

# Bayesian Analysis (BA)

- 1) Preamble: BA is path to truth quantified  
 a scientific method quantified.
- 2) Unquantified, or Qualitative Bayesian Analysis  
Priors & Posteriors
- 3) Bayes Theorem & a bit of History
- 4) Proof of Bayesian Analysis is in the ideal limit — which can be approached arbitrarily closely  
 → Truth (b) Bayes Odds Ratio & Bayes Factor
- 7) Multinomial Theorem & Multinomial Probability
- 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

2)

# 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

not like our computer packages

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

Mostly important theories are like this.

at least improvement, The Scientific Method, Qualified

E.g., Newtonian Physics,

↳ It is believed to be the macroscopic, <sup>limit</sup> of QM low velocity <sup>limit</sup> of Special relativity weak gravity limit of General relativity

One perspective is it an approximate theory

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.

13

Other examples 2<sup>nd</sup> law  
of thermodynamics

Natural selection evolution

↳ unlike Newtonian Physics  
these can be proven  
by mathematical logic

— and so can Bayesian Analysis  
↳ aside from pure rigorous  
quibbling.

## 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.

4) These are your prior probabilities  
or priors

Then new experience happens  
in that situation

and you vaguely update  
your priors

to your posteriors (which sounds  
awful and often  
is) { i.e.  
posterior  
probabilities

Your updating is also  
qualitative and of variable  
accuracy

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

Nature via Natural selection  
evolution has its own  
way 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

9

↳ especially you are new at the game — or the 'rules' have changed → each interview causes you to update your priors 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 <sup>more</sup> general conclusions about updates occur

↳ qualitative Bayesian analysis approaches the scientific method.

6)

Bayesian analysis itself  
is the sci Method Quantified  
(I argue)

### 3) Bayes Theorem

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))

Consider 3 events  $A, B, K$ .

$K$  is background knowledge that  
 $\Sigma$  introduces as extra item because  
 $\Sigma$  needs it later. It isn't needed for proof  
Bayesian Theorem itself.

Use  $A \cup B \cup K$  as joint event as a  
shorthand for union =  $\cup$  symbol

order  
has no  
meaning

in math  
which is too klutzy for me in the  
context.

# Conditional probability

$$P(A|K) = \frac{N_{AK}}{N_K}$$

Probability of A given K

Now

$$P(A|B|K) = \frac{N_{ABK}}{N_K} = \frac{N_{ABK}}{N_{BK}} \frac{N_{BK}}{N_K}$$

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

a factorization rule which is proven by frequentist approach. But you could assume it as just an axiom — but I think that is pointless.

Clearly

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

∴ Bayes Theorem in symmetric form easy to remember

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

Asymmetric form  $(P(A|B) P(B) = P(B|A) P(A))$

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

$$P(B|A) = \frac{P(A|B) P(B)}{P(A)}$$

↪ or suppressing K

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

$$P(B|A) = \frac{P(A|B) P(B)}{P(A)}$$

Q.E.D.

A frequentist definition of Probability (Trotta p. 9)

Trotta p. 9 argues that need this for probability formalism.

But I argue that you need it for probability formalism to have any meaning outside itself.

True  $N_{AK}$  and  $N_K$  may be vague often, but they must exist in some ideal sense. I dealization is the rule as Galileo argued for

understanding just leaving it implicit.

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

## Root of Bayesian analysis

↳ Bayesian analysis was greatly  
elaborated in 1939 BA,  
(Scientific Inference 3rd ed 1972  
by Harold Jeffreys (1891-1989)

(no relation - but he manages to  
look like my father anyways)

## What is Bayesian analysis?

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 easier.

In hard cases <sup>or real</sup> hard, because  
there is a lot of slop in assigning  
priors.

No one much used Bayesian analysis  
before 1990s. Since then it has  
had growing vogue.

Why was it unused before, [ 7 ]

Well, it 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 ~~are~~ so 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

10

speaking in an absolute sense  
or having some truth and some  
falseness definitely in  
less absolute sense.

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 absolute sense.

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

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

# 4) Proof of Bayesian Analysis ||

in Ideal limit

Say we have background  
or initial knowledge  $K_0$   
about some aspect of reality

and a set theories  $\{T_i\}$  (I of them)  
about that aspect. We assume theories  
discrete for simplicity  
often they form a continuum  
if they are distinguishable

[Of course, the set is another source  
part of  $K_0$ ]

Separated by free  
parameter values.  
In fact such continuous  
of theories are usually  
considered

And we have probabilities of being  
true (to our knowledge) our theory will

$$P(T_i | K_0) \quad \text{our initial priors}$$

How do you assign  $P(T_i | K_0)$

- maybe  $K_0$  gives you an idea
- the barest idea is the Principle of indifference

all  $P(T_i | K_0)$  are equal.

free parameter  
But we won't consider that complication now. Later.  
1995  
p. 776  
Not explicit  
Will

12)

Then you acquire data in  
~~a~~ a sequence  $D_1, D_2, D_3, \dots, D_e$

And your knowledge increases

$$K_1 = K_0 D_1$$

$$K_2 = K_1 D_2 = K_0 D_1 D_2$$

↓

$$K_e = K_{e-1} D_e$$

⋮

$K_{\infty}$  = You know everything  
but don't be  
disappointed if  
you don't get  
there.

