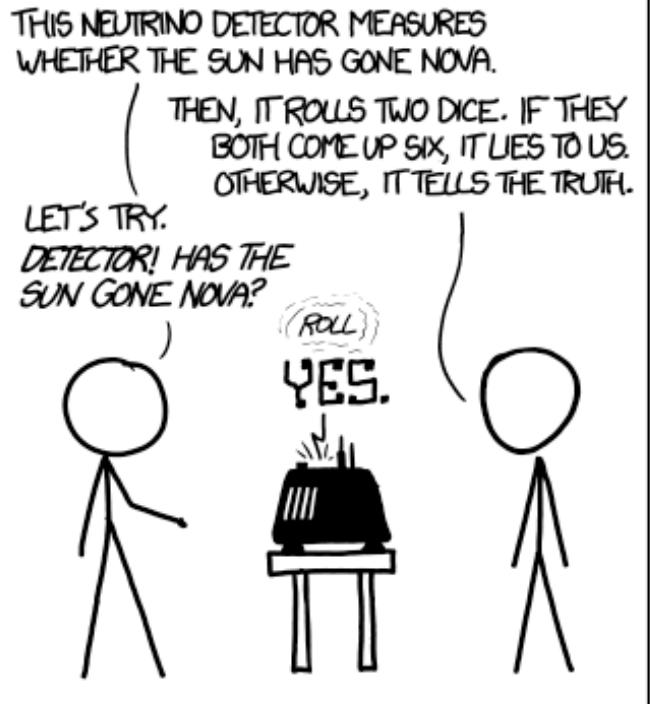


## Pre-class assignment

- 1) In Bayes' Theorem, what is the utility of including the prior and evidence probabilities? (In other words, why is it useful to include them?)
- 2) What is the difference between the frequentist and Bayesian approaches toward statistical inference, and how does it manifest in the way these approaches are used?

# DID THE SUN JUST EXPLODE?

(IT'S NIGHT, SO WE'RE NOT SURE.)



FREQUENTIST STATISTICIAN:

THE PROBABILITY OF THIS RESULT HAPPENING BY CHANCE IS  $\frac{1}{36} = 0.027$ . SINCE  $p < 0.05$ , I CONCLUDE THAT THE SUN HAS EXPLODED.

A stick figure stands next to a small barbecue grill-like device.

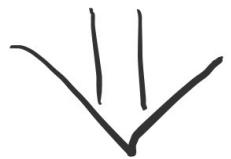
BAYESIAN STATISTICIAN:

BET YOU \$50 IT HASN'T.

A stick figure stands next to a small barbecue grill-like device.

Bayes' theorem :

$$P(x|y) = \frac{P(y|x) \times P(x)}{P(y)}$$



$$\begin{aligned} [\text{probability of } x \text{ given } y] &= \frac{[\text{prob. of } y \text{ given } x] \times [\text{prob. of } x]}{[\text{prob. of } y]} \\ \text{or} \end{aligned}$$

$$[\text{posterior}] = \frac{[\text{Likelihood}] \times [\text{prior}]}{[\text{evidence}]}$$

In terms of a model (w/params  $\Theta$ ) and data ( $d$ ):

$$p(\Theta|d) = \frac{p(d|\Theta) p(\Theta)}{p(d)}$$

or

[probability that model  
parameters  $\Theta$  describe  
this data  $d$ ]

=

Likelihood  
[probability of  
data given model  
parameters]

prior  
[prob of  
model  
parameters]

[probability of data]

posterior

Evidence

Question You receive a test for a certain disease. 2% of people in your demographic have the disease, and the test is 90% reliable. If the test is positive, how likely are you to actually have the disease?

Bayesian approach:

$$P(\text{disease given positive test}) = \frac{P(\text{test given disease}) \times P(\text{disease})}{P(\text{test being positive})}$$

$P(T|D)$                            $P(D)$   
↓                                    ↓  
 $P(D|T)$                            $P(T)$

$$P(D|T) = \frac{P(T|D) P(D)}{P(T)}$$

$$P(T|D) = 0.9 = \text{accuracy of test}$$

$P(D) = 0.02$  = chance of a person in your demo having this disease

$$P(+)= (\text{prob. of true positive}) + (\text{prob. of false positive})$$

$$= P(+|D) P(D) + P(+|\text{not } D) P(\text{not } D)$$

$$= 0.9 \times 0.02 + (1-0.9) \times (1-0.02)$$

$$\boxed{P(D|+)} \approx 0.1552 \sim 15.5\%$$

Followup: how does the probability of you actually having the disease change as a function of test reliability? (Try reliability values of 0.8, 0.95, 0.99, 0.999)

<u>reliability</u>	P(DIT)
0.75	0.0577
0.8	0.0784
0.9	0.1552
0.95	0.2794
0.99	0.6689
0.999	0.9532
0.9999	0.9951
0.99999	0.999

Context: COVID-19  
antigen tests

One  $\approx$  80%  
 reliable, and  
 false positive and  
 negative results  
 may differ.

COVID-19 molecular  
 tests are closer  
 to 99-99.9%  
 reliable.

Question: If I have photometric observations of a galaxy that might be a high-redshift galaxy but could be other things (e.g., a closer, red galaxy), how do I assess the Bayesian probability of the high- $z$  hypothesis?

Some hints:

- 1) What do you know about galaxy evolution over time?
- 2) How do you include both true positives and false positives?

probability  
of high-z hypothesis

Likelihood that  
high-z galaxy model  
matches data

prior based on  
galaxy distribution as  
 $f(\text{redshift})$

$$P(\Theta | d) = \frac{P(d | \Theta) p(\Theta)}{P(d)}$$

= P(prob. of data  
match to model)

$$P(d) = P(\text{data fit to high-z model}) p(\text{high z galaxies}) + P(\text{data fit to low-z model}) \times p(\text{low-z galaxies})$$

true positives

false positives

plugging in some (hand-wavy and unreliable) numbers:

$$p(\text{dl high } z) = 0.7 \quad [\text{high-redshift model matched to data}]$$

$$p(\text{dl low } z) = 0.4 \quad [\text{low-redshift model matched to data}]$$

$$\begin{aligned} p(\text{high } z \text{ gal}) &= 0.1 \\ p(\text{low- } z \text{ gal}) &= 0.9 \end{aligned} \quad \left\{ \begin{array}{l} \text{probability representing} \\ \text{underlying galaxy} \\ \text{distribution (more galaxies at low redshift!)} \end{array} \right.$$

} note: these do  
not need to  
add up to 1.  
Why?

$$P(\text{high } z_{\text{galaxy}} \mid \text{data}) = \frac{0.7 \times 0.1}{0.7 \times 0.1 + 0.4 \times 0.9}$$

$$= 0.16279 \rightarrow \approx 16.3\% \text{ galaxy is at high redshift!}$$

notes:

- underlying sample matters: more galaxies at low redshift are observable! (so even a low prob. of error for a given galaxy impacts the outcome significantly)
- model uniqueness/dark fit matters a lot: this is why few-band photo-z is awful, many-band helps w/excluding false positives

photo-z  
PDF for a  $z \sim 6$  blue galaxy w/  $\leq 5$  HST bands might look a lot like this:

