

Bayesian Analysis (BA)

(A name I prefer to the usual Bayesian Inference)

- 1) Preamble: BA is path to truth quantified or 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)
~~to Truth~~ 5) Bayes Odds Ratio & Bayes Factor
- 7) Multinomial Theorem & Multinomial Probability (p. 21) (p. 23)
- 8) An Example of Bayesian Analysis
A toy example: The Dice Problem
- 9) Marginalization & Occam's Razor
shortcuts to Bayesian analysis
- 10) BIC = Bayesian Information Criteria
AIC = Akaike Information Criterion

7002

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

But then our computers

packages

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

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

act intuitively, except maybe the theory of everything or 2nd law of thermodynamics

E.g., Newtonian Physics,

It is believed to be the macroscopic, low velocity, weak gravity limit of Special relativity and Newtonian physics

Another perspective is as an approximate theory

absolutely exact, but maybe there is quibbling

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

in everyday life.

Other examples 2nd law

17003
Maybe
exactly what
Well the
system must be
big enough

of Thermodynamics

(a law in qualitative sense by no means)

Natural selection evolution

↳ unlike Newtonian Physics
these can be proven

by mathematical logic (not law? Well
not any toy which
you
buy)

— and so can Bayesian Analysis

↳ aside from pure rigorist 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 ^{things} 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

Next
time
you
think
of
stakes
relative

Then new experience happens
in that situation
and you vaguely update

~~your~~ priors
to your posteriors
(which sounds
awful and often
(s))

the
posterior
probabilities

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 impacts
can change the rules of success
suddenly → like what happened to
non-avian dinosaurs

Ex. Job interviews

~~17005~~
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

improvement if just

restricted to

each new experience

~~of a unique situation~~

~~from~~ has very restricted generality.

~~But~~

But as ^{more} general conclusions about updates occur

↳ qualitative Bayesian analysis approaches the scientific method.

7006

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

7.3)

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

tdo 2677
3
2
1
of

Consider 3 events A, B, K.
We have a diagram / knowledge that
I introduce an extra term because
I need it later. It isn't needed for ^{the} proof
Bayesian theorem itself.

I use $A \cup B \cup K$ as joint event as an
shorthand for union = \cup symbol
in math
which is too klutzy for no in this
context

order
has no
meaning.

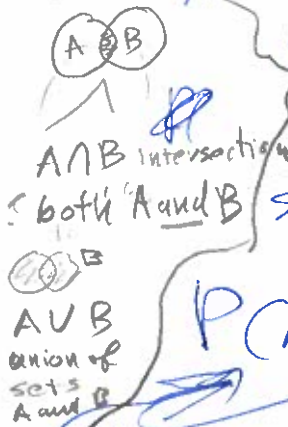
2025 May 01

17007

Conditional Probability, and Probability Product Rule

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

ABK is ~~joint event~~ ^{intersection} all 3 states BAK Background Knowledge that is BA jargon is often called conditioning knowledge on context, (Lorato 2014 p. 3)



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

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

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

Bayes theorem in symmetrical form

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

Asymmetric form

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

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

Some assert this is an axiom or ^(truth?) can be so taken. But I believe ~~probabi~~ only has meaning in a frequentist sense even if N_{ABK}/N is only vaguely known or hypothesized imaginative. In fact, BA means you can improve from vaguely known to "truth". N_{ABK} & N exist in some ideal limit as Galileo would argue.

7008

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

Root of Bayesian analysis

Bayesian analysis was greatly elaborated in 1999 OK.
(Scientific Inference 3rd ed 1973)
by Harold Jeffreys (1891-1989)
(no relation - but he manages to look like my father anyways)

What is Bayesian analysis?

illy, clual, dke, ideal, ps, xonal, p 1, = not, ne, enyou, ems to, int step 1, th the, re that, ep 1's, ll, actually, et, 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, hard, because there is a lot of stop 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.

No one much used Bayesian analysis before 1990s, since then it has had growing vogue.
one in single pens. New are d this data, now s anal take way that

Why was it unused before?

17009

Well, ~~it~~ ^{its} ~~real~~ 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 ^{not ~~working~~ leaving aside partially}

ive theories

or having some truth and some falseness definitely in

less absolute sense ^{not leaving aside partially true theories}

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

It's probability to own knowledge

Example

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

It's one or the other in absolute sense.

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

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

2026 Jan 27

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 Loredot & Wolport (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 on adding to C_0 . The set $\{T_i\}$ does NOT have to be complete even to C_0 as it is, but in practice

7012)

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, 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

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

2026 Jun 24

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 $\{T_i\}$

initial probabilities
= initial priors

$$P(T_i | C_0)$$

conditional probability of T_i given context C_0

theory T_i

Context

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!

