_{l-1})}{P(T_j | C_{l-1})}$$

$K_{ij}^{(l)}$  = evidence for  $T_i$  / evidence for  $T_j$

evidence is synonym for marginal likelihood (Wiki: Marginal likelihood)

Of course, in One-Step Bayesian analysis  $l = 1$

In One-Step BA, the initial probabilities  $P(T_i | C_0)$  are usually just chosen by the principle of indifference

17031

because the set of theories  $\{T_i\}$  has already been chosen

by qualitative BA to be the likeliest theories,

With the choice  $P(T_i | C_0)$  all equal, the <sup>to your knowledge</sup> relative judgment of adequacy of the theories all comes down to the Bayes factor  $K$ ;

the ratio of evidences (i.e., the ratio of marginal likelihoods)

which ratio or its logarithm!

I think the ratio of evidences sometimes is just called the "evidence"

or the "Bayesian evidence" itself.

Jeffreys offered an Evidence Table judging the "value" of the evidence based on his examples or empirical testing (I think). Different versions have been offered, but I do NOT know of any optimum table.

17032

Evidence

However, the Jeffreys table simplified by Kass & Raftery (1995, p.777) (probably because the original was speciously over-precise) is widely cited and may be the Arducial table.

Jeffreys Evidence Table for the Judgment of Bayesian Evidence (Bayesian K factor)

Simplified by Kass & Raftery (1995, p.777)

| $K_{ij}^{(2)}$ | $\log(K_{ij}^{(2)})$ | Significance of Evidence                     |
|----------------|----------------------|--|
| $< 1$          | $< 0$                | Negative (in favor of <del>...</del> $T_j$ ) |
| 1 a 3.2        | 0 a 0.5              | Insignificant                                |
| 3.2 a 10       | 0.5 a 1              | Significant                                  |
| 10 a 100       | 1 a 2                | Strong                                       |
| $> 100$        | $> 2$                | Decisive                                     |

ii)  $K_{ij}^{(2)} = \frac{P(D_2 | T_i, C_{2-1})}{P(D_2 | T_j, C_{2-1})}$  where usually  $C_{2-1} = C_0$

and  $D_2 = D_1$

for one-step Bayesian Analysis

ii) Note: Systematic errors could nullify the evidence

2026 Jan 24

17033

Why is Evidence Table so cautious in being decisive?

Well, generally one is cautious about being decisive in science (e.g., the 5 $\sigma$  rule for claiming discovery). The evidence (using the term in its everyday sense) has to be strong because there can be so many mistakes in data and analysis. There's also Carl Sagan's (1939-1996) aphorism: "Extraordinary claims require extraordinary evidence" (ECREE), but others said similar things earlier (e.g., David Hume (1711-1776)). However, to be specific to Bayesian analysis, in marginalization (see p. 170.53), one must choose priors on the free parameters of theories and

17034

that choice is fraught  
with uncertainty and  
vast variation in what  
is thought to be good.

A poor choice of priors  
can make a bad theory look good.

In fact, I think people are  
dismissive of  $K_{ij}^{(e)} \leq 10$  and  
may NOT even consider  $K_{ij}^{(e)} = 100$   
decisive.

But if  $K_{ij}^{(e)} \geq 1000$  that may  
be decisive even if the  
choice of priors is very poor.

Also a philosophical difference in  
choosing priors can arise

- i) Choose them theory-independently
- ii) Or choose prior (on the  
range of free parameters)  
that is physically plausible for  
a theory

NOT priors  
on a theory  
but on the  
range of  
allowed  
for the  
free  
parameters  
for a  
theory  
in

marginalization

See  
p. 170

17.6

# Multinomial Probability Distribution & The Multinomial Theorem

a) There are probabilities,  $P_1, P_2, P_3, \dots, P_I$  and  $\sum P_i = 1$

Say you had  $I$  bins from which to draw events or in statistical mechanics terms in which to put ~~part~~ classical particles classical but identical

|       |       |       |       |       |     |       |
|-------|-------|-------|-------|-------|-----|-------|
| row 1 | $P_1$ | $P_2$ | $P_3$ | $P_4$ | ... | $P_I$ |
| row 2 | $P_1$ | $P_2$ | $P_3$ | $P_4$ | ... | $P_I$ |
| row 3 | $P_1$ | $P_2$ | $P_3$ | $P_4$ | ... | $P_I$ |
| ...   | ...   | ...   | ...   | ...   | ... | ...   |
| row N | $P_1$ | $P_2$ | $P_3$ | $P_4$ | ... | $P_I$ |

Do N selections with replacement in stats jargon (Wiki sampling going down the rows)

~~Point~~ Pointer-II calls these arrangements but these are equal probability in that book

You get sequences: e.g.

- seq 1:  $P_1 P_1 P_1 \dots P_1 = P_1^N$
- seq 2:  $P_1 P_2 P_1 \dots P_1 = P_1^{N-1} P_2$
- seq 3:  $P_2 P_1 P_1 \dots P_1 = P_1^{N-1} P_2$

probability of this one sequence. These are 2 sequences of equal probability.

interchanging the identical subscripts does not give a new sequence

~~interchanging~~ permuting interchanging the 2 and 2 indices of sequence 2 gives sequence 3 which is different sequence, (also a different distinct sequence)

17036

(2026 Jan 24)

b) I have never found a good way in words for explaining why permuting identical indices gives NO distinct sequence.

All I can say is look at the procedure on p. 17034,

Note, that procedure is actually identical to expanding the multinomial  $(P_1 + P_2 + \dots + P_I)$  by power  $N$ ; i.e.,

$$(P_1 + P_2 + \dots + P_I)^N \text{ which is why}$$

we call the Multinomial distribution the multinomial distribution. The expansion is the multinomial theorem.

The distinct sequences are called combinations (Wiki: combination)

Combinations are identical if they

conform to set  $\{n_i\}$   $\{n_i\}$

where  $\{n_i\}$  means

- $n_1$   $P_1$ 's
- $n_2$   $P_2$ 's
- $n_3$   $P_3$ 's
- $\vdots$
- $n_I$   $P_I$ 's

The order of the  $P_i$  does NOT cause non-identicality.