Each data acquisition  
allows you to  
update priors to posteriors which  
then become priors for the  
next data acquisition

Since we are thinking ideally, we  
assume  $\{T_i\}$  is complete.

$$\text{ie, } P = \sum_i P(T_i | K_0) = 1$$

↳ if we are wrong we waff on  
may not build out.

Consider data acquisition [13]  
 I completed and recall Bayes Theorem  
 (see p. 7)

$$P(T_i | K_e) = P(T_i | K_{e-1} D_e) \quad \left\{ \begin{array}{l} \text{Prion} \end{array} \right.$$

$$\begin{array}{l} \uparrow \\ \text{Posterior} \end{array} = \frac{P(D_e | T_i K_{e-1}) P(T_i | K_{e-1})}{P(D_e | K_{e-1})}$$

We assume you can calculate the  $P(D_e | T_i K_{e-1})$  - the probability of  $D_e$  given theory  $T_i$  and  $K_{e-1}$

What is this?

$$\begin{aligned} &P(D_e | K_{e-1}) \\ &= \sum_i P(D_e | T_i K_{e-1}) P(T_i | K_{e-1}) \\ &= \sum_i \frac{N_{D_e T_i K_{e-1}}}{N_{T_i K_{e-1}}} \frac{N_{T_i K_{e-1}}}{N_{K_{e-1}}} \end{aligned}$$

Using the Frequentist definition of probability,  
 But what does that mean here,

You went did a trial (a complex trial) ab reality/ ~~the~~ <sup>N<sub>tr</sub>=2</sup> times and got 3 ~~the~~ <sup>theories</sup>

$N_{K_0} = 1$  since you assume your theories  $N_{K_0} = \frac{2}{2}$  by principle of indifference

14)

I'm a vague implicit  
and hard to calculate  
sense you get

$$\frac{N_{T_{i, k_{e-1}}}}{N_{k_{e-1}}}$$

} N's  
needed  
to check  
in

and the ratio

is easier by far than trying  
to analyse what the individual

N's are — and  
probably pointless too.

I'm a <sup>sort of</sup> multiverse sense

$N_{k_{e-1}}$  universes and

$N_{T_{i, k_{e-1}}}$  theory  $T_i$  is true  $N_{T_{i, k_{e-1}}}$

times

Do I believe this. Sort of but  
it really takes a more careful

# Marginalization

12 Oct 2011 Nov 2 '12

21  
(Nov 15-20)  
Yot

What if your theories are NOT discrete; i.e. theories  $\{T_i\}$

have free parameters collectively symbolized  $\theta$ ?

Bayes Odds ratio

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

Bayes k factor

Prior ratio

You don't believe you have all theories, even all interesting theories, and so the denominator of Bayesian Analysis is ~~usually~~ usually ~~measurable~~ and impossible to know — but that exists in principle is vital to Bayesian analysis being ~~the~~ true theory. So you just do Bayes Odds Ratio.

Almost always  $l=0$  and these are the initial prior since one only does one explicit Bayesian Analysis iteration. All past knowledge  $K_0$  is implicit earlier iterations.

But what do you do about free parameters and the priors for a theory with them?

Do you choose the parameters by maximum likelihood

$$L_{\text{max}} = P(D_e | T_i(\theta_{\text{max}}) | K_{e-1}), \quad \cancel{L_{\text{max}} = L(\theta_{\text{max}} | D_e)}$$

~~If a theory is true, maximum~~

Maximum likelihood gives best

choice assuming a theory is true,

$$\lim_{D_e \rightarrow \infty} \{ \theta_{\text{max}} \} = \theta_{\text{true}}$$

but not if the theory is not true.

If you assume equal priors for two theories, <sup>and</sup> one with far more free parameters, (theory i)

$$\text{then } K_{\text{factor}} = \frac{P(D_e | T_i(\theta_{\text{max}}) | K_{e-1})}{P(D_e | T_j(\theta_{\text{max}}) | K_{e-1})}$$

might be spuriously large

since the  $\theta_{\text{max}}$ 's are fitted to the  $D_e$ .

you could correct by choosing unequal priors to put theories on an equal footing but what is your guidance.

Usually no. In any case  $\theta_{max}$  may be useless choice of  $\theta$  if  $T_i$  is not true

The path is marginalization.