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

7017)

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., $P(T_i | C_e)$)

[2026 Jan 24]

[17015]

which are your priors for step $l+1$.

c) Updating Priors to Posteriors
using Bayes theorem at step l

Consider joint probability and factorize

$$\begin{aligned}
 & P(T_i, C_{l-1}, D_l) \\
 & \quad // \quad \quad \quad // \\
 & P(T_i | C_{l-1}, D_l) P(C_{l-1}, D_l) \quad \quad P(D_l | T_i, C_{l-1}) P(T_i | C_{l-1}) \\
 & = P(T_i | C_l) P(D_l | C_{l-1}) P(C_{l-1}) \quad = P(D_l | T_i, C_{l-1}) P(T_i | C_{l-1}) P(C_l)
 \end{aligned}$$

$$\therefore P(T_i | C_l) P(D_l | C_{l-1}) = P(D_l | T_i, C_{l-1}) P(T_i | C_{l-1})$$

$P(C_{l-1})$
cancel out, but
they are also 1
since we have

(i.e. know)
 C_{l-1}

Posterior
for step l

which is just the
symmetric Bayes theorem
from p. 17017 recovered

$$P(T_i | C_l) = \frac{P(D_l | T_i, C_{l-1}) P(T_i | C_{l-1})}{P(D_l | C_{l-1})} \left\{ \begin{array}{l} \text{Prior} \\ \text{for} \\ \text{step } l \end{array} \right.$$

$P(D_l | T_i, C_{l-1})$ is the
marginal likelihood of the data D_l

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 have 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

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

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

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

But
this is 1
since
we
know it.

really just a best estimate to your 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 priors $P(T_i | C_{e-1})$ are automatically normalized as noted on p. 17019:
 $\sum_i P(T_i | C_e) = 1$.

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

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

Weighted mean of the marginal likelihood

Except for format niceness

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

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

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

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.

7018 |

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_e
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

17019

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 |

✓) 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.

i) 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 transformation. Famously, Einstein was led to special relativity by, among other things, the belief that

7022

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.

Follow a herd, 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 29

17023

vii) Does BA guarantee that you will get to a true theory even if only limited to vast C_2 ?

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 random, 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;

7024

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

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

(Wiki: 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?

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

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

$$= \frac{10^{6000}}{10^{18}} \text{ s} \approx 10^{6000} \text{ s} * \left(\frac{1 \text{ year}}{\pi * 10^8 \text{ s}} \right)$$

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

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 jay 24

This example is a pathological case, (1702)
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)

2026 Jan 27

17029

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 non-rigorously is a valid justification

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

7028

17.5 Bayesian Analysis: What People Actually Do
Including Bayes Odds Ratios, and
Bayes Factor But Deferring
Marginalization to p. 170

So far as I know, no one in practical work does more than one explicit BA step of ideal BA in a paper. } One-step Bayesian analysis

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 datasets.

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.

17029

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
relative to adequacy of theories in your set $\{T_i\}$.

(7030)

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

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

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

From p. 17017

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

$K_{ij}^{(l)}$ The (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

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

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

= evidence for T_i / evidence for T_j

evidence is synonym for marginal likelihood (Wik: 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 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

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 fiducial table.

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

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

$K_{ij}^{(e)}$	$\log(K_{ij}^{(e)})$	Significance of Evidence
< 1	< 0	Negative (in favor of T_i T_j)
$1 \text{ a } 3.2$	$0 \text{ a } 0.5$	Insignificant
$3.2 \text{ a } 10$	$0.5 \text{ a } 1$	Significant
$10 \text{ a } 100$	$1 \text{ a } 2$	Strong
> 100	> 2	Decisive

$$K_{ij}^{(e)} = \frac{P(D_e | T_i, C_{e-1})}{P(D_e | T_j, C_{e-1})}$$

where usually $C_{e-1} = C_0$

$$\text{and } D_e = D_1$$

for one-step Bayesian Analysis

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 σ 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), one must choose priors on the free parameters of theories and

7034

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.

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

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

marginalization

See p. 170

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

17.6

Multinomial Probability Distribution & The Multinomial Theorem

1) 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 indist.

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 jarg. (With some go in down the row)

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

You get sequences: e.g.

- seq 1 $P_1 P_1 P_1 \dots P_1 = P_1^N$ probability of this one sequence
- seq 2 $P_1 P_2 P_1 \dots P_1 = P_1^N P_2$
- seq 3 $P_2 P_1 P_1 \dots P_1 = P_1^N P_2$

These are 2 sequences of equal probability

interchanging the identical subscripts does not give a new sequence

permuting interchanging the 1 and 2 indices of sequence 2 gives sequence which is different 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 $\{n_i\}$

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-identity.

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.

2026 Jun 24

1703

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

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

N elements in the selection

$\therefore N!$ 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 $n_2 = N - n_1$ and $C(\{n_i\}) = \frac{N!}{n_1!(N-n_1)!} = \binom{N}{n_1}$ which is the binomial coefficient.

and also the multinomial coefficient of the multinomial theorem.

7038

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.

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

When you actually do the multinomial expansion, it is clear (maybe) that the probability of configuration $\{n_i\}$ being obtained is $P(\{n_i\}) = C(\{n_i\}) \prod P_i^{n_i}$

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

There are the counts of combinations conforming to set $\{\epsilon_{n_i}\}$

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

which must be an integer

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

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

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

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

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

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

Probability of getting any one of the combination

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

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

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 quantum identical not classical identical

26
7070

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

Normalization?