In the jargon of statistical mechanics

(at least according to Pointon-10-12), the set  $\{n_i\}$

is called a configuration and the

number of conforming combinations is weight  $C(\{n_i\})$

of the configuration.

How does one calculate weight  $C(\{n_i\})$  and the probability of selecting a conforming combination?

c) Weight:  $N! = C(\{n_i\}) \prod_i n_i!$

For example consider sequence of elements

$p_1, p_2, p_1, p_1, p_6, p_1, \dots, p_{I-3}$

N elements in the selection

∴  $N!$  permutations

i.e., N ways of selecting a first element, N-1 ways of selecting a second element, etc.,

but permuting the  $n_i p_i$ 's among themselves creates NO new combination (i.e., distinct sequence).

Therefore, you can factorize  $N!$  as above and get

$$C(\{n_i\}) = \frac{N!}{\prod_i n_i!}$$

which is the weight of the configuration  $\{n_i\}$

Note, if there are just two elements  $p_1$  and  $p_2$ , then

$$n_2 = N - n_1 \text{ and } C(\{n_i\}) = \frac{N!}{n_1! (N - n_1)!} = \binom{N}{n_1}$$

and also the multinomial coefficient of the multinomial theorem.

which is the binomial coefficient.

17038

d) Probability of Selecting a Combination conforming to configuration

Now the  $P_i$ 's are the probabilities of events  $i$  (see p. 17035).

So the probability of selecting a specific combination conforming to  $\{n_i\}$  is  $\prod P_i^{n_i}$

But what is the total probability of selecting any combination conforming to  $\{n_i\}$  (i.e., the probability of getting configuration  $\{n_i\}$ )?

One guesses  $P(\{n_i\}) = C(\{n_i\}) \prod P_i^{n_i}$ ,

but I can't see any proof just based on what we've said so far

However the probability of getting any combination conforming to any combination or in other words of getting any sequence is 1.

$$1 = 1^N = (P_1 + P_2 + P_3 + \dots + P_I)^N = (\sum_i P_i)^N$$

When you actually do the multinomial expansion,

it is clear (maybe) that the probability of configuration  $\{n_i\}$  being

obtained is  $P(\{n_i\}) = C(\{n_i\}) \prod P_i^{n_i}$

All sequences produced uniquely, therefore all combinations with their correct weighting (see p. 17037)

These the count of combinations conforming to set  $\{\epsilon_{n_i}\}$

$$C(\{\epsilon_{n_i}\}) = \frac{N!}{\prod n_i!}$$

which must be an integer

The probability of combinations have equal value of ~~only one reference~~ and this value ~~giving~~  $\{\epsilon_{n_i}\}$  is  $\prod P_i^{n_i}$

Also probability of any sequence conforming to set  $\{\epsilon_{n_i}\}$

∴ the probability of getting the combination arrangement distribution is conforming to set  $\{\epsilon_{n_i}\}$

$$P(\{\epsilon_{n_i}\}) = C(\{\epsilon_{n_i}\}) \prod P_i^{n_i}$$

Probability distribution

number of combinations conforming to set distribution  $\{\epsilon_{n_i}\}$

Probability of getting any one of the combinations

~~of the distribution of combinations~~  
of set  $\{\epsilon_{n_i}\}$

$$= \frac{N!}{\prod n_i!} \prod P_i^{n_i}$$

e) Now 
$$P(\{\epsilon_{n_i}\}) = N! \prod \frac{P_i^{n_i}}{n_i!}$$

As we know from statistical mechanics this is the distribution for classical identical particles.

In a QM sense, wave functions must obey the symmetrization principle. Boson/Fermion wave functions must be symmetric/antisymmetric

The probability distributions for bosons & fermions are different

Bose-Einstein statistics

Fermi-Dirac statistics

we often say they are because the particles are identical but I think we just mean identical not classical identical

26  
40

and the particle statistics arise from the symmetrization principle in QM which nature demands to prevent infinite degeneracy of multi-particle states or so it seems.

Normalization?

Proof

$$1 = \sum_i p_i$$

sequences of what we demand for particles with spin

$$1 = \left( \sum_i p_i \right)^2$$

sequences of 2 particles

$$1 = \left( \sum_i p_i \right)^N$$

multinomial theorem

with polynomial coefficients and variables  $ax^2 + bx + c$  etc. (wik)

A) If  $i = 2$  you have the binomial theorem

$$(p_1 + p_2)^N = \sum_{r=0}^N \binom{N}{r} p_1^r p_2^{N-r}$$

Binomial coefficient

and  $P(\sum n_i) = \binom{N}{r} p_1^r p_2^{N-r}$

for higher multinomial

~~the formulas~~ are harder

to get formula easy formula for. Are no easy ones hint.

The key point is that in ~~multiply out the expansion into and collect~~ you get just the sequences and collect the coefficients like terms, the multinomial probability distribution

e.g.  $(p_1 + p_2)^2 = p_1 + p_1 p_2 + p_2 p_1 + p_2^2 = p_1 + 2p_1 p_2 + p_2^2$

Of course, if you actually started collecting sequences of ~~putting~~ putting  $N$  particles into  $I$  states, it is only in the limit of ~~large~~  $l \rightarrow \infty$  sequences that you would recover the multinomial probability distribution

Unit: Digression

9) Further Digression on Binomial Theorem

Binomial Thm

Special case: Binomial theorem & Binomial probability distribution

Prove by induction if affected by Paravola.

$$1 = (p + q)^n = \sum_{k=0}^n \frac{n!}{k!(n-k)!} p^k q^{n-k}$$

for  $q = 1 - p$

Moments of the binomial distribution.

