

Bayesian Statistics

The first principle is that you must not fool yourself and you are the easiest person to fool.

- Richard Feynman



Bayesians and Frequentist don't agree on Planes and Elections

- Frequentist give probability of event
- Bayesians give confidence of an event happening $[0,1]$
- Are frequentist methods incorrect then?

No.

Frequentist methods are still useful or state-of-the-art in many areas. Tools such as least squares linear regression, LASSO regression, and expectation-maximization algorithms are all powerful and fast. Bayesian methods complement these techniques by solving problems that these approaches cannot, or by illuminating the underlying system with more flexible modeling.

DID THE SUN JUST EXplode?
(IT'S NIGHT, SO WE'RE NOT SURE.)

THIS NEUTRINO DETECTOR MEASURES
WHETHER THE SUN HAS GONE NOVA.

THEN, IT ROLLS TWO DICE. IF THEY
BOTH COME UP SIX, IT LIES TO US.
OTHERWISE, IT TELLS THE TRUTH.

LET'S TRY.

DETECTOR! HAS THE
SUN GONE NOVA?

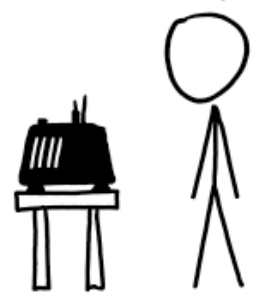
ROLL
YES.

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.

BAYESIAN STATISTICIAN:

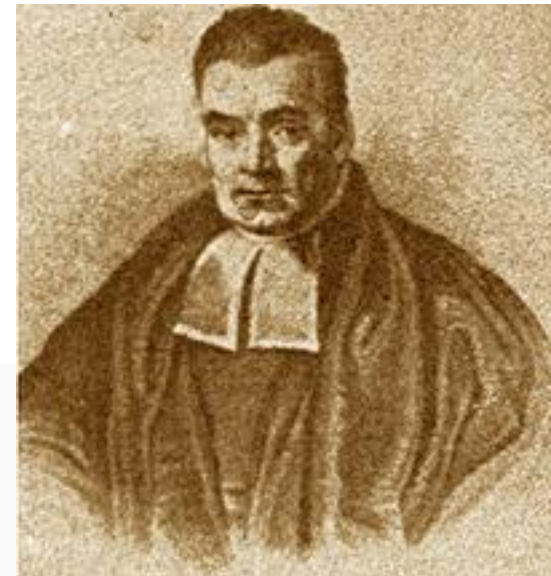
BET YOU \$50
IT HASN'T.



History

- Bayes never published, Richard Price found it among his notes after Bayes' death in 1761, re-edited it, and published it. Unfortunately, virtually no one seems to have read the paper, and Bayes' method lay cold until the arrival of Laplace
- Pierre-Simon Laplace, came to believe that probability theory held the key, and he independently rediscovered Bayes' mechanism and published it in 1774.
- In 1781, Richard Price visited Paris, and word of Bayes' earlier discovery eventually reached Laplace. Laplace eventually refined it to the form we see today.
- Bayes method died again until the invention of machines, most famously being used by Alan Turing to crack the Enigma codes.
- Computational advances in the 90s launched the 'Bayesian revolution' in a long list of fields. This is only partly because Bayes' Theorem shows us the mathematically correct response to new evidence. It is also because Bayes' Theorem works.

Bayes Theorem



Likelihood

Probability of collecting this data when our hypothesis is true

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

Prior

The probability of the hypothesis being true before collecting data

Posterior

The probability of our hypothesis being true given the data collected

Marginal

What is the probability of collecting this data under all possible hypotheses?