(Implements Occam's razor by effectively eliminating possibility of parameter of uncertain value and uncertain values)

Expand

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

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

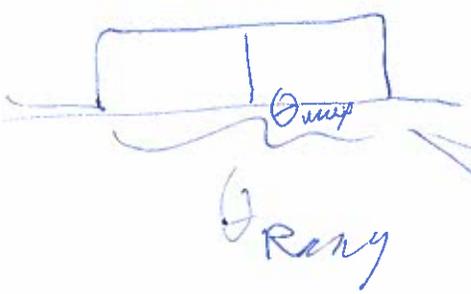
where  $P(\theta)$  is a probability density for the priors.

Usually, one just guesses

Can do more complicated things but need guidance

$$P(\theta) = \begin{cases} \frac{1}{\theta_{range}} & \text{for } \theta \text{ in } \theta_{range} \\ 0 & \text{for } \theta \text{ out of } \theta_{range} \end{cases}$$

The range that you think at all possible



Maximum likelihood parameters are NOT covered, since you don't know  $T_i$  is true

24

— as long as you set your range fairly

theories with unequal sets of parameters are on an equal footing and you can use equal priors again.

So marginalization implements Occam's razor effectively. It puts interesting theories on an equal footing.

eliminates pluralities

Of course for complex theories with many free parameters  $\theta$  and petabytes or exabytes of data to fit/produce

Make an mpf

$$P(D_o | T_i, k_o)$$

$$= \int P(D_o | T_i(\theta), k_o) p(\theta) d\theta$$

can be an immense supercomputer problem.

So it's not exactly choosing formal priors  $P(T_i | k_o)$  that is the uncertain part — Just set equal for interesting theories. — the uncertain part is choosing  $P(\theta)$  (a sort of prior)

Make an unfair choice  
and your Bayer  $k$  factor  
could be off by orders of magnitude.

→ This why the Jeffreys scale  
sets  $\log k = 2$ , ( $k=100$ )  
as decisive

and in cosmology  $k \approx 10$   
is considered completely undecisive

$\Sigma$  d. given.



# Bayesian Analysis

Analysis, Inference (more standard)

2017 Jul 16

1

## Ideal Procedure

### 1) Bayes' Theorem

Bayesian analysis is a large field of applied theories but the theorem is general to all probability analysis of Cosmology for Ch.22

Proof from frequency analysis (Frequentist analysis)

What can probability mean if it doesn't correspond to number of events of a population?

Is there any other rigorous (non-vague) way? Yes/No/Maybe

I'd say no at the moment.

A & B are <sup>general</sup> events  
(events = events, e.g., data, theories, beliefs, ...)

K is background event of set of theories

*THIS IS CONFUSING ORDER OF THINGS, SEE NEW P. 344*

$$P(A|K) = \frac{N_{AK}}{N_K}$$

$$P(B|K) = \frac{N_{BK}}{N_K}$$

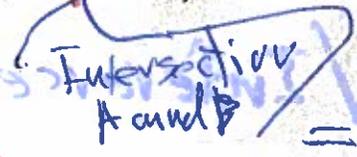
All background knowledge in Bayesian analysis

The theorem is a general simple rule of probability - beyond

$$P(A|B|K) = \frac{N_{ABK}}{N_{BK}} = \frac{N_{ABK}}{N_{BK}} \frac{N_{BK}}{N_K}$$

*Knowledge of all outcomes possible - all data sets*

*Background knowledge? How to phrase*



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

### Symmetric Form of Bayes' Theorem

$P(AB|K) = P(A|BK) P(B|K) = P(A|B) P(B|K) P(K)$   
Alternative picture  $\rightarrow$  causal

$$P(AB|K) = P(A|BK) P(B|K) = P(B|AK) P(A|K)$$

omitting K as background as usual  $\leftarrow$  Memorial form

$$P(AB) = P(A|B) P(B) = P(B|A) P(A)$$

### Asymmetric Form

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

omitting K as usual

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)} \quad \& \quad P(B|A) = \frac{P(A|B) P(B)}{P(A)}$$

2) It should be obvious how to go from  $K$ -initial form to  $V$  form.

But it wasn't initially to me, Now? sort of.

Note Bayes' Theorem is just a universal rule of probability.

Bayesian analysis as I understand it is the ideal mathematical scientific method which best/most obviously applies to intrinsically probabilistic theories

It's also in qualitative way how we learn and all bits too it seems - including unconsciously

Which most people don't do - But neural net machine learn does

## 2) Bayesian Analysis / Inference

has to be an ideal or how can you trust it

Iterative procedure ideally

Background knowledge  $K_0, K_1, \dots, K_n$   
 Fresh Data acquired  $D_0, D_1, \dots, D_n$

$$K_e = D_e K_{e-1}$$

The intersection yields new knowledge

a union not a product intersection

a true limit that you can approach

It's not exhaustive. Next time all  $P \rightarrow 0$ . it's wrong.  $P \rightarrow C$  except for one but don't know this one

Not exhaustive include trust exp

Set of theories  $\{T_i\}$  (a discrete set for can initial discussion)

Ideally ~~exhaustive~~ which is

ideal limit of all knowledge about a system Reach

3] easily reached when so defined as in student problems

$\exists$  ideally always.