Proof

$$1 = \sum_i p_i$$

represented by

which we demand for probability particles with 0 spin

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

represented by Binomial theorem of 2 reference of N

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

multinomial theorem

With polynomial one coefficient and variable $ax^2 + bx + c$ etc. (WK)

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

for higher multinomial

~~the formulae~~ are harder to get formulae easy formulae for. Are no easy ones & that.

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

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

2025dec05

27
17041

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

Unit: Digression

9) Further Digression on Binomial Theorem

Binomial Theorem
Special case: Binomial theorem & Binomial probability distribution,
Prove by induction if afflicted by paranoia.

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

for $q = 1 - p$

$$= \sum_{k=0}^n \binom{n}{k} p^k q^{n-k}$$

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)

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

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

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

$$\sum_{k=0}^n k^2 \binom{n}{k} x^k q^{n-k} \Big|_{x=p} = \left(x \frac{d}{dx} \right)^2 (x + q)^n \Big|_{x=p}$$

the n choose k formula

$$b) l=0$$

$$M_0 = 1$$

~~218~~
1702

$$l=1 \quad M_1 = \left. \frac{d}{dx} n(x+q)^{n-1} \right|_{x=p} = np$$

$$l=2 \quad M_2 = \left. \frac{d^2}{dx^2} [n(x+q)^{n-1}] \right|_{x=p}$$
$$= np \left[\frac{d}{dx} (x+q)^{n-1} + n(n-1)(x+q)^{n-2} \right]_{x=p}$$

$$= np [1 + p(n-1)]$$

$$= np + p^2 n(n-1)$$

$$\sigma^2 = \langle (x - \bar{x})^2 \rangle = \langle x^2 - 2x\bar{x} + \bar{x}^2 \rangle$$

$$= \langle x^2 \rangle - \bar{x}^2 = np + p^2 n(n-1) - p^2 n^2$$

$$= \cancel{p^3 n(n-1)} = np - p^2 n$$

$$= np(1-p)$$

Bev-53
correct

Probably Now totally Worthless

2021 Nov 22

$P_i = \frac{1}{I}$ $\sum P_i = 1$ or 2 skewed
 This simplified Navette and got $\{n_i\}$

1709

Theory 1

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

$\sum n_i = N$

This is for I even issue $P_i < 1$ $i < 2$ Total = $\frac{I}{2} + \frac{I}{2} = \frac{3}{2}$

Theory 3

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

As so often in Astronomy, rod has told is there are only 2

Theory 2

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

possible true theories

Very easy to post dict (PGDIT) not like Jupyter system with huge data

Die Problem for Zeus II

Normal problem

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



Theory 2 $P_i = \frac{2 - \text{mod}(i, 2)}{6}$

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

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

Principle of indifference (Barnes p. 6)

So throw complex enough, but not decide?

Really did this in 2019

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

$P(D|T_1) = \frac{10!}{3! 2! 4! 1! 1! 1!} \cdot \left(\frac{1}{6}\right)^{10} = 2.50057 \dots \times 10^{-3}$

$P(D|T_2) = \dots \left(\frac{1}{3}\right)^{10} 1^6 \cdot 2^2 = 6.938 \dots \times 10^{-4}$

mat 6 1/2 any 4 2

30 $P(k) = \frac{P(D|H_1) * P(H_1) + P(D|H_2) * P(H_2)}{P(D|H_1) * P(H_1) + P(D|H_2) * P(H_2)}$

44

$k = 3.609$

Probability given your knowledge - Before even No assumptions

So Prior odds = 1

Bayes factor $k = 3.609$

Posterior odds = 3.609

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

eg fun calculator (using)

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

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

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

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

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

$P_1(1) = 0.7828$, $P_2(2) = 0.21799$ to posterior \rightarrow 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		

$P(D|T_1) = \frac{10!}{3!2!2!1!2!} (\frac{1}{6})^{10}$

$P(D|T_2) = \dots (\frac{1}{6})^{10} (\frac{1}{6})^6 (2)^4$

The House never loses, And we are the House

17.7

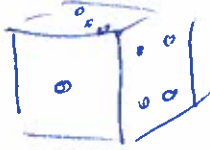
(2025max01)

17045 (31)

a Toy problem

Die Problem Solved by Bayesian Analysis

One die



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

Possible events $I = 6$

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

Theory 2 $P_i = \begin{cases} \frac{2 - \text{mod}(i, 3)}{9} \\ \frac{1}{9} & \text{if } i \text{ is } 1 \text{ or } 4 \\ \frac{2}{9} & \text{if } i \text{ is } 2 \text{ or } 5 \end{cases}$
 $\sum_i P_i = 3(\frac{1}{9}) + 6(\frac{2}{9}) = 1$

Prior probabilities

$$P(T_1 | C_0) = \frac{1}{2}$$

$$P(T_2 | C_0) = \frac{1}{2}$$

Using principle of indifference just let them equal since we don't know which is true. But we could have assigned them any way and eventually the BA would've found truth.