Bayes Theorem

Likelihood

Probability of collecting this data when our hypothesis is true

Prior

The probability of the hypothesis being true before collecting data

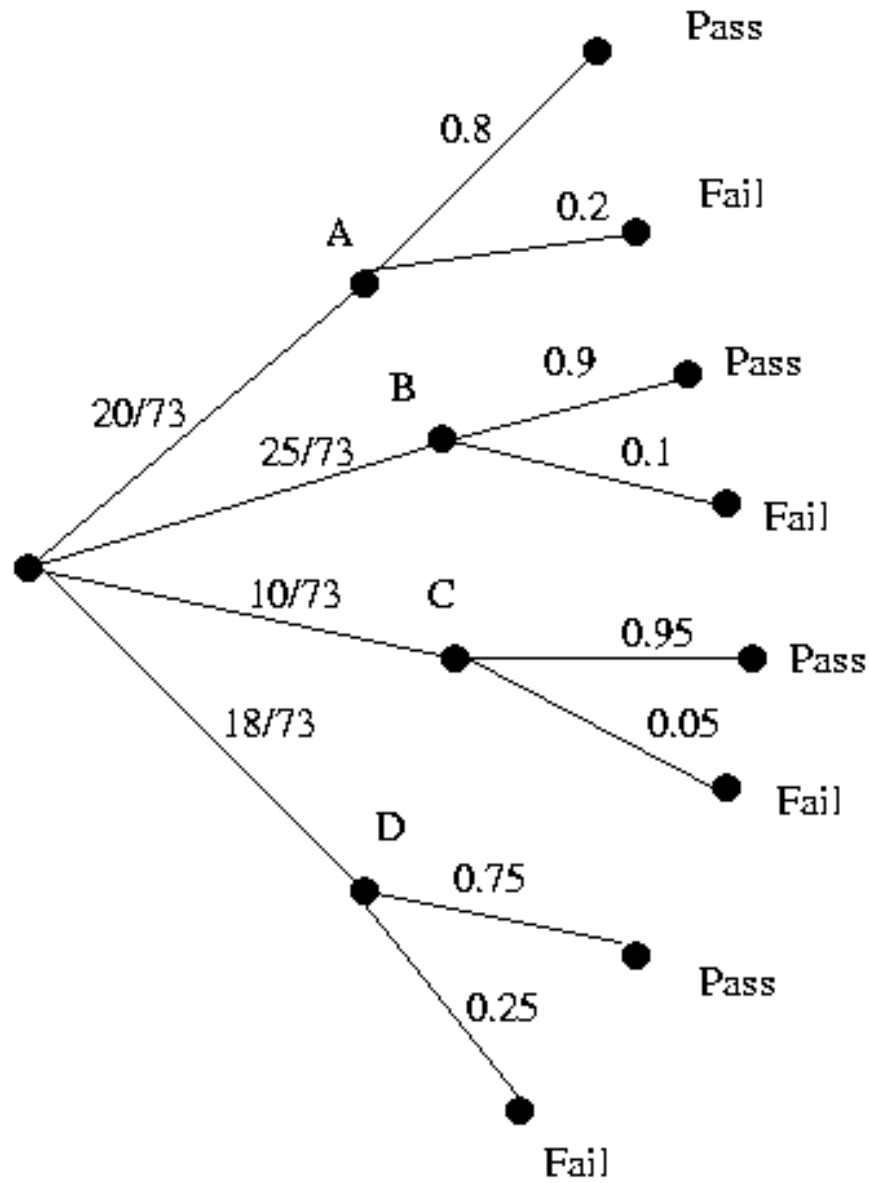
$$P(H|D) = \frac{P(D|H) P(H)}{P(D)}$$

Posterior

The probability of our hypothesis being true given the data collected

Marginal

What is the probability of collecting this data under all possible hypotheses?



1. Probability you will pass?
2. What is the probability that you will pass the class if you are in section D?

(read carefully here)

3. What is the probability that if you pass, you are in section D?

Notice that the likelihood stays the same as the Frequentist and the Marginal is essentially a normalization factor.

Therefore the prior is the name of the punch line.