— enough  $K_0$  and there can only be a finite set with non-zero initially probability

3

$$\left\{ \begin{array}{l} \forall T_i \in K_0 \\ P(T_i | K_0) > 0 \\ \text{Balls in } K_0 \end{array} \right.$$

Only one theory is true  
 $P(T_i) = 1$   
 $P(T_j) = 0$

But we don't have that probability about absolute truth  
 ↓  
 use probability we use on our knowledge  
 $P(T_i | K_0)$  is probability over knowledge

$K_n$  is not all knowledge  
 Just all knowledge relative to truth of  $E_i$

Once you have  $K_0$ , you know with certainty which theory is true.

in some cases it may be small & finite — so, but sometimes

$P(T_i | K_0)$  initial set of probabilities

Using Bayes Theorem

1 way to his Not wrong

$$l=1 \quad P(T_1 | D_1, K_0) = \frac{P(D_1 | T_1, K_0) P(T_1 | K_0)}{P(D_1 | K_0)}$$

likelihood of getting this data from this theory      likelihood  $\Rightarrow$  in many cases maximizes with respect to  $T_i$

2 general

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

Now what is  $P(D_e | K_{e-1})$  really?

4

It is not  $P(D_e)$

Always favor  $\Sigma$  believe

If you have  $D_e$ , ~~the~~  $P(D_e) = 1$ .

It's what the universe / reality actually gave you.

Remember ideally  $\Sigma T_i$  is exhaust

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

Frequentist understanding

What meaning?

Probability of  $D$  given your background knowledge

→ your understanding in terms of your theory

It's NOT 1.

See →

$$= \sum_j \frac{N_{D_e T_j K_{e-1}}}{N_{T_j K_{e-1}}} \frac{N_{T_j K_{e-1}}}{N_{K_{e-1}}} \quad (P_{K_{e-1}} = 1)$$

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

So

$$P(T_i | K_e) = P(T_i | D_e K_{e-1}) = \frac{P(D_e | T_i K_{e-1}) P(T_i | K_{e-1})}{\sum_j P(D_e | T_j K_{e-1}) P(T_j | K_{e-1})}$$

Posterior

Which sounds we're discussing posteriors

$\sum P(T_i | K_e) = 1$  unless denominator is your theory

50 asked on site if journal

Weighted average of the  $P(D_e | T_j K_{e-1})$  we  $\sum P(T_i | K_e) = 1$

$\rightarrow$  the  $P(T_i | K_{e-1})$  must be normalized  $\rightarrow$  Averaging  
 but that's not  $\rightarrow$

$$P(D_e | T_i, K_{e-1}) \left\{ \begin{array}{l} \geq \\ \leq \end{array} \right\} \left\langle P(D_e | T_j, K_{e-1}) \right\rangle$$

Choosing wisely on favored theories should help  
 ↓  
 Rather bad redup ↓  
 Minimizes statistics by doing the right experiment

then the Posterior { increase stays some }  
 { decrease }  
 (Normalization built into loop as all  $P(D_e | T_i, K_{e-1})$ )

you choose you new data sets wisely, convergence

to  $P(T_i | K_L) = \begin{cases} 1 & i = I \\ 0 & i \neq I \end{cases}$

quickly.  
 If not, the convergence can be slow and even terribly misleading.

Choose wrongest, bad luck, make blunder

The "true" theory can have its probability go down. (see p. 7)

If all data is just statistical with  $L = \infty$ , then  $L = \infty$ , { But if you are given some non statistical info,  $L$  finite. }

not zero probabilities

so  $T_I$  acts true with your  $K_e$

It may never do so if  $\epsilon$  is super small

in which case you give up when  $P(T_I | K_e) = 1 - \epsilon$  and just say  $T_I$  is true unless it fails at some later date  
↑ tiny amount

6

now save to p.17  
save die

problem for  
example, (see  
p.17)

~~$P(T_i | K_0)$~~

If all data statistical,  
then  $L = \infty$  formally.

But say after a million  
rolls, you study  
symmetry of the die  
and learn each  
face is symmetric.

$$P(D_i | T_i; K_{l-1}) = \begin{cases} 1 & \text{for equal probabilities} \\ & \text{face } T_i \\ 0 & \text{otherwise} \end{cases}$$

$D_i$  faces symmetrically

and the iteration stops short.

This is what Rutherford meant  
about doing wrong  
experiment

Trotter p.9

But if your true theory is intrinsically statistical some iteration is necessary.