A old trick

Define function

$$f(x) = (x + q)^n = \sum_{k=0}^n \binom{n}{k} x^k q^{n-k}$$

Just called binomial coefficient. (With)

$$M_x = \frac{d}{dx} f(x) \Big|_{x=p}$$

$$\left( \frac{d}{dx} \right)^r f(x) \Big|_{x=p}$$

~~$$f(x) = (x + q)^n$$~~
~~$$f(x) = (x + q)^{n-1}$$~~

$$\sum_{k=0}^n k^r \binom{n}{k} x^k q^{n-k} \Big|_{x=p} = \left( x \frac{d}{dx} \right)^r (x + q)^n \Big|_{x=p}$$

the  $n$  choose  $k$  formula

28) 117042 |  
l=0

$$M_0 = 1$$

~~28~~

$$l=1 \quad M_1 = \lambda n (x+q)^{n-1} \Big|_{x=p} = np$$

$$\begin{aligned}
 l=2 \quad M_2 &= n \lambda \frac{d}{dx} \left[ \lambda (x+q)^{n-1} \right] \Big|_{x=p} \\
 &= np \left[ \cancel{x+q}^{n-1} + n(n-1)(x+q)^{n-2} \right] \Big|_{x=p} \\
 &= np [1 + p(n-1)] \\
 &= pn + p^2 n (n-1)
 \end{aligned}$$

$$\begin{aligned}
 \sigma^2 &= \langle (x - \bar{x})^2 \rangle = \langle x^2 - 2x\bar{x} + \bar{x}^2 \rangle \\
 &= \langle x^2 \rangle - \bar{x}^2 = pn + p^2 n (n-1) - p^2 n^2 \\
 &= \cancel{p^3 n (n-1)} = pn - p^2 n \\
 &= pn(1-p)
 \end{aligned}$$

BeV-53  
correct

probably now totally worthless

2021 nov 22

$P_i = \frac{1}{\Omega_i}$   $\sum P_i = 1$  as it should  
 This simplified Navier and got  $\{n_i\}$

~~29~~  
17093

Theory 1

$$P(\{n_i\}) = \frac{N!}{n_1! \dots n_I!} \prod_{i=1}^I \left(\frac{1}{I}\right)^{n_i}$$

$$= \frac{N!}{n_1! \dots n_I!} \left(\frac{1}{I}\right)^N$$

$\sum n_i = N$   
 This is for I write issue  $P_i = 1$  odd PK2  $\sum$  Total  $= \frac{I}{2} \cdot 1 + \frac{I}{2} \cdot 2 = \frac{3I}{2}$

Theory 2

$$P(\{n_i\}) = \frac{N!}{n_1! \dots n_I!} \prod_{i=1}^I \left(\frac{i}{I(I+1)}\right)^{n_i}$$

$$= \frac{N!}{n_1! \dots n_I!} \frac{1}{\left(\frac{I(I+1)}{2}\right)^N} \prod_{i=1}^I i^{n_i}$$

As so often in Astronomy, not has told is there are only

Theory 3

$$P(\{n_i\}) = \frac{N!}{n_1! \dots n_I!} \prod_{i=1}^I \left[2 - \text{mod}(i, 2)\right]^{n_i}$$

2 possible true theories

Very easy to point out PCDIT) not like Taylor system with info data

5) Die Problem for Zeus II

Normal problem

$I = 6$   
 Theory 1  $P_i = \frac{1}{6}$

Theory 2  $P_i = \frac{2 - \text{mod}(i, 2)}{9}$

$P_1 + P_2 + \dots + P_6 = 1$

$P(I_1) = \frac{1}{2}$   
 $P(I_2) = \frac{1}{2}$

Principle of indifference (Barnes p. 6)  
 $P(I_1) = P(I_2) = 1$

So ~~10~~ throw complex enough but ~~is~~ device?

Really deal this in 2019

|   |   |   |   |   |   |   |   |   |    |
|---|---|---|---|---|---|---|---|---|----|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| 1 | 5 | 1 | 1 | 2 | 5 | 2 | 6 | 3 | 4  |

$P(DIT) = \frac{10!}{3! 2! 2! 1! 1! 1!} \cdot \left(\frac{1}{6}\right)^{10} = 2.50057 \times 10^{-3}$   
 $P(DIT) = \frac{10!}{1! 1! 1! 1! 1! 1!} \cdot \left(\frac{1}{6}\right)^{10} = 6.938 \times 10^{-4}$

math  
6  
10  
and  
4  
2/2

30  
17044

$$P(D|I_1, k_0) = \frac{P(D|I_1, k_0) \times P(I_1|k_0)}{P(D|I_1, k_0) \times P(I_1|k_0) + P(D|I_2, k_0) \times P(I_2|k_0)}$$

$$k = 3.609$$

Probability given your knowledge Before even No experiment

So Prior odds = 1

Bayes factor  $k = 3.609$

Posterior odds = 3.609

If these two theories really are exhaustive, one can calculate their probabilities

See p. 12 for update formula with update to no summing posterior

$$P(1) = \frac{3.609}{4.609} = 0.7828, \quad P(2) = \frac{1}{4.609} = 0.21799$$

If one carried on the experiment to  $N \rightarrow \infty$  for a die,

$$P(1) \rightarrow 1, \quad P(2) \rightarrow 0$$

But my die is not perfect and it can be loaded Face 4 opens!!!

I can load it. It's a real Vegas die

$$P_1(1) = 0.7828, \quad P_2(2) = 0.21799 \dots \text{the posteriors} \rightarrow \text{priors}$$

Do another set of 10

10, 9, 1, 7, 6, 5, 4

|   |   |   |   |   |   |   |   |   |    |
|---|---|---|---|---|---|---|---|---|----|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| 6 | 1 | 3 | 6 | 6 | 2 | 5 | 3 | 5 |    |

$$P(D|I_1) = \frac{10!}{3!2!2!1!2!} \left(\frac{1}{6}\right)^{10}$$

$$P(D|I_2) = \left(\frac{1}{6}\right)^{10} (1)^6 (2)^4$$

The House never loses, And we are the House

17.7

(2025 max 01)

17045

31

a Toy problem

ideal

# Die Problem solved by Bayesian Analysis

One die



Possible events  $I = 6$  (BA)

Probabilities Normalized

for 2 theories of how a die works on trials (i.e., throws of the die)

Theory 1  $P_i = \frac{1}{6}, \sum_i P_i = 1$

Theory 2  $P_i = \frac{2 - \text{mod}(i, 2)}{9}$   
 $\frac{1}{9}$  odd ~~odd~~  
 $\frac{2}{9}$  even ~~even~~  
 $\sum_i P_i = 3(\frac{1}{9}) + 6(\frac{2}{9}) = 1$

Prion probabilities

$P(T_1 | \mathcal{C}_0) = \frac{1}{2}$   
 $P(T_2 | \mathcal{C}_0) = \frac{1}{2}$

Using principle of indifference just let them be equal since we don't know which is true.