Bayes Theorem

Likelihood

Probability of collecting this data when our hypothesis is true

Prior

The probability of the hypothesis being true before collecting data

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

Posterior

The probability of our hypothesis being true given the data collected

Marginal

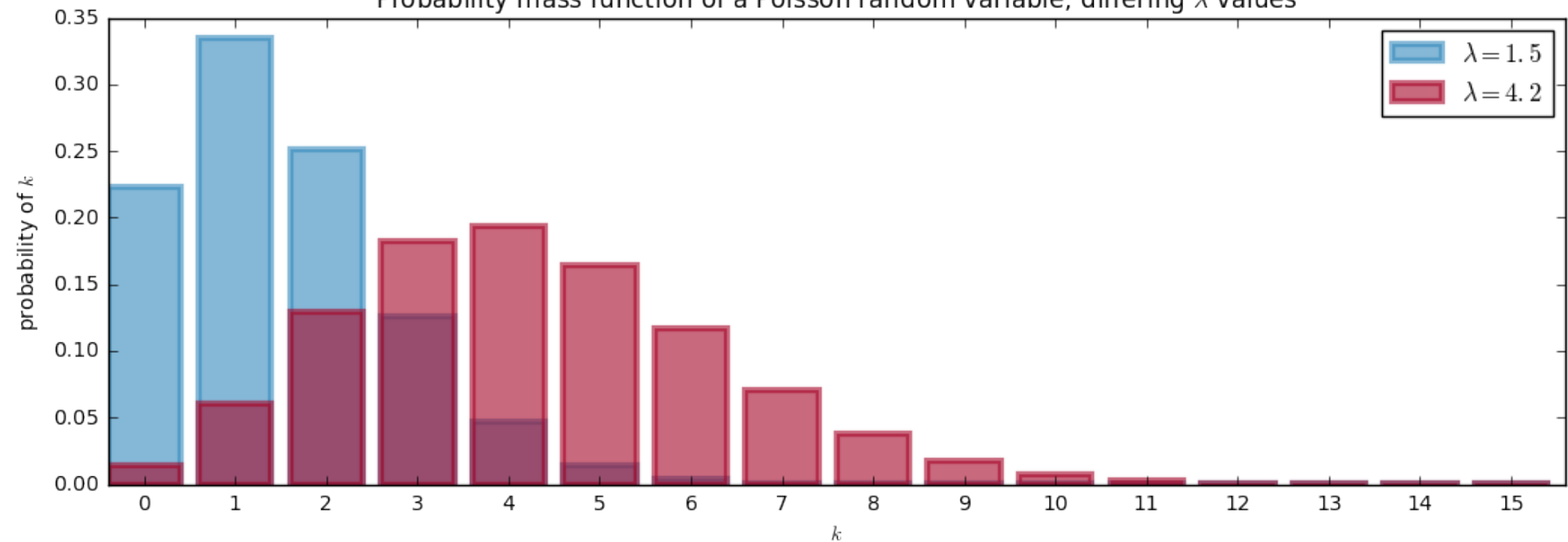
What is the probability of collecting this data under all possible hypotheses?