So I did 10 throws (in 2019 actually)

Data = D_4

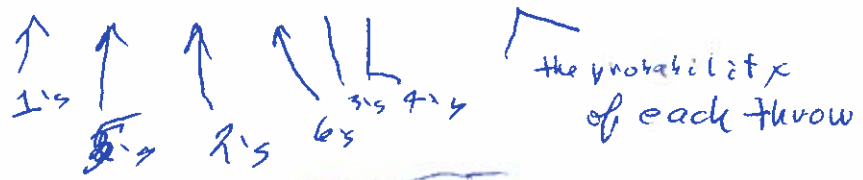
Throw	1	2	3	4	5	6	7	8	9	10
event	1	5	1	1	2	5	2	6	3	4

Calculated Marginal Likelihoods

likelihoods

$$P(D_1 | T_1, C_0) = \frac{10!}{3! 2! 2! 1! 1! 1!} \left(\frac{1}{6}\right)^{10} = 2.90097 \times 10^{-8}$$

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 sequences

$$P(D_1 | T_2, C_0) = \frac{10!}{9!} \left(\frac{1}{3}\right)^{10} \times (1)^6 \times (2)^4 = 6.738 \times 10^{-8}$$

= 86 (odd event) (even event)

32
17046

(2025/01/01)

First Posterior odds ratio (OR_{Post})

Prior odds ratio (OR_{Prior})

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

Bayes Factor
Bayes K factor = $K_{1,2}$

~~$K_{1,2}$~~ = $\frac{2.50057 \dots \times 10^{-3}}{6.938 \dots \times 10^{-4}}$

~~But this not
my background knowledge~~
 ~~$K_{1,2}$~~

= 3.604

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)

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.

(2029 may 01)

17047

What of the formal improvement from prior probabilities to posteriors

$$P(T_i | \mathcal{E}_e) = \frac{P(D_e | T_i | \mathcal{E}_{e-1})}{\langle P(D_e | T_i | \mathcal{E}_{e-1}) \rangle} P(T_i | \mathcal{E}_{e-1})$$

posterior ratio of likelihoods to average likelihood prior

$$P(D_e | \mathcal{E}_{e-1}) = \sum_j P(D_e | T_j | \mathcal{E}_{e-1}) P(T_j | \mathcal{E}_{e-1})$$

Our best estimate of $P(D_e | \mathcal{E}_{e-1})$

Probability of getting data D_e given context \mathcal{E}_e (background knowledge).

An exactly true expansion

if all our $P(D_e | T_j | \mathcal{E}_{e-1})$

and $P(T_j | \mathcal{E}_{e-1})$

are exactly right and were averaged over all

theories T_j with nonzero probability.

Note, Normalization

$$\sum_i P(T_i | \mathcal{E}_e) = \frac{\langle \dots \rangle}{\langle \dots \rangle} = 1$$

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 them 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

$$\lim_{L \rightarrow \infty} P(T_1 | E_L) \rightarrow 1$$

$$\lim_{L \rightarrow \infty} P(T_2 | E_L) \rightarrow 0$$

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
event	6	1	3	6	6	2	1	5	3	5

likelihoods

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

$$P(D_2 | T_2, k_1) = \dots (\frac{1}{7})^{10} (1)^6 (2)^9 = ?$$

but that is common to both likelihoods

16 just on top

didn't compute there

2025 may 01

~~75~~
17049

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

OR_{Pre} = $K_{12}^{(1)}$

OR_{post}

accidentally
the same
3.604

This is the
previous OR_{post}
= 3.609₂

$$= \frac{B_{12}^{(2)}}{B_{12}^{(1)}}$$

13 strong evidence (exp 17044)

Formally this is posterior Odds ratio, but the distinction to Bayes factor is I think not always nothing

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

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

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

But these sequence of throw could accident have been unrepresentative of the true distribution. So theory could still be true.

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

36
17050

But is there another way?

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

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.

Defn. D_3 is exact 6-fold symmetry

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

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

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

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

ote
his
w
may
ken
L
fd
so
itu.
u
re
terpreted
readly

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

17051 ~~27~~

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

~~But if what if you -~~

In fact, the ^{ideal} BA with
steps of data acquisition
is not really ever done
in practise — or very
rarely for some special case

See
discussion
on
p.17028

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

$$R_{ij} = \frac{P(D_i | T_i; C_0)}{P(D_i | T_j; C_0)}$$

and they slot ^{various} ~~various~~ sets
of data or combinations
of data in.

38
7052

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 Λ -CDM model works so well that even small systematic ~~pro~~ errors can make Bayesian analysis indeterminate in practice and that maybe where we are - plus other data (2025)

In 2025 DESI DR2 meet.
In 2026 April, Λ -CDM is still viable in the view of many

Does DESI DR2 ^{plus} rule against Λ -CDM by 3 σ (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~~ ^{by the} data itself

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

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

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 $T_i(\theta)$

except for trivial toy problems like the die problem

θ is just a symbol for the set of free parameters which can be many and can differ between the theories, but ~~we are~~ ^{we are} being general and just need a symbol

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

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

$$\begin{aligned}
 P(T_i | C_e) &= \frac{P(D_e | T_i, C_{e-1}) P(T_i | C_{e-1})}{\sum_j P(D_e | T_j, C_{e-1}) P(T_j | C_{e-1})} \\
 &= \frac{P(D_e | T_i, C_{e-1})}{\langle P(D_e | T_j, C_{e-1}) \rangle} P(T_i | C_{e-1})
 \end{aligned}$$