But we could have assigned them any how and eventually the BA would've found truth.

So I did 10 throws (in 2019 actually)

Data =  $D_1$

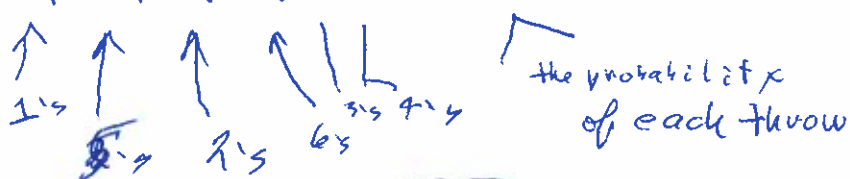
|       |   |   |   |   |   |   |   |   |   |    |
|-------|---|---|---|---|---|---|---|---|---|----|
| Throw | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| event | 1 | 5 | 1 | 1 | 2 | 5 | 2 | 6 | 3 | 4  |

Calculated Marginal Likelihoods

Marginal likelihoods

$P(D_1 | T_1 \mathcal{C}_0) = \frac{10!}{3! 2! 2! 1! 1! 1!} (\frac{1}{6})^{10} = 2.90097... \times 10^{-3}$

Need multinomial theorem in fact which there is a proof of in my notes (see p. 17035)



The number of ways this sequence could've happened summing over orders of sequence S

$P(D_1 | T_2 \mathcal{C}_0) = \frac{10!}{1! \dots} (\frac{1}{9})^3 \times (1)^6 \times (2)^1 = 6.738... \times 10^{-4}$

same  $\rightarrow$  odd event  $\rightarrow$  even events

2025 May 01

32  
17046

First Posterior odds ratio (OR<sub>Post</sub>)

Prior odds ratio (OR<sub>Prior</sub>)

$$\frac{P(T_1 | E_1)}{P(T_2 | E_1)} = \frac{P(D_1 | T_1, E_0)}{P(D_1 | T_2, E_0)} \frac{P(T_1 | E_0)}{P(T_2 | E_0)}$$

Bayes Factor  
 Bayes K Factor =  $K_{12}$

$$K_{12} = \frac{2.50057 \dots \times 10^{-3}}{6.938 \dots \times 10^{-4}} = 3.604$$

~~But this not my background knowledge~~  
 ~~$K_{12}$~~

OR Bayesian evidence  
 OR the  $\log_{10}(k)$  is called the evidence.

which is also the Posterior odds ratio, given equal probabilities in prior odds ratio

See p. 17032 for a different version

Jeffrey's Table based on some empirical testing (I think)  
 (See p. 17032)

| k        | Evidence strength |
|----------|-------------------|
| < 1      | negative support  |
| 1 - 3    | worth mentioning  |
| 3 - 10   | substantial       |
| 10 - 30  | strong            |
| 30 - 100 | very strong       |
| > 100    | decisive          |

Of course, this is just a fiducial set of values and systematic errors could nullify the evidence.

In our case theory 1 is substantially favored over theory 2.

What of the formal improvement from prior probabilities to posteriors

$$P(T_i | \mathcal{C}_\ell) = \frac{P(D_\ell | T_i \mathcal{C}_{\ell-1})}{\langle P(D_\ell | T_i \mathcal{C}_{\ell-1}) \rangle} P(T_i | \mathcal{C}_{\ell-1})$$

posterior

↗

ratio of likelihood  
to average ~~likelihood~~ marginal likelihood

marginal

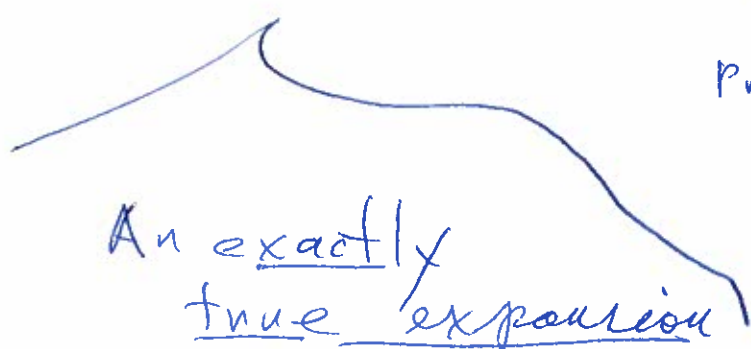
prior

$$P(D_\ell | \mathcal{C}_{\ell-1}) = \sum_j P(D_\ell | T_j \mathcal{C}_{\ell-1}) P(T_j | \mathcal{C}_{\ell-1})$$



Our best estimate of  $P(D_\ell | \mathcal{C}_{\ell-1})$

Probability of getting data  $D_\ell$  given context  $\mathcal{C}_{\ell-1}$  (background knowledge).



An exactly true expansion

if all our  $P(D_\ell | T_j \mathcal{C}_{\ell-1})$  and  $P(T_j | \mathcal{C}_{\ell-1})$

are exacty right and were averaging over all

theories  $T_j$  with nonzero probability.

Note, Normalization

~~$$\sum_i P(T_i | \mathcal{C}_\ell) = \frac{\langle \dots \rangle}{\langle \dots \rangle} = 1$$~~

and so even if the likelihoods and priors were only relative probabilities, the posteriors are normalized.