Probability Distributions

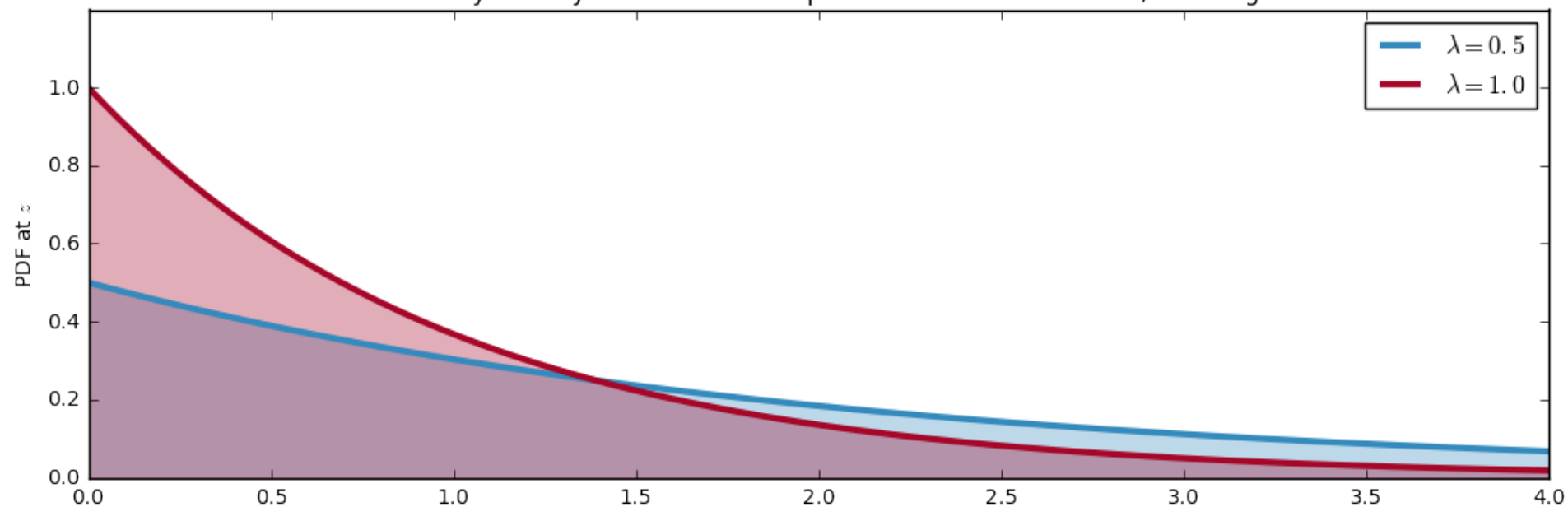
We can divide random variables into three classifications:

- Discrete: Discrete random variables may only assume values on a specified list. Things like populations, movie ratings, and number of votes are all discrete random variables. Discrete random variables become more clear when we contrast them with...
- Continuous: Continuous random variable can take on arbitrarily exact values. For example, temperature, speed, time, color are all modeled as continuous variables because you can progressively make the values more and more precise.
- Mixed: Mixed random variables assign probabilities to both discrete and continuous random variables, i.e. it is a combination of the above two categories.