The iteration stalls for you complete the cycle

64

Example of a case where you may never get to the truth <sup>in practice</sup>

Good ones

Completely deterministic as to source

but completely random to receiver

# Random number generators

Random or completely deterministic

It may not matter.

But a profound question. Yes profound

But in some respects it does NOT matter

in other respects, it does because we want to know truth

(maybe → are we ready for the truth)

For cycle repeat  $4.3 \times 10^{600}$  flo

Bill Press, Numerical Recipes (1992?)

n. 191 - 199 but seems to have omitted some of the great discussion

$\times 10^{18}$  flops  $\times \frac{3 \times 10^{16}}{64}$

$\sim \frac{4}{x} \times 10^{570}$  flo

maybe a quantum computer

we a flops

N flops

by cut of

Random

$4.3 \times 10^{600}$

$10^{600}$

$\times \frac{16}{3 \times 10^{16}}$

Good ones I recall

## Random Number Generators

completely deterministic relative to source & will recycle

Maybe a quantum computer can solve the cycle?

but random relation receiver Mersenne Twister Pseudo Random

→ 219937 - 1  $\triangleq 4.3 \times 10^{600}$

66

So you just ~~take~~ were given the string of numbers

After many statistical tests → Maybe a brilliant specially designed test

You conclude they are random

could pull all the determinism, but for most purposes

in truth "you're wrong"

Absolutely wrong effectively right

→ absolutely deterministic but for virtually any purpose → Monte Carlo etc.

wouldn't make

they work just as if they were

Philosophically

→ Realistically completely deterministic or with intrinsic random element as in QM

Random in time or randomness of time in "initial" condition

How can you tell?

Could make any difference to universe evolution or human history?

Maybe Not

# 3) Bayes Odds Ratio & Bayes Factor

It's a fact: getting the sum of probabilities of theories

Bayesian evidence  
Trotter-56  
But Ethical & never  
likelihood  
NOT  
Max.

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

(see p. 4)

is hard. Since rarely one seldom exhaustively knows all theories. ~~We don't want to test all survey~~

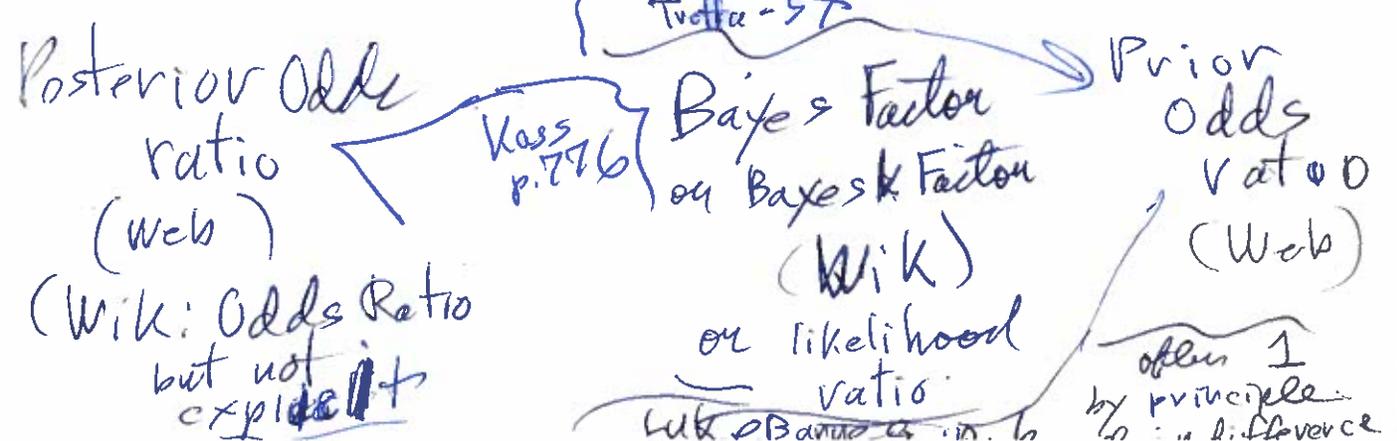
Practical after only 1 formal iteration  
Do

One often just wants to compare two competing theories → so Odds ratio

Relative Probabilities  
Wissenschaft

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

likelihoods → often maximize to fit parameters  
But in Bayes inference many values over parameters