34  
17048

2029 May 01

$P(D_1 | T_1, k_0)$   
 $P(D_1 | T_2, k_0)$

In our case

$$P(T_1 | E_1) = \frac{(3.6 \dots) * \frac{1}{2}}{(3.6 \dots) * \frac{1}{2} + (1) * \frac{1}{2}}$$

$$= \frac{3.604}{4.604} \approx 80\%$$

I divided these as a simplification

$$P(T_2 | E_1) = \frac{1}{4.604} \approx 20\%$$

If my die were perfect ~~(true)~~ and not loaded,

it can be loaded

$$\lim_{r \rightarrow \infty} P(T_1 | E_r) \rightarrow 1$$

$$\lim_{r \rightarrow \infty} P(T_2 | E_r) \rightarrow 0$$

If it is a real vegas die. It can be loaded, First rule of gambling: the House never loses and we are the House

But I didn't believe this, so I did another data acquisition

Data D2

|                |   |   |   |   |   |   |   |   |   |    |
|----------------|---|---|---|---|---|---|---|---|---|----|
| throw          | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
| marginal event | 6 | 1 | 3 | 6 | 6 | 2 | 1 | 5 | 3 | 4  |

likelihoods  $P(D_2 | T_1, k_1) = \frac{10!}{3! 2! 2! 1! 2!} (\frac{1}{6})^{10} = ?$

No event of 4  $P(D_2 | T_2, k_1) = \dots (\frac{1}{9})^{10} (1)^6 (2)^4 = ?$

didn't compute these

17048

2025 may 01

$$\frac{P(T_1 | E_2)}{P(T_2 | E_2)} = \frac{P(T_1 | E_1)}{P(T_2 | E_1)}$$

OR<sub>Pre</sub> =  $K_{12}^{(1)}$   
17049

OR<sub>post</sub>

accidentally the same

3.604

This is the previous

OR<sub>post</sub> = 3.609

$$= \frac{P^{(2)}}{P^{(1)}}$$

≈ 13

strongy evidence (exp 72) 17044

Formally this is posterior Odds ratio but the distinction to Bayes factor is I think not always made

$$P(T_1 | E_2) = 0.9285 \dots$$

$$P(T_2 | E_2) = 0.07148 \dots$$

But those sequence of throws could accidentally have been unrepresentative of the true distribution. So theory 2 could still be true

I did calculate these probabilities so theory 1 is looking strongy and theory 2 is weak.

But if I just keep doing throw sequence of 10 and only true probabilistic results, Only ~~Even~~ the sequence of 10 ~~would~~ would one theory go to Probability 1.  
in the limit of infinite throws ab. 10

~~36~~  
17050

But is there another way?

Remember Ernest Rutherford's Rule  
"If you use statistics,  
you've done the wrong  
experiment."  
(p. 17025)

So instead of a sequence of throws,

I examine my die  
and see it has a 6-fold symmetry.

Now it's not perfect symmetry.

But do I care about my particular die?

No! I care about the ideal die.

The die all dies  
aspire to be.

The ideal die has exact  
6-fold symmetry.

Data  $D_3$  is exact, 6-fold symmetry

$$P(D_3 | T_1, K_2) = 1$$

$$P(D_3 | T_2, K_2) = 0$$

$$P_{B2}^{(3)} = \frac{1}{0} = \infty$$

$$\text{and } P(T_1 | K_3) = \frac{P(D_3 | T_1, K_2) P(T_1 | K_2)}{P(D_3 | T_1, K_2) P(T_1 | K_2) + 0} = 1$$

Note  
This  
a  
new  
theory  
taken  
as  
data  
— so  
data  
can be  
interpreted  
broadly

$$P(T_2 | C_3) = 0$$

17051 ~~37~~

What if  $D_r$  acquisitions were throws of 1

$$P(D_1 | T_1, C_0) = \frac{1!}{1!} \frac{1}{6} = \frac{1}{6}$$

(see p. 17045)

So the Bayesian Analysis of this trivial toy problem has found truth — about ideal dies.

$$P(D_1 | T_2, C_0) = \frac{1!}{1!} \begin{cases} \frac{1}{9} \text{ odd} \\ \frac{2}{9} \text{ even} \end{cases}$$

So the ideal Bayesian analysis would still be decisive in the limit  $\ell \rightarrow \infty$ .

(see p. 17049)

~~But if what if you -~~

In fact, this <sup>ideal</sup> BA with steps of data acquisition is not really ever done in practise — or very rarely for some special case

Rather people just collect a set of theories that that they can assign equal probability to by principle of indifference and then calculate the Bayes factor

$$B_{ij}^{(1)} = \frac{P(D | T_i, C_0)}{P(D | T_j, C_0)}$$

and they slot <sup>various</sup> ~~various~~ sets of data or combinations of data in.

See discussion on p. 17028

~~38~~  
17052

2029max01

In cosmology they try this combination of data and then that and try to figure out whether there are systematic errors in one set or another and do all kinds of finicky tests.

The problem is the  $\Lambda$ -CDM model works so well that even small systematic ~~pro~~ errors can make Bayesian analysis indeterminate in practice and that maybe where we are.

In 2025 DESI DR2 came out. In 2026 April,  $\Lambda$ -CDM is still viable in the view of many

Does DESI DR2 <sup>plus other data (2025)</sup> rule against  $\Lambda$ -CDM by 3 $\sigma$  (which is BA evidence).

Yes — if the systematics are under control. But some argue they are NOT.

But also BA in real cases requires Marginalization.

# Marginalization

a) Most theories Bayesian Analysis is applied to have free parameters that must be set ~~at~~ <sup>by the</sup> data itself

Now a theory with free parameters could be regarded an infinite continuum of theories

data interpreted broadly → it could a theory about some other aspect of reality than the ~~aspect~~ you are studying.

But that is not a useful perspective.

It's better to treat theories with free parameters as different from discrete theories and needing a different treatment to a degree.