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

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

$$\underbrace{P(D_e | T_i, C_{e-1})}_{\text{marginal likelihood or evidence}} = \int \underbrace{P(D_e | T_i(\theta), C_{e-1})}_{\text{likelihood}} \underbrace{P(\theta | T_i, C_{e-1})}_{\text{probability density of } \theta} d\theta$$

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

probability
density of θ
given T_i and C_{e-1}

Actually, my basic reference (Kass & Raftery 1995, p. 176) 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 (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 θ and ranges of θ , which doesn't mean $T_i(\theta)$ is true just that it is adequate for the data

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,

0 otherwise

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 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

17097

and probably you have always
forward calculate from $T(\theta_i)C_{e-1}$
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.

Yours truly knows little of MCMC
though I'm somewhat knowledgeable
about Monte Carlo radiative transfer.

But after that 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

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

7048

So overall calculating $P(D_e | T_i; C_{e-1})$
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 _____).

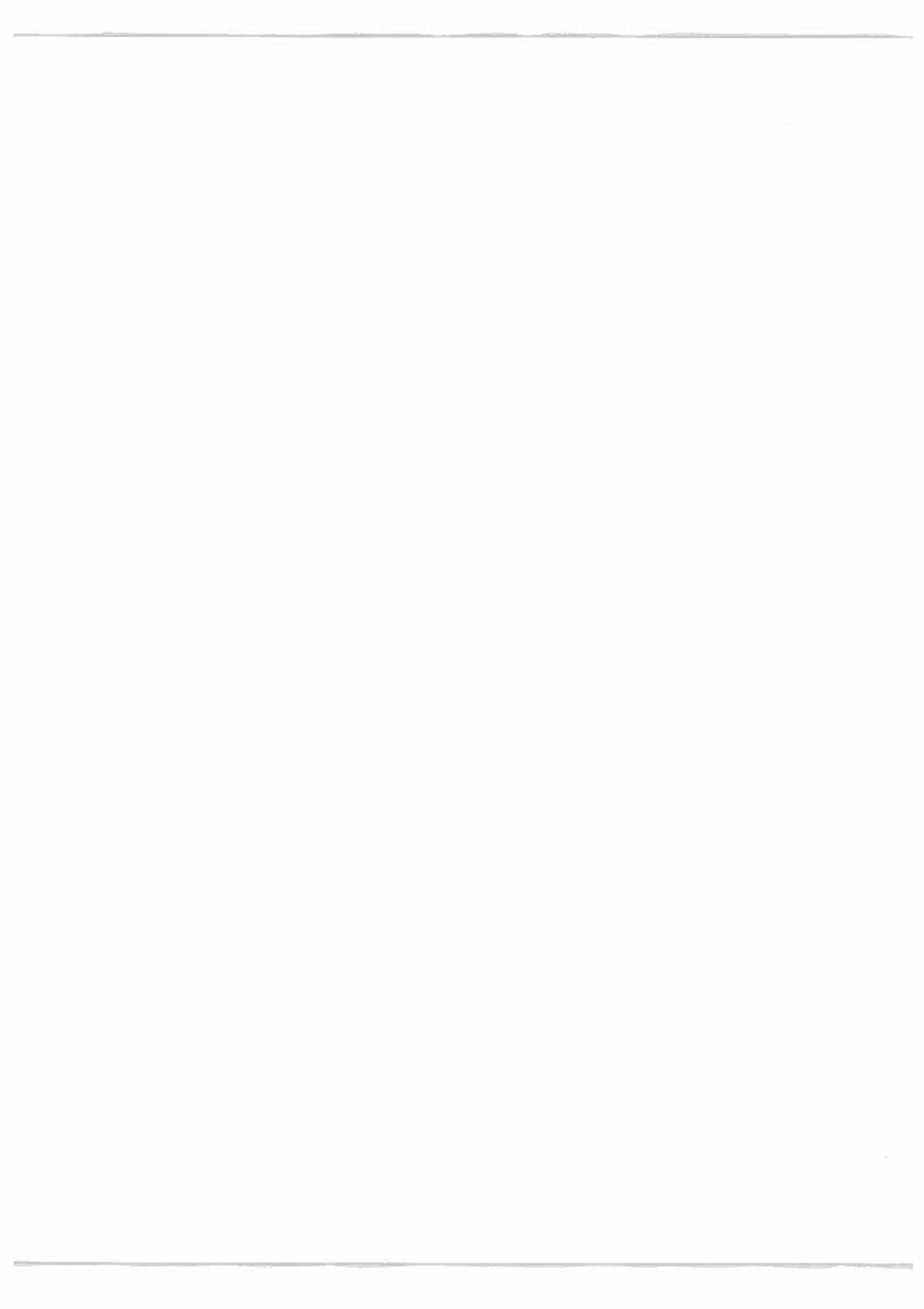
But here BA is used to try
to find the most adequate theory
without fine-tuning parameters
to fit the theory to
particular data sets
(were you could just be fitting noise
or systematic error).

