# BAYES, BANJOS, BAD TEST SCORES, AND BIG SCIENCE

Brian W Junker Department of Statistics Carnegie Mellon University Professional Statistics: <u>brian@stat.cmu.edu</u> Everything Else: <u>junkerb@gmail.com</u>



Sunday Forum, 1UU Pg

2

### **Topics**

- Bayes
  - · Probability on a Dartboard, Conditional Probability
  - The Reverend Thomas Bayes (not a little-known Unitarian!)
  - · Bayes' Theorem
- Banjos
  - Getting the beat in music
- Bad Test Scores
  - Scoring the SAT and similar tests
  - What to do with a zero or a 100%??
- Big Science
  - Big data in science, scalable computation, and Bayes
  - Bayesian Brains?

## Probability on a Dartboard

- Suppose I throw a dart at a circular target on a square board.
- Let's let Event A = "dart lands inside the circle"
  - Each dart is an independent trial on which "A" can happen
  - If darts land equally likely anywhere in the square, then the probability of "A" will be proportional to the ratio of the area of the circle to the area of the square



### Aside: Monte Carlo estimation of $\pi$

• If  $P(A) = \pi/4$  then  $\pi = 4 \times p$ 

```
    If we have estimated
```





then  $\pi pprox$  4 imes 0.78 = 3.12

- This method is called "Monte Carlo estimation". It was first named and applied to difficult problems by the <u>Manhattan Project</u> folks
  - The estimate gets better and better the more darts we throw, by the "Law of Large Numbers"
  - (but that is a different presentation!)



## **Conditional Probability**

 The idea of conditional probability does not depend on darts. If A and B are two events, then we define

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

#### This is entirely consistent with our estimate of P(A|B):

$$P(A|B) \approx \frac{\#(\text{green darts})}{\#(\text{green & red darts})}$$

$$= \frac{\#(\text{green darts})/\#(\text{all darts})}{\#(\text{green & red darts})/\#(\text{all darts})}$$

$$\approx \frac{P(A \& B)}{P(B)}$$

$$P(A|B)$$

## The Reverend Thomas Bayes

- 1701?-1760
- 2<sup>nd</sup> Gen Presbyterian Minister & Amateur Mathématician
- Published
  - Divine Benevolence, or an Attempt to Prove That the Principal End of the Divine Providence and Government is the Happiness of His Creatures (1731)
  - An Introduction to the Doctrine of Fluxions, and a Defence of the Mathematicians Against the Objections of the Author of the Analyst (anon.,1736)
- Wrote but did not publish
  - An Essay towards solving a Problem in the Doctrine of Chances



Sunday Forum, 1UU Pgh

- This almost certainly is not Bayes(!)
- It is the only claimed likeness
- Not a Unitarian (or Universalist)!

### Before Bayes' "Essay"

- Official statistics have been kept since 2200BC (China)
- By the mid-1700's in Europe
  - Statistical Tables were kept for taxation and other state purposes
  - Observations were combined to reduce measurement error (especially in Astronomy)
    - Calculations were developed to show how much better the mean was, than an individual observation, as a measurement
    - Least squares was used as a curve-fitting tool
    - Gauss would later (ca. 1795) develop the normal distribution to place least squares on more principled footing
- Bernoulli & deMoivre knew how to take initial conditions B and derive probabilities of events A from them <u>(probability)</u>
- BUT they could not take a set of events A and infer the initial conditions B from them <u>(inverse probability)</u>



 Although the algebra is easy (next slide!) it was a big breakthrough to be able to write this down.

9

### Bayes' Theorem on Inverse Probability

- Suppose we know P(A), P(B), and P(A|B).
- We want the "inverse probability" P(B|A).

$$P(B|A) = \frac{P(B\&A)}{P(A)}$$
 Definition of  $P(B|A)$   
$$= \frac{P(A|B) \cdot P(B)}{P(A)}$$
 Since  $P(A|B) = P(A\&B)/P(B)!$ 

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

(Did you blink?)



### Bayes' Billiard Table (he didn't call it that!)

• Let's play a game!