Probability mass function of a Poisson random variable; differing λ values



Probability density function of an Exponential random variable; differing λ



Bayes Theorem

Likelihood

Probability of collecting this data when our hypothesis is true

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

Prior

The probability of the hypothesis being true before collecting data

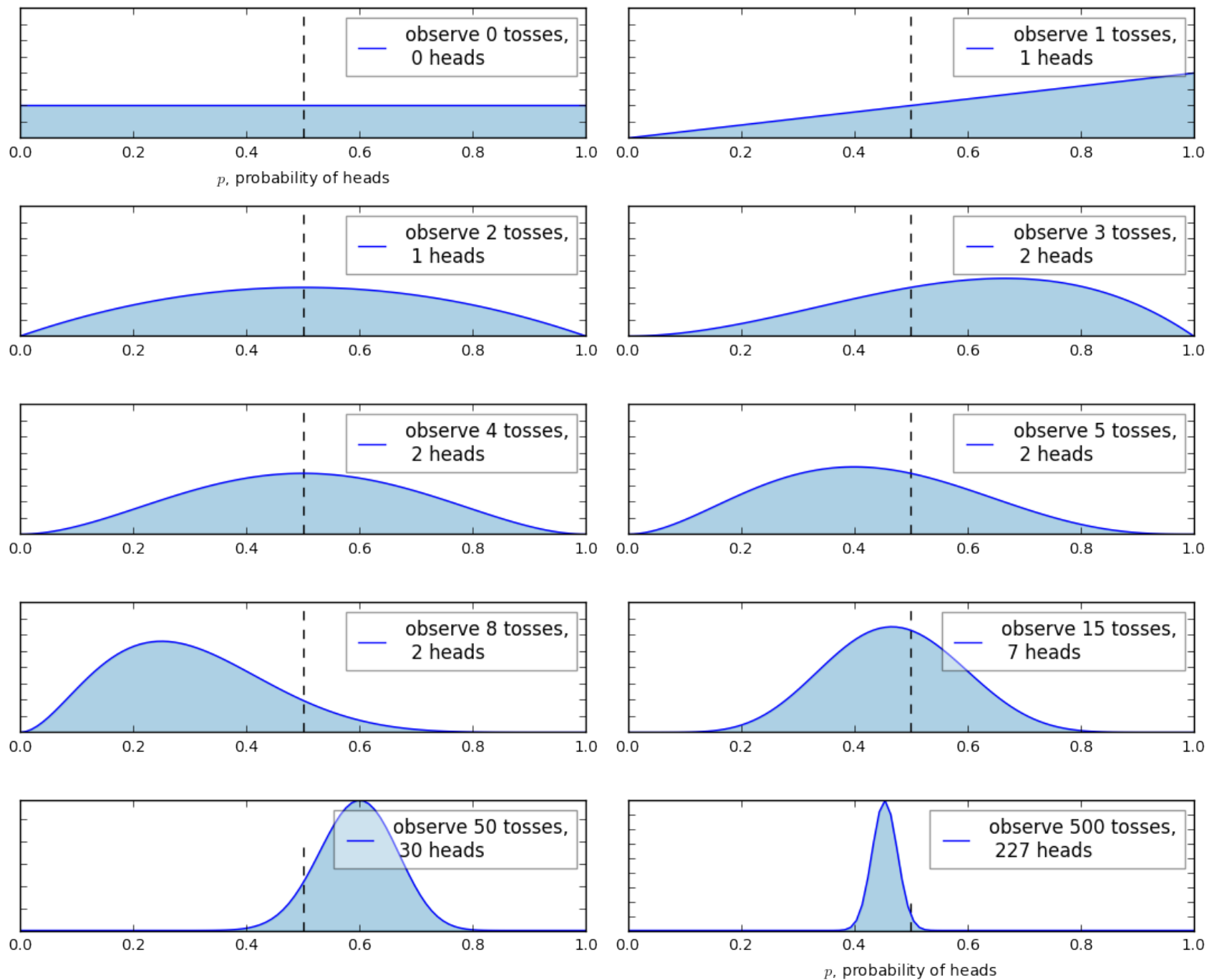
Posterior

The probability of our hypothesis being true given the data collected

Marginal

What is the probability of collecting this data under all possible hypotheses?

Bayesian updating of posterior probabilities



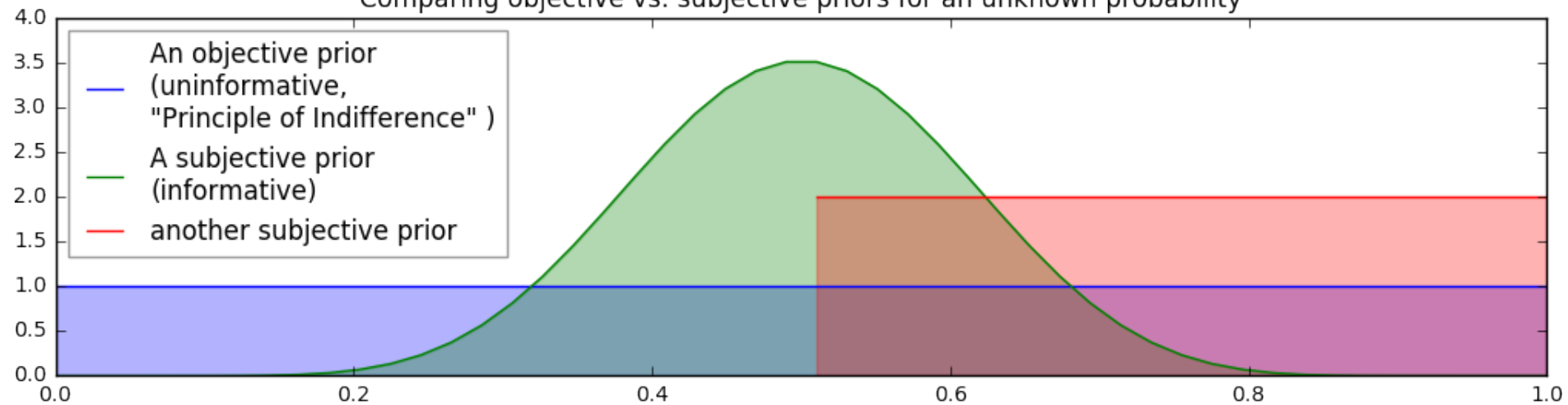
Subjective vs Objective priors

- objective priors, which aim to allow the data to influence the posterior the most, and subjective priors, which allow the practitioner to express his or her views into the prior.