8

Joffrey, ~~scale~~ <sup>log scale</sup> for k factors Kass p. 77

Based on some empirical analysis

- < 0
- 0 - 1/2
- 1/2 - 1
- 1 - 3/2
- 3/2 - 2
- 2 → up

negative for theory i

barely worth mention

significant

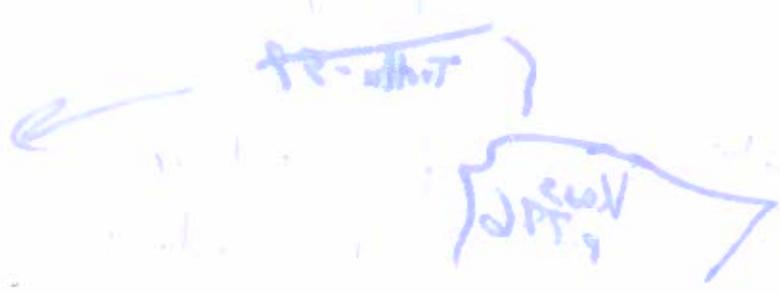
strong

very strong

decisive

estimate with available information

Why? Subject to later revision you do not expect and admit perhaps seeking in a direct sense.



# 4) Die Problem (AKA Dice Problem) 9

In general.

omit in lecture for non-urgent problem

Polynomial theorem



J faces labeled

$$i = 1, \dots, J$$

$P(T_1)$

Theory 1

$$P_i = \frac{1}{J}$$

$P(T_2)$

Theory 2

$$P_i = \frac{i}{\sum_{i=1}^J i} = \frac{i}{J(J+1)/2}$$

Priors

modulo function (wik)

Gauss Proof

$P(T_3)$

Theory 3

$$P_i = \frac{2 + \text{mod}(i, 2)}{J}$$

so evens twice as likely as odds

assume J even

- 1  $J = J+1$
- 2  $J-1 = J+1$
- ...
- J  $1 = J+1$

$$\frac{1}{2}(J+1) + \frac{1}{2}(J-1) = \frac{3}{2}J$$

Say N throws =  $N = \sum_i n_i$

with outcomes  $n_1, n_2, \dots, n_J$

$$P(n, \{ \}) = \frac{N!}{n_1! n_2! \dots n_J!} \prod P_i^{n_i}$$

Multinomial theorem

used to get

Throws 1 1 2 4 1

only one ordering occurs

$$C \frac{n_1! \dots n_J!}{n_1! \dots n_J!} = N!$$

among identical throws

Those that look identical only occur once identical permutations actually occur

Combinations Wik: Combinatorics



# 4) Multinomial Theorem & Multinomial Probability Distribution

- A digression to let me set up an example application of Bayesian analysis

Say you have a set of  $I$  events:  $1, 2, \dots, I$   
 with probabilities  $P_1, P_2, \dots, P_I$

and you then took a (sample) of  $N$  ~~events~~ events  
 and got e.g.,  $P_2, P_3, P_4, P_7, \dots, P_{I-1}, P_I$   
 event count  $\rightarrow 2, I-3, 4, 7, \dots, I-1, N$   
 $\left. \begin{matrix} P_2 \\ P_3 \\ P_4 \\ P_7 \\ \dots \\ P_{I-1} \\ P_I \end{matrix} \right\} \begin{matrix} \text{Product} \\ \text{is} \\ \text{the} \\ \text{prob} \\ \text{of this} \end{matrix}$

~~How many distinct combinations are there?~~  
~~distinct permutations~~  
~~statistics jargon~~

combination: collection without order distinction  
 permutation: collection with order distinction  
 distinct permutations: those you can tell apart by appearance

For  $N$  objects of  $I$  kinds  
~~For our trials~~

$$N! = \prod_{i=1}^I n_i! = C(\{n_i\}) \prod_{i=1}^I n_i!$$

$\underbrace{N!}_{\text{Permutation of each set of events}} = \underbrace{\prod_{i=1}^I n_i!}_{\text{combination of } N \text{ objects}} = \underbrace{C(\{n_i\})}_{\text{Number of distinct permutations corresponding to combination } \{n_i\}} \underbrace{\prod_{i=1}^I n_i!}_{\text{Number of permutations among identical objects.}}$

Multinomial Coefficients with

106

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

no. of distinct permutations of collection  $\{n_i\}$

Now consider our trials.

All distinct permutations

~~can turn up~~ ~~are distinct overall events.~~

But only ~~individual ones~~ = set of indistinct ones are just one overall event. and are individual events

... 1 ... 2 ...  
... 2 ... 1 ...

Two overall events.

You draw from the 1 bin, then the 2 bin and then the reverse

Two ways of doing this

<u>Count</u>	1	2	3	4
event	<del>1</del>	<del>1</del>	2	6
prob	$P_1 \times P_1$	$\times P_2$	$\times P_2$	$\times P_2$

... 1 ... 1 ...  
... 1 ... 1 ...

Just one overall event.