Say you have a set of theories  $T_i(\theta)$

which perhaps ~~are~~ almost always ~~for~~ except for trivial / toy problems like the die problem

All of which are continuous variables (i.e., real numbers ... or numbers)

$\theta$  is just a symbol for the set of free parameters which can be many and can differ between the theories, but ~~we are~~ <sup>we are</sup> being general and just need a symbol  $\theta$

b) Marginalization

According to Wikipedia (Marginal distribution)  
 "marginalization" comes from summing entries  
 in a table column in a margin of the table.  
 Marginalization meant using the sums rather  
 than the entries.

For use in Bayesian analysis, "integration"  
 might be more descriptive of what  
 is done, but "integration" may be too  
 general to give the special use.

Recall from p. 17017

$$P(T_i | C_e) = \frac{P(D_e | T_i, C_{e-1}) P(T_i | C_{e-1})}{\sum_j P(D_e | T_j, C_{e-1}) P(T_j | C_{e-1})}$$

$$= \frac{P(D_e | T_i, C_{e-1})}{\langle P(D_e | T_j, C_{e-1}) \rangle} P(T_i | C_{e-1})$$

Recall the denominator is usually  
 only of formal interest (see p. 17017).

Now to go "continuum family of theories" which we  
 now call just one theory with "free parameters" which  
 are usually real or complex number variables, we expand

$$P(D_e | T_i, C_{e-1}) = \int P(D_e | T_i(\theta), C_{e-1}) P(\theta | T_i, C_{e-1}) d\theta$$

Marginal likelihood  
 or evidence

likelihood  
 = probability of  $D_e$   
 given  $T_i(\theta)$  and  $C_{e-1}$

probability  
 density of  $\theta$   
 given  $T_i$  and  $C_{e-1}$

Actually, my basic reference (Kass & Raftery 1995, p. 776) ambiguously refer to both  $P(D_e | T_i(\theta) C_{e-1})$  and  $P(\theta | T_i C_{e-1})$  as "densities".

Maybe there is choice on which is a probability and which is a probability density.

But the product  $P(D_e | T_i(\theta) C_{e-1}) P(\theta | T_i C_{e-1})$  is clearly a probability density.

In any case, I like to think of likelihood  $P(D_e | T_i(\theta) C_{e-1})$  as a probability NOT a probability density.

If  $T_i(\theta)$  is an exactly deterministic theory with exact setting (ie, these are initial conditions not included in parameters) and  $D_e$  is exact (ie, error free), then  $P(D_e | T_i(\theta) C_{e-1}) = 1$

which it may do for multiple  $\theta$  and ranges of  $\theta$ , which doesn't mean  $T_i(\theta)$  is true just that it is adequate for the data, 0 otherwise.

However, this specification is an over-idealization for practice because

- (i)  $D_e$  will be drawn from some uncertainty distribution,
- (ii)  $T_i(\theta)$  may have some uncertainty in initial conditions,
- and (iii) many, maybe most,

And there is no case of  $P(D_e | T_i(\theta) C_{e-1}) \in (0, 1)$

17056

interesting theories for Bayesian analysis will be probabilistic at least for reason (ii)

↳ This is true for cosmology where the initial density fluctuations (i.e., primordial density fluctuations (Wik)) are randomly determined by theory and a random selection is used to initial large-scale structure formation.

Thus, in practice, we must expand likelihood in an integral

$$P(D_e | T_i(\theta) C_{e-1}) = \int \underbrace{P(\dots)}_{\text{likelihood}}$$

a likelihood probability density

Then  $P(D_e | T_i(\theta) C_{e-1})$  will range over  $[0, 1]$

Calculating  $P(D_e | T_i(\theta) C_{e-1})$  for interesting cases is a major chore (i.e., bore).

$D_e$  can be petabytes (e.g., DESI final data will be a petabyte ( $10^{15}$  bytes) (Google AI))

2026 Jan 24

17057

and probably you have always  
to forward calculate from  $T(\theta)C_{e-1}$   
to  $D_e$ . I did a reverse  
calculation for my die example (see p. 17045),

but that was a very easy case.

One often uses Markov chain Monte Carlo (MCMC) to predict data as a function of inputs. A big calculation for big data.

I calculated  $P(D_e | T(\theta)C_{e-1})$  from the data  $D_e$ .

Yours truly knows little of MCMC

though I'm somewhat knowledgeable about Monte Carlo radiative transfer.

But after that, MCMC, you still have to calculate the marginal likelihood

$$P(D_e | T_i C_{e-1})$$

which usually requires another huge multi-dimensional integral

weighted by the prior probability density for free parameters

$$P(\theta | T_i C_{e-1}).$$

17058

So overall calculating  $P(D_i | T_i, C_{i-1})$   
is a huge computation  
for interesting cases.

Key point: you are marginalizing  
(i.e., integrating over  
ranges of free parameters:  
i.e., the priors).

If you know a theory is true  
(i.e., fully adequate), then  
you maximize the likelihood  
to find the best parameters  
(see p. 170 \_\_\_\_\_).

But here BA is used to try  
to find the most adequate theory  
without fine-tuning parameters  
to fit the theory to  
particular data sets  
(where you could just be fitting noise  
or systematic error).

2026 Jan 24

17059

If the Bayesian evidence ratios favor a particular theory by huge factors (e.g.,  $\geq 100$ ) for plausible priors, then that theory is the most adequate no matter what free parameters are chosen for each theory and may be the true theory.

Marginalization implements Occam's razor mathematically in that it does NOT favor theories with many free parameters just because they can be tightly fit to particular datasets using maximum likelihood.