2026 Jan 24

17059

If the Bayesian evidence ratios favor a particular theory by huge factors (e.g., ≥ 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.





b) A word about likelihoods with free parameters

edu 124
of likelihood
if probability
it is
- come
on Bayes
M.
ep12
+ is NOT
probability
density

$$P(D_n | T_i(\theta) K_{n-1})$$

New data

theory T_i
with its set
of free parameters

Background
knowledge
which
~~is~~
is usually
left
implicit

What if you know
or assume

theory $T_i(\theta)$ is
true, but you don't
know its parameter values

See p 49
for more
on maximum
likelihood

Take general rule

is. Maximize the likelihood

Symbolically $\frac{\partial}{\partial \theta} P(D_n | T_i(\theta) K_{n-1}) = 0$

and solve for set θ

Remember the theory may just be
statistical

and also the data D_n may
statistical error

but we assume no systematic error
(but in practice there can be)

With potentials
of data
many
parameter
things
may be
a super-
computer
calculator
eg.
stochastic
gradient
descent

2025 May 01

(41)

But because of the statistical aspect
the D_e may be unrepresentative
of the true parameter
values.

So you can get very bad
values

But if D_e — { large } → infinite
 { decisive } all
then eventually data
you will
get the true parameter values.

But again if your theory is
not right the
parameter values may
be wrong
or meaningless in that
they may NOT be the
true parameter.

Also if your theory has many
parameters, you may simply
get an excellent fit
to a set of data
and excellence argues for it
being right → but this
may wrong → you may be
fitting noise for
example.

92

marginalization
 ideally doesn't worry
 about ~~best~~ fitting parameters,
 but about the truth of
 theories,
 in fact, marginalization
 is marginalizing over
 parameter values

In a real way, Marginalization
 embodies Occam's razor
 putting theories on an equal
 basis no matter how
 many free parameters
 they have.

c) Marginalization

For theories $T_i(\theta)$, instead
 of prior probabilities you
 have prior probability densities.

$dP(T_i | k_{e-1})$

~~$dP(T_i | k_{e-1})$~~ = $P(T_i(\theta) | k_{e-1}) d\theta$

Respite what
 Laredo 2024 n.2
 say,
 They misspeak.

The likelihood is
 still a probability
 NOT a probability
 density!

differentiated
 symbolizing of
 all parameters

$dP(T_i | k_e) = P(T_i(\theta) | k_e) d\theta$

$= \frac{P(D_e | T_i(\theta) k_{e-1}) P(T_i(\theta) | k_{e-1}) d\theta}{\sum_{j_m} P(D_e | T_j(\theta_m) k_{e-1}) P(T_j(\theta_m) | k_{e-1}) d\theta_m}$

The sum on m can just be turned in to integrals

$$P(T_i(\theta) | K_x) d\theta = \frac{P(D_e | T_i(\theta) K_{e-1}) P(T_i(\theta) | K_{e-1}) d\theta}{\sum_j \int P(D_e | T_j(\theta) K_{e-1}) P(T_j(\theta) | K_{e-1}) d\theta}$$

mean likelihood again.

Different integral but sum economized on symbols and indic.

So in an idealized BA, you could get ~~reconst~~ constructed updated probability density in a series of steps but that's probably never done.

One just integrated over θ

~~$P(T_i(\theta) | K_e)$~~
have symbolically

$$P(T_i | K_e) = \frac{\int P(D_e | T_i(\theta) | K_{e-1}) P(T_i(\theta) | K_{e-1}) d\theta}{\langle \text{mean likelihood} \rangle}$$

θ integrated away

in fact, the series stops usually at $l=1$, (always)

$\sum P(T_i | K_e) = 1$
'normalized automatically'

94)

And who cares about the denominator.

We just want the evidence ^{Bayesian}

$$B_{ij} = \frac{P(T_i | K_e)}{P(T_j | K_e)} = \frac{\int P(P_e | T_i(\theta) | K_e) P(T_i(\theta) | K_e) d\theta}{\text{same for } T_j}$$

Integrate away \rightarrow marginalize the parameters

Really this Bayes posterior odds ratio in some ~~in the~~ old-fashioned formulae sense (Kass & Raftery 1995)

but I think everyone just calls it the Bayes factor or Bayes evidence or the log of it the Bayes evidence.

d) Pivone

What is the difficulty since everything looks simple formally.

What is $P(T_i(\theta) | K_{e-1})$?

usually ~~is~~
 $e-1 = 0$

What are the priors?

2029 May 01

95

The commonest approach is
the flat prior



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

You can do
more tricky
distributions
but you'd
have to have
a good reason

$\theta_1 \Rightarrow$ best guess at absolute lowest
value θ could be

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

In fact, people estimate θ_1 and θ_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

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

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

→ It was easy for
 the die problem (see
 p. 31),
 but it's not easy
 for the DESI DR2.

Use advanced multi-dimensional
 integral methods

— probably Markov Chain Monte Carlo

(with which I have no experience
 though I know something
 about Monte Carlo Radiative
 transfer)

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

[47]

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) ~~But~~ The Bayesian Evidence

$$B_{ij} = \frac{P(T_i | K_e)}{P(T_j | K_e)} \quad \text{see p. 44}$$

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 ~~parameter~~
parameter values are.