- I will randomly place a black ball on a pool table
  - This is the <u>initial condition</u> (it is <u>random!</u>)
- You get to blindly place a ball on the table 1000 times
  - I will only tell you if it is to the left or to the right of my ball
  - These are the <u>observable events!</u>
- Can you guess where my ball is, from
  - p = # of your balls to the left of my black ball
  - q = # of your balls to the right of my black ball
- R: billiards()

13 January 2013

Sunday Forum, 1UU Pgl

14

### **Bayes' Billiard Table**

- Easy to "guess" the black ball is at 9/(9+17) = 0.35
- With his theorem, Bayes could also derive
  - A probability curve showing where the black ball was likely or not likely to be
  - A <u>recipe</u> for an interval where the black ball was likely to be!
- Bayes could measure <u>uncertainty</u> in the guess!





- <u>Bayes</u> gave a better partial solution, but only in terms of the billiards example
  - He shared the Essay with essentially no one while he was alive
  - Posthumously edited, updated and published by <u>Richard Price</u>
- <u>Pierre-Simon Laplace</u> independently discovered these ideas and published them in a complete, general form
- And so, of course, we call it Bayes' Theorem.

### **Bayesian Decision-Making**

- <u>States of nature</u>: B<sub>1</sub>, B<sub>2</sub>, ..., B<sub>N</sub>
- Observable event: A
- <u>Decision</u>: Which B<sub>k</sub> given, A?
- <u>Prior distribution</u>: P(B<sub>1</sub>), P(B<sub>2</sub>), ..., P(B<sub>N</sub>)
- Likelihood function:
  - $P(A|B_1), P(A|B_2), ..., P(A|B_N)$
- <u>Posterior distribution</u>: Inverse probability of each B<sub>k</sub>.:
  - $P(B_1|A) = P(A|B_1) \cdot P(B_1) / P(A)$ •  $P(B_2|A) = P(A|B_2) \cdot P(B_2) / P(A)$ 
    - ...
       P(B<sub>N</sub>|A) = P(A|B<sub>N</sub>) · P(B<sub>N</sub>) / P(A)
- 13 January 2013

Sunday Forum, 1UU Pgh

18

Decision: Choose the Bk

with maximum  $P(B_{k}|A)$ 

### **Bayesian Decision-Making**

 Each posterior probability has same factor 1/ P(A); can omit this and get the same decision:

•  $P(B_k|A) = P(A|B_k) \cdot P(B_k) / P(A) \propto P(A|B_k) \cdot P(B_k)$ 

"posterior is proportional to likelihood times prior"

- Interpretation: Bayes' Theorem lets you combine
  - Past experience about prevalence of each B<sub>k</sub> (prior)
  - Consistency of the data A with each B<sub>k</sub> (likelihood)

in making a decision (in a way that minimizes expected loss [or maximizes expected gain]...)



### A likelihood for beats: P(music | time sig)



- But some are more likely than others!
- Temperly develops P(music | time sig) based on this.





### **Music and Probability**

- Temperly applies the same Bayesian principles to
  - Recognizing the rhythmic grid of a music
  - Recognizing the key from a monophonic melody
  - Recognizing the key from a polyphonic melody
  - Etc

#### • He then discusses the consequences of this for

- Musical expectation and error detection
- Recognizing different musical styles
- Scope for deviation from a musical style
- Etc.
- (fun book! \$20 at Amazon.com...)



#### Sunday Forum, 1UU Pgl

2

### Bad Test Scores

- Standardized tests are a big industry in the US; see Zwick (2002) for an overview.
- Many tests are not scored on simply the "number right".