17060

### c) Bayesian Evidence with Marginalization

This aspect is formally just as on p. 17030, but now with marginalization.

$$OR_{postij} = \frac{P(T_i|C_e)}{P(T_j|C_e)} = \frac{P(D_e|T_i, C_{e-1}) P(T_i|C_{e-1})}{P(D_e|T_j, C_{e-1}) P(T_j|C_{e-1})} = \frac{\int P(D_e|T_i(\theta), C_{e-1}) P(\theta|T_i, C_{e-1}) d\theta}{\int P(D_e|T_j(\theta), C_{e-1}) P(\theta|T_j, C_{e-1}) d\theta} \frac{P(T_i|C_{e-1})}{P(T_j|C_{e-1})}$$

↑  
 Posterior odds ratio

Bayes K factor of marginal likelihoods

OR<sub>priij</sub>  
 prior odds ratio

In one-step Bayesian analysis  $l=1$  always, but in a single paper the data sets included in  $D_e = D_1$  is often varied to see what changes. The data sets usually have varying and uncertain value.

2026 Jan 24

17061

You use Jeffreys Evidence Table  
(see p 17032) or some variation  
in order to judge the evidence.

And then, it seems, go on to  
an indecisive discussion of what  
are the best data sets and  
end conclusively/inconclusively,

Alas, competing theories for  
the cosmic scale factor  $a(t)$   
in cosmology are all very close  
to each other and the data sets  
(e.g., luminosity distances,  
angular diameter distances)

though outstanding in size and having  
been beaten down statistical  
and systematic uncertainties  
compared to the past,  
are still indecisive in Bayesian analysis

17062

in the overall judgment  
of the community as far as  
yours truly can tell  
as of 2026 May.

It seems some data sets  
are inconsistent at  
the level needed for  
decisiveness suggesting  
there are uncounted for  
systematic uncertainties <sup>(Hengt et al.,  
2026 Feb 05)</sup>  
and/or no considered  
theory is adequate.

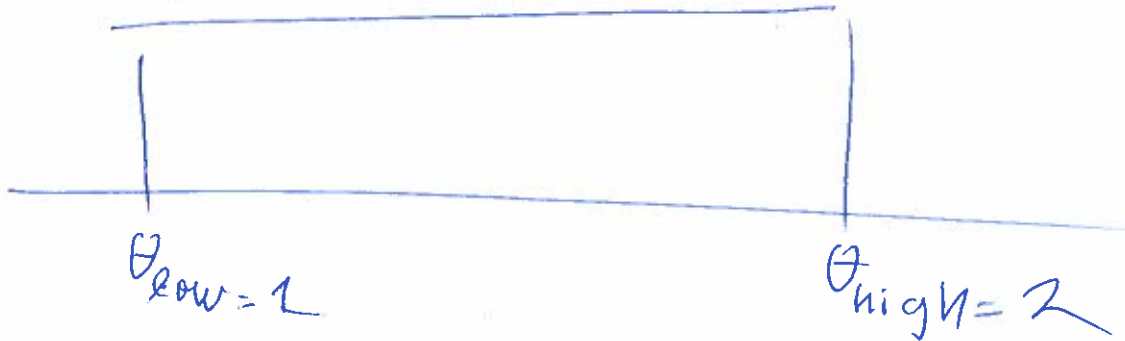
a) Priors: The Prior Probability Densities  
for the Theory Parameters

So NOT the initial prior probabilities  
for the theories  $P(T_i | C_{e=0})$  which are  
usually chosen to be all equal for interesting theories  
(by the principle of indifference), but  $P(\theta | T_i | C_{e-1})$   
from p. 17054. But choosing the  
good priors is where physical intuition is needed and  
that varies widely.

2029 May 01

17063 ~~95~~

The commonest approach is the flat prior



$$P(T_i(\theta) | \mathcal{C}_{i-1}) = \begin{cases} \frac{1}{\theta_2 - \theta_1} & \text{for } \theta \in [\theta_1, \theta_2] \\ 0 & \text{otherwise} \end{cases}$$

You can do more tricky probability densities, ~~distributions~~ but you'd have to have

a good reason which could be theoretical (i.e., based on

$\theta_1 \Rightarrow$  best guess at absolute ~~lowest~~ <sup>some</sup> value  $\theta$  could be favored theory

$\theta_2 \Rightarrow$  best guess at absolute highest  $\theta$  could be.

In fact, people estimate  $\theta_1$  and  $\theta_2$  in all kinds of ways. It is the ART of Bayesian analysis.

Choice is Qualitative Bayesian analysis.

Ok, you can do that  
but what's the problem  
then?

Doing the integral

$$\int P(D_e | T_i(\theta) k_{e-1}) P(T_i(\theta) | k_{e-1}) d\theta$$

see  
p. 17057  
discussion

You have to predict for maybe  
~~maybe~~ petabytes of  
of ~~data~~ data. data  
the likelihood  
from a multidimensional  
theory and then do a multidimensional  
integral

→ It was easy for  
the die problem (see p. 17045),  
but it's not easy for  
e.g., ~~for the~~ DESI DR2.

Need to use advanced multi-dimensional  
integral methods  
— probably Markov Chain Monte Carlo (MCMC)  
(with which I've no experience  
though I know something  
about Monte Carlo Radiation  
-transfer)

So, in fact, state of art Bayesian analysis on huge data is enormous computationally

Probably, Machine learning is helping.

Also there are ~~simplex~~ quick methods of estimating Bayesian evidence: e.g.

BIC = Bayesian information criterion

AIC = Akaike Information criterion

which I think well approximate Bayesian evidence in certain useful limits but I know little of them.

e) The Bayesian Evidence

~~R~~<sub>ij</sub> =  $\frac{P(T_i | K_e)}{P(T_j | K_e)}$  see p. 44

You've marginalized over the parameters. The evidence is based on how good the theory is no matter how many parameters there are and (if your priors are good) no matter what the ~~parameters~~ values

48 | 17066 |

Because of the slop  
in the det of priors  
and systematic errors,  
No one is impressed  
by evidence  $K_{ij} \geq 10$   
maybe not even  $K_{ij} = 100$ ,  
but if it's  $K_{ij} = 1000$ ,  
then theory  $T_i$  is probably  
a lot better than  $T_j$   
and  $T_j$  is probably  
ruled out,  
unless you have  
made bad mistakes.