48

Because of the slop
in the out of priors
and systematic errors,
No one is impressed
by evidence $B_i \approx 10$
maybe not even $B = 100$,
but if it's $B_i = 1000$,
then theory T_i is probably
a lot better than T_j
and T_j is probably
ruled out.

A) Maximum Likelihood Redux

What is it from a Bayesian analysis perspective? A Mythical interpretation discussion

Recall from p. 43

$$P(T_i(\theta) | K_e) d\theta = \frac{P(D_e | T_i(\theta) | K_{e-1}) P(T_i(\theta) | K_{e-1})}{\int P(T_i(\theta) | K_e) d\theta}$$

mean likelihood

$$\sum_i \int P(T_i(\theta) | K_e) d\theta = 1 \text{ normalized automatically}$$

But say you take T_i as true and just normalize

$$1 = \int P(T_i(\theta) | K_e) d\theta = A \int P(D_e | T_i(\theta) | K_{e-1}) P(T_i(\theta) | K_{e-1}) d\theta$$

Then you could read ^{normalization constant} $\times P(T_i(\theta) | K_{e-1}) d\theta$

$$P(\theta) = P(T_i(\theta) | K_e) = A P(D_e | T_i(\theta) | K_{e-1}) P(T_i(\theta) | K_{e-1})$$

as probability density for obtaining ~~for obtain parameters~~

Regard D_e as a measurement of θ parameter values θ .

50] But what is the measurement for obtaining the actual

θ_{true} ?

Obtaining ~~all~~ K_{max}

All knowledge relevant to the theory,

But what does the probability density describe?



Maybe an ensemble of universes where T_i is true in

~~every~~ ^{all of them}

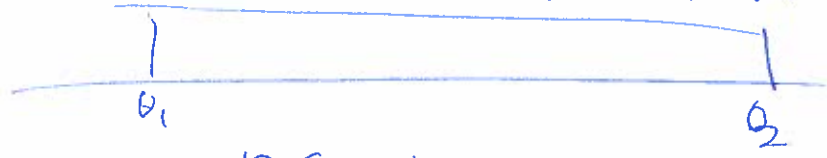
but with ~~a~~ continuum of varying θ ,

You don't know what universe you are in until you get K_{max}

$P(\theta)d\theta$ is probability for ~~your~~ your universe being ~~in~~ in $d\theta$.

~~Posterior~~ probability density

Assume known $P(T_i(\theta) | K_{e-1})$ is flat between bounds on θ as a ~~simple~~ sweeping simplification

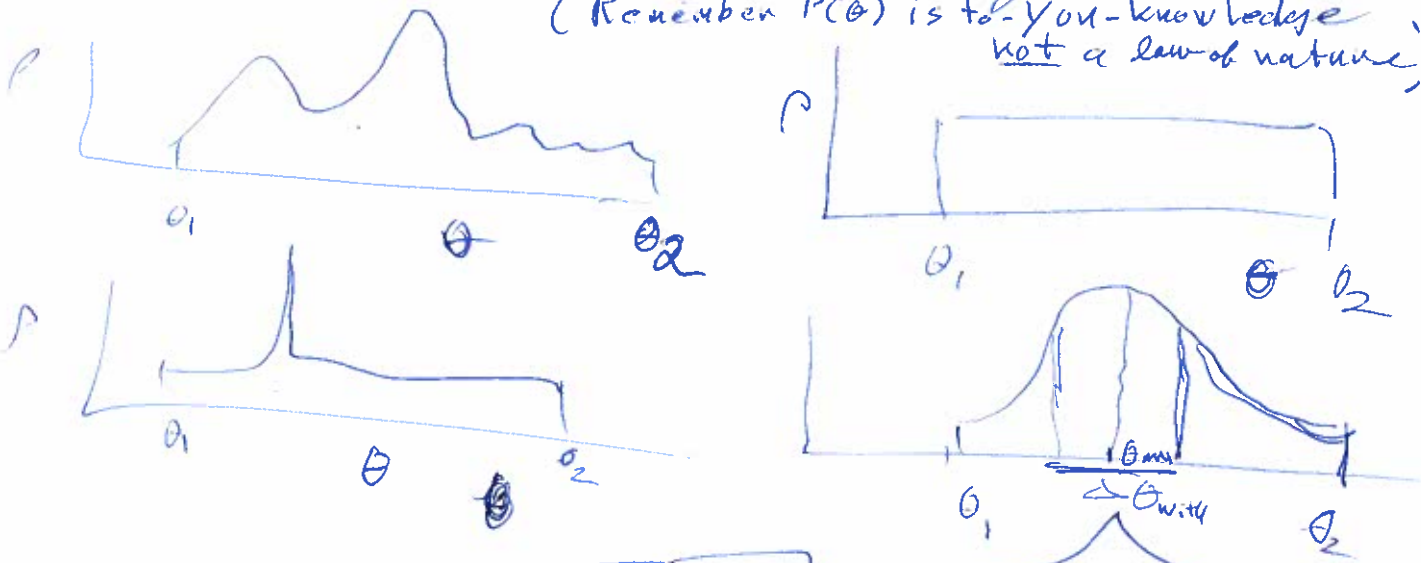


$$P(\theta) \propto P(D_e | T_i(\theta), K_{e-1})$$

2025 May 01

51

$P(\theta)$ could be anything in general
 (Remember $P(\theta)$ is to-you-knowledge, not a law of nature)



But suppose it is like this.