Here you draw from the 1 bin then 1 bin

Just one way of doing it!

the probability of this event

Imaginary interchange of 1's doesn't double the probability. (selection with replacement in stats argon)

Probability of the ~~combination~~ combination  $\sum u_i$  10c

$$P(\sum u_i \Xi) = C(\sum u_i \Xi) \prod_i P_i^{n_i}$$

every distinct permutation of  $\sum u_i \Xi$  has this product probability

∴ Multinomial Distribution

$$P(\sum u_i \Xi) = \frac{N!}{\prod_i u_i!} \prod_i P_i^{u_i}$$

Prove  $= N! \prod_i \left( \frac{P_i^{u_i}}{u_i!} \right)$

Now  
total probability  
conserved

$$\sum_i P(\sum u_i \Xi) = 1 \quad \text{It should be}$$

Proof The probability of some event happening is 1

Well  $1 = \sum_{i \neq} P_i$  on first event of N years

$\left. \begin{array}{l} P_i \sum_j P_j \\ \text{then} \\ \text{sum on } i \end{array} \right\} 1 = \left( \sum P_i \right)^2$  on second event of N sequence

$1 = \left( \sum P_i \right)^N$  on Nth event of N sequence

$P_i \left( \sum_j P_j \right)^2$  then sum for three events and so on

(10d)

1 = expansion is isomorphic to our country of the  $\mathbb{C}(\{n_i\})$

Multinomial theorem

$$= \sum_i \mathcal{P}(\{n_i\})$$

with each one-to-one correspondence

with  $p_1, p_2, p_3, \dots, p_r$

replacing  $1, 2, 3, \dots, r$

I'm not totally convinced of the cogency, but the result is true, QED

Binomial Theorem

special case: Binomial theorem

& Binomial probability distribution

Prove by induction if affected by Paradox.

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

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)

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

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

the  $n$  choose  $k$  formula

$$l=0$$

$$M_0 = 1$$

$$\boxed{10e}$$

$$l=1$$

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

$$l=2$$

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

$$= np \left[ \cancel{(x+q)^{n-1}} + n(n-1)(x+q)^{n-2} \right]_{x=p}$$

$$= np [1 + \cancel{p(n-1)}]$$

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

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

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

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

$$= np(1-p)$$

BeV-53  
correct

10A

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

11

Theory 1

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

$\sum n_i = N$

This is for I units  
 I units  
 issue  
 $P_i = 1$   
 i odd  
 $P_i = 2$   
 Total  
 $= \frac{I}{2} \cdot 1 + \frac{I}{2} \cdot 2 = \frac{3I}{2}$

Theory 3

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

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

As so often in Astronomy, god has told us there are only 2

Theory 2

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

possible true theories

Very easy to post dict (PCDIT) not like complex systems with messy data

5) Die Problem for Zeus II

Normal problem

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

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

$\left. \begin{matrix} \frac{1}{9} \text{ i odd} \\ \frac{2}{9} \text{ i even} \end{matrix} \right\}$

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

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

So 10 throws complex enough, but ~~not~~ decide? (N=10)

Really did this

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

math

$P(D/I_1) = \frac{10!}{3! 2! 2! 1! 1! 1!} \cdot \left(\frac{1}{2}\right)^{10} = 2.50057 \dots \times 10^{-3}$   
 $P(D/I_2) = \left(\frac{1}{2}\right)^{10} 1^6 \cdot 2^2 = 6.238 \dots \times 10^{-4}$

$$k = 3.604$$

Probability given your knowledge - Before even, No assumption

So Prior odds = 1

Bayes factor  $\left\{ \begin{array}{l} k = 3.604 \\ \text{Posterior odds} = 3.604 \end{array} \right.$

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

$$P(1) = \frac{3.604}{4.604} = 0.7828 \dots \quad P(2) = \frac{1}{4.604} = 0.2172 \dots$$

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!!!

$$P_1(1) = 0.7828 \dots \quad P_2(2) = 0.2172 \dots \quad \text{No posteriors} \rightarrow \text{priors}$$

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	4	5	3	5

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

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

So in math 2  $\leftarrow$  But follow odds  $= k \times \frac{P(1, k_1)}{P(1, k_1)} = k^2 \left| 13 \right.$   
 we get  $k = 3.604$  again  $\approx 13$

because some number of even & odd  
 But favors theory 2 twice in a row

$P_2(4k) = .9285$   $P_2(2k) = .07148$

So we keep advancing to truth  
 But no! Instead of rolling die,  
 we investigate its symmetry

Rutherford Rule

What if you can't or don't know how or just need to do a simple experiment.

$D_3 \Rightarrow$  6-fold symmetry

$P(D_3 | 1k_2) = 1, P(D_3 | 2k_2) = 0$

$P(1 | k_3) = \frac{P(D_3 | 1k_2) P(1 | k_2)}{P(1 | k_2) \cdot P(D_3 | 1k_2) + P(2 | k_2) P(D_3 | 2k_2)}$   
 $\leftarrow P(1, k_2) + P(2, k_2)$  are now irrelevant

See part but factors in other order.