#### Instead

- a likelihood is formed,
- and a parameter (kind of an "ideal test score") is estimated from the likelihood, from each student's test score

### Building the likelihood

• The "test score" likelihood is a product of functions:





This is "maximum likelihood estimation" (MLE)

likelihood

25

### Problem: What about 0% or 100% right?



# **Bayesian Regularization** Incorporating prior "beliefs" makes the problem of finding the ideal score manageable. It doesn't make infinite scores impossible, just very unlikely. Much practical Bayesian statistics is like this. The prior distribution "regularizes" the problem enough to make it solvable. This "regularization" is how I first became acquainted with **Bayesian methods!** 30 Advantages of Bayes in Science Regularization Borrowing strength from past experience (or concurrent data-rich settings) to improve inferences Language of probability Common approach to formulating problems and solutions Tools to find optimal solutions Vast expressive power to formulate many scientific problems in probabilistic terms

### **Bayesian Successes**

- In small to medium-sized problems, there has been a Bayesian revolution
  - Enabled by cheap, fast computers with capacious memory
- Some examples:
  - Fair methods for scoring tests and extracting as much information as possible from them
  - Analysis of complex federal surveys
  - Analysis of face to face (or online!) social networks
  - Pooling information from multiple raters (e.g. in radiology or job performance), and learning about individual differences among raters
  - Describing and understanding trajectories toward Alzheimers' Disease
- Other examples:
  - · Inferring treatment success from historical medical records
  - Nate Silver's 538 (and PEC, and ...) predictions!
  - Spam detection, email/word completion, ...

13 January 2013

Sunday Forum, 1UU Pgl

32

### **Bayesian Limitations**

- "Computational scalability"
  - We know how to solve small problems but when the same problem gets large, we can't do the computation
- Thomas Bayes experienced this:
  - For small p, q, any freshman in calculus could solve the problem of estimating where the black ball was
  - For large p, q, this solution is tedious, and too slow for a human to carry out
- With modern computers, we can extend the reach of Bayes theorem to modestly large problems
  - N = # observations
  - P = # features per observation
- N ~ 1,000,000, P ~ 1,000 is sort of the Bayes ballpark

### **Big Data and Bayes**

- Many current "Big Science" & "Big Commerce" projects involve truly big data (large N, large P!), e.g.
  - Sloan digital sky survey (N=500,000,000 objects)
  - Online recommender systems (N=300,000 of purchases/day)
  - Google search logs (N=400,000,000 searches/day)
  - Large Hadron Collider at Fermilab (N=10<sup>3</sup>--10<sup>9</sup> or more collisions per second)
  - Microarray (gene expression) data (N=10<sup>3</sup> observations, P=10<sup>6</sup> features)
- Three problems
  - Computation may be too slow (Thomas Bayes' problem for supercomputers!)
  - Likelihood may be too complex to be represented in the computer
  - Curse of dimensionality: as P grows, the number of probes needed to search a volume of radius R is proportional to R<sup>P</sup>



Sunday Forum, 1UU Pgr

34

## Strategies for Big Data

- Reduce P: Dimension Reduction & Feature Extraction
- Reduce N: Subsampling
- Focus computation on a narrow inferential goal
  - Don't try to represent the whole likelihood or the whole posterior
- Scalable computation take at most one pass through the data!
- Hastie, Tibshirani & Freedman (2009) survey some basic methods
  - Some methods are "provably correct"; others are just heuristic

### Aside: Bayesian Brains?

- "All people are Bayesians, some just haven't figured that out about themselves yet"
  - I kind of doubt it
- The same problems with representing models, and doing exact computations on them, exist for computers or human brains
  - Human brains are not especially fast or efficient at computation
- Most likely, humans adopt heuristics that reduce to (approximately) Bayesian thinking, for small problems.
  - In this sense, we probably handle inference the way big data inferences area handled

### References

- Stigler, S. (1980). The history of statistics: the measurement of uncertainty before 1900. Cambridge MA: Harvard University Press.
- Temperly, D. (2007). *Music and probability.* Cambridge MA: MIT Press.
- Zwick, R. (2002). *Fair game? The use of standardized tests in America.* New York: Routledge-Farmer.
- Hastie, T., Tibshirani, R. & Freedman, J. (2009). The elements of statistical learning (2<sup>nd</sup> ed). NY: Springer. See also <u>http://www-stat.stanford.edu/~tibs/ElemStatLearn/</u>

36