## A) Maximum Likelihood

This is well known technique  
for determining the best/likeliest  
parameters for a theory  $T_i(\theta)$ .  
If theory is NOT true/adequate,  
the parameter values determined may  
not be very good or they may have  
no meaning relative to the true/adequate theory

2026 Jan 24

17067

To find the maximizing parameters  
solve

$$\frac{\partial P(D_{\ell} | T_i(\theta) C_{\ell-1})}{\partial \theta} = 0$$

for  $\theta$ .

But this is a symbolic equation  
since  $\theta$  is a set of parameters.

So  $\theta \rightarrow \{\theta_n\}$  where there  $N$  parameters

You  
must  
solve

$N$   
simultaneous  
equations  
in  
general

$$\frac{\partial P(D_{\ell} | T_i(\dots, \theta_n, \dots) C_{\ell-1})}{\partial \theta_n} = 0$$

And in nontrivial cases

$P(D_{\ell} | T_i(\dots, \theta_n, \dots) C_{\ell-1})$  is

NOT analytic, and so you  
forced to use some numerical method  
(e.g., stochastic gradient descent,  
which is analogous to  
stochastic gradient descent)

17068

You search up gradients,  
but jump around stochastically  
so as NOT to get trapped  
in local maxima.

For large data sets  $D_e$ ,  
this can be computationally  
demanding.

And systematic error  
on data sets unrepresentative  
of the true distribution  
can lead to bad parameter values

And also if your theory is  
very inadequate, but has  
a large number of parameters,  
your likeliest parameters could  
be just giving a good fit  
to noise

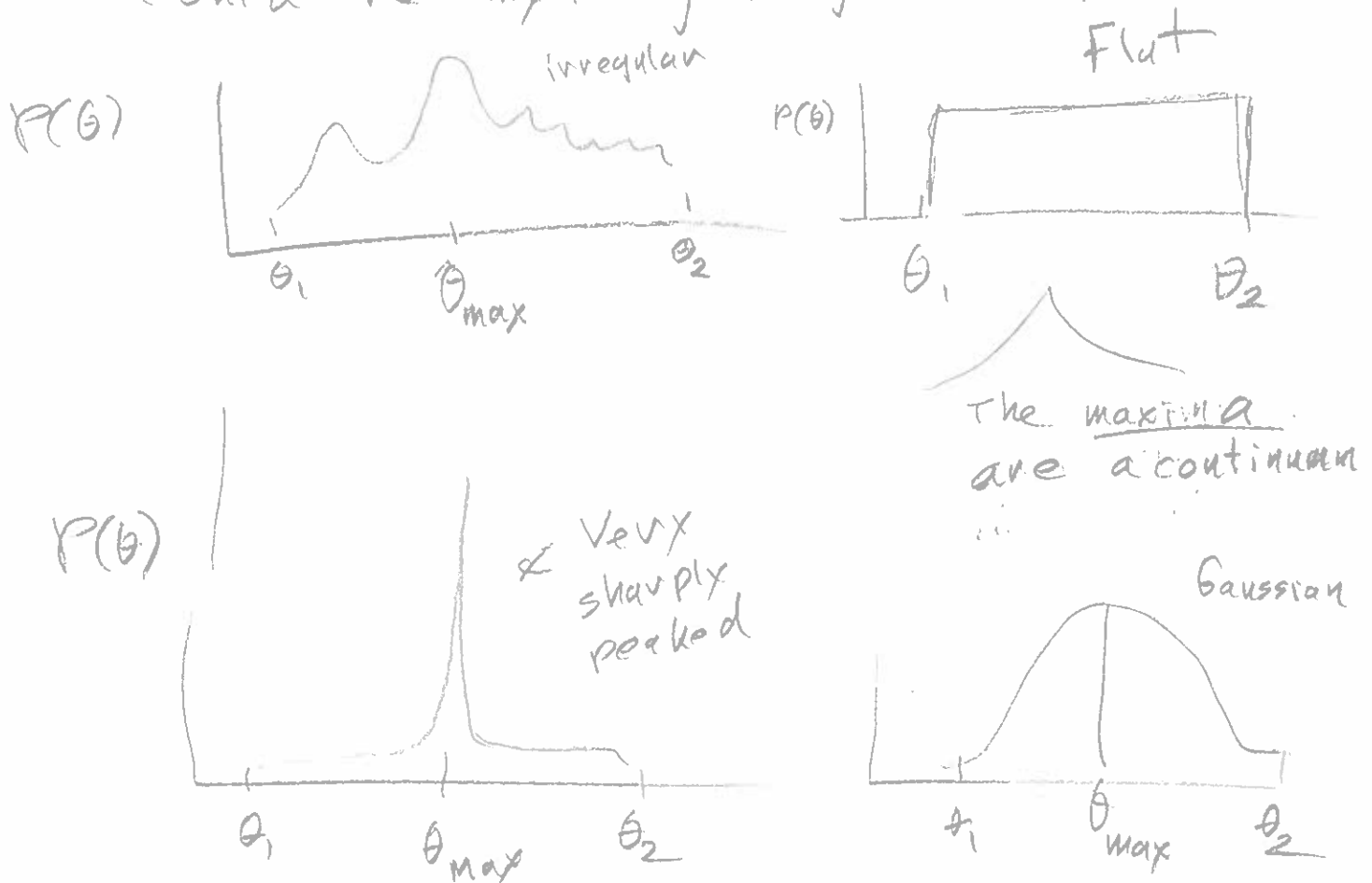
However, if your theory is true/adequate  
and  $D_e$  is sufficiently large, the  $\{\theta_n\}$   
must asymptotically approach the true values  
since if you have all possible data the distribution  
of data must be the true distribution,

2026 Jan 24

17069

Note,  $P(D_2 | T_1(\theta) C_0) = P(\theta)$

could be anything in general,



So you may or may not have an easy job finding  $\theta_{max}$

and  $\theta_{max}$  may or may not well approximate  $\theta_{true}$

if your theory is true/adequate.

Of course, you should follow Rutherford rule and try to get decisive data that allows you to find  $\theta_{max}$  most easily