$= 1$

$P(2 | k_3) = \frac{P(D_3 | 2k_2) \cdot P(2 | k_2)}{P(1 | k_2) \cdot P(D_3 | 1k_2) + P(2 | k_2) P(D_3 | 2k_2)}$

$= 0$

So Rutherford was right!! Truth is  
 "If you need statistics, you did the wrong experiment!"

19)

But if the system is intrinsically probabilistic.

(Die is)

and you cannot find its distribution directly (as for <sup>you can</sup> Die)

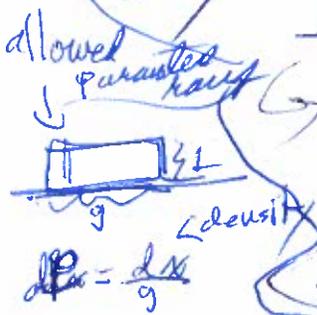
then Bayesian analysis makes sense.

- Cosmology
- social science
- epidemiology

Marginalization process

5

### Marginalization



Say your theory has a free parameter

it's a continuum infinity of theories.

$\langle P_1 | T_1 \rangle \langle P_2 | T_2 \rangle \dots \langle P_n | T_n \rangle$

But best to regard as one theory with free parameters

$T(x)$   
↑  
integrate over all allowed range - some physical data used some lev

If you think your theory

maybe common choice? Truth is choice for set of free run

set of

$P(D_2 | T, K_{e-1})$  v.i. /  $P(D | \theta)$  15

is right put it to

We'll ~~do it~~ expand on this in later

the data using maximum likelihood to determine

Very ~~to maximize~~

(global one or local if

that are possible

you can be fooled by either or both

Local maximum & global maxima } several theories all with free parameters

free parameters

All data relevant to theory

as ~~θ~~ → all  
 $\theta \rightarrow \theta_{exact}$

if your theory is exactly right, all data should be fit

the best fit one may just because

it has lots of free parameters

Chi-Square will decrease as freedom increases =  $\chi^2$  =  $\chi^2$  =  $\chi^2$

But also

if your theory has enough free parameters

and you are fooled

wonder of path

What if two different theories fit absolutely all data → some equivalent

Even if theory is wrong? in maximum theory

Need Occam's razor

compare Non-best fits of theories

Bayesian always over free parameters

And put ~~complex~~ many & few parameters theories on a level playing field.

over good range the prior variance

and rate theories independent of parameter fit  
 independent of parameter fit just in actual practice rather than in principle

Trotter-57 says a level more parameters analyzed - maybe

been by some people

16) - rate on best fit  
 has prior range of parameter.

Some a nuisance parameters

can't be measured maybe

not of interest (WIK)  
 at least at moment

and deciding on the theory

So now

$$P(T_i | K_{e-1})$$

$$= \int P(T_i(\theta_i) | K_{e-1}) P(\theta_i) d\theta_i$$



Often Topkat prior

But Iveta 29, 30

a probability distribution itself

But Replace v.4 formula.

$\theta_i$  parameter of theory

weight function

$$P(\theta) =$$

1 where possible

$$\sum_j P(D_e | T_j(\theta_j) K_{e-1}) P(T_j(\theta_j) | K_{e-1}) P(\theta_j) d\theta_j$$

otherwise where you know the parameter can't be

$$\sum_j \int P(D_e | T_j(\theta_j) K_{e-1}) P(T_j(\theta_j) | K_{e-1}) P(\theta_j) * d\theta_j$$

sum over all theories

Marginalize over all parameter sets

$$P(T_i | K_e) P(\theta_i) =$$

$$P(T_i | K_e)$$

Flat prior not best!

If you have vast sets of data

biological petabytes

$1 \text{ pb} = 10^{15} \text{ bytes}$

LHC in 5 pb per year

couple processes  
3.4 pb per day  
in 2009

Integration can be huge even for our modern standards  
Just ~~doing~~ ~~it~~ ~~with~~ ~~factories~~

Are shortcuts of varying utility of which I do not know

1) Physiological information with little (WIKI) done  
eg Schwann cell from Bayes - an approach

Finding  $P(\theta) | T_i(\theta, k)$

for the problem simple but petabytes produced follow a complex theory

↓

I intuit that here you need Markov chain Monte Carlo MCMC

Cotta-35 and a lot of skill

derive from information theory

# Akaike Information Crit. (Wike)

relative AIC quality of model (from information theory)

smaller better

Not absolute

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

for number of parameters

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

↑ k ↑  
↑ number of estimated parameters

↑ L ↑ AIC ↓  
↑ likelihood of model

probably better than comparison

Best of set

$$[AIC_{min} - AIC] / 2$$

like k or parameter odds

has 1, others lower

## BIC NOT

similar to BIC NOT somehow related

different results

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

probably better than comparison  
value of Chi Square  
which become indices of p < 0.5 or no