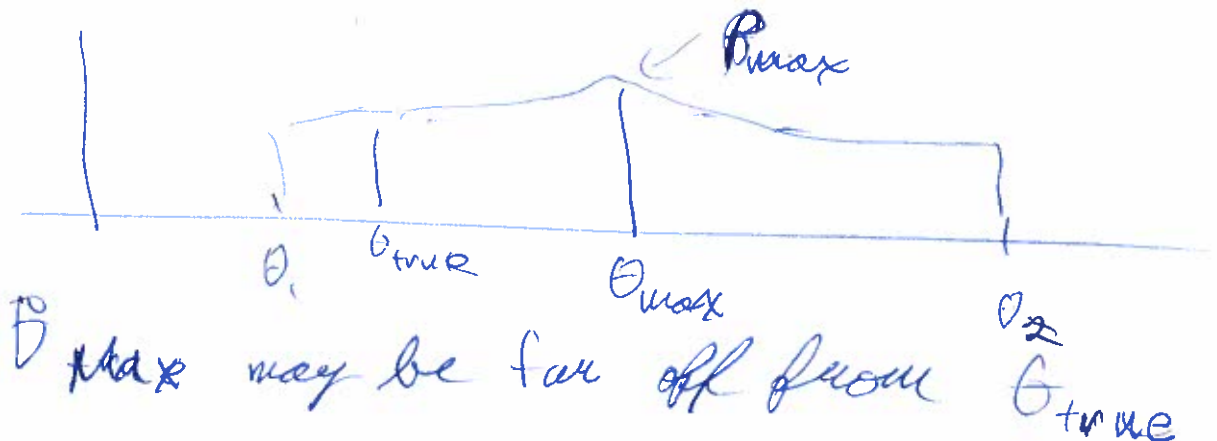
Most probability is clustered
 within $\Delta\theta_{width}$ of θ_{max} .

~~θ~~ $[\theta_{max} - \Delta\theta_{width}, \theta_{max} + \Delta\theta_{width}]$

Say
 $P(\theta_{max}, \Delta\theta) = 0.9$
 or 0.99
 or 0.999

θ_{max} is a good approx to θ_{true}
 if $\Delta\theta_{width}$ is rather small.

But $P(\theta)$ may not be so favorable.

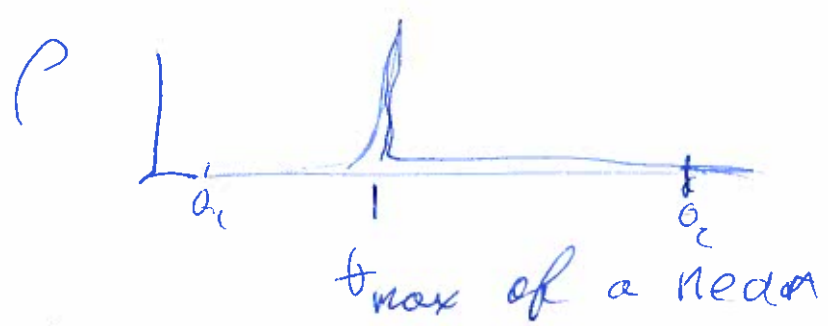


θ_{max} may be far off from θ_{true}

So if $P(\theta)$ is good
 θ_{max} should be close
 to θ_{true} ,
 but it can be very bad.

How do you make $P(\theta)$ good?
 find D_e decisive
 following Rutherford rule

$$P(D_e | T_j(\theta), \alpha_e)$$



- almost only the true θ
 can give ~~the~~
 data D_e .
 Dirac Delta function.

And if your theory isn't true?
 You may never find the decisive D_e
But you can always try ~~one~~
 theory T_j

(2019 nov 26)

↳ derive from information theory

Bayesian Information Crit. (BIC)

relative AIC

quality of model

↳ from information theory

Not absolute

$$AIC = 2k - 2 \ln L$$

↳ number of parameters the same

↑ k ↑
Keep number of estimated parameters

↑ L ↑, AIC ↓

$$\ln \left(\frac{L_{max}}{L_{max, AIC}} \right) \approx k \text{ or other value added}$$

max likelihood of model

Best of set of ~~all~~

$$\frac{[AIC_{min} - AIC]}{2}$$

$$= \frac{[2k - 2 \ln L - (2k_{min} - 2 \ln L_{min})]}{2}$$

like k or parameter added

here 1, others lower

BIC NOT

Similar to ~~Bayesian~~ some

somehow related

different applications

They BIC & AIC must work well in some ideal limits & so useful but can't be universally used or people

seller star
possibly often have superior
used in space
Pr 374
with some
since $P < 0.5$
34
no