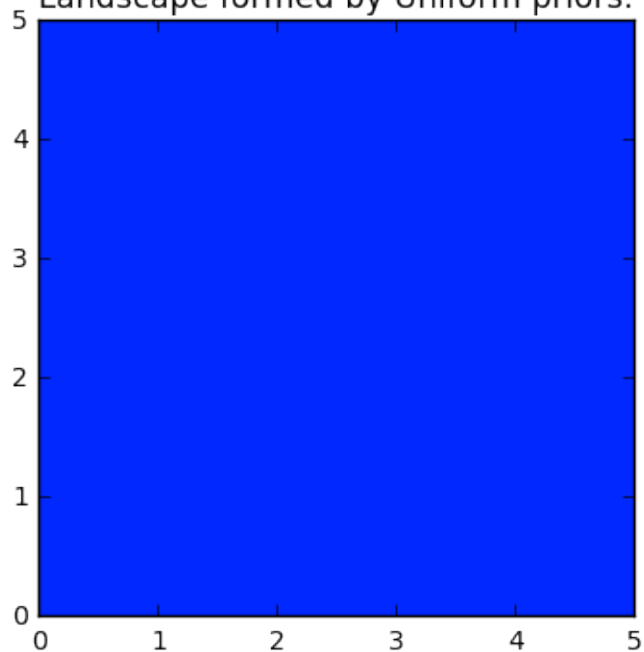
Example is the coin flip, is a flat prior "The Principle of Indifference", which is a uniform distribution over the entire possible range of the unknown. Using a flat prior implies that we give each possible value an equal weighting.

- Subjective Priors On the other hand, if we added more probability mass to certain areas of the prior, and less elsewhere, we are biasing our inference towards the unknowns existing in the former area. This is known as a subjective, or informative prior

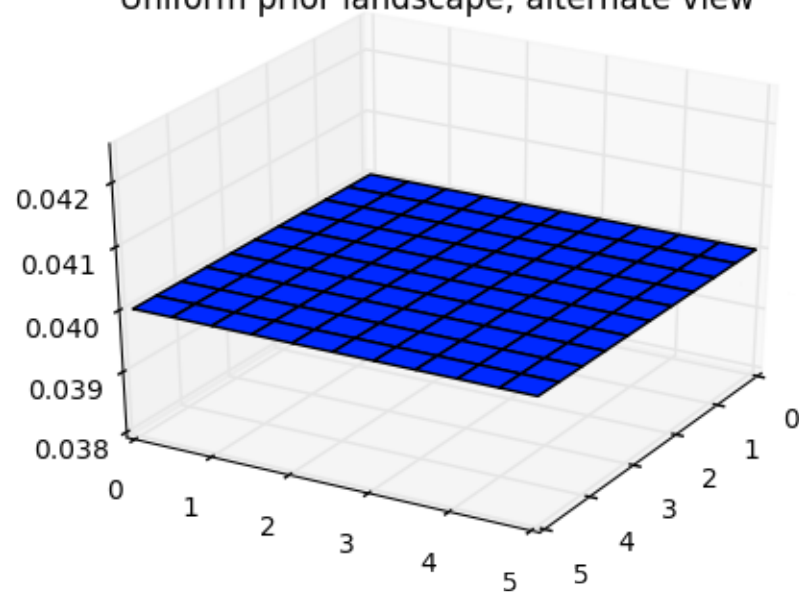
Comparing objective vs. subjective priors for an unknown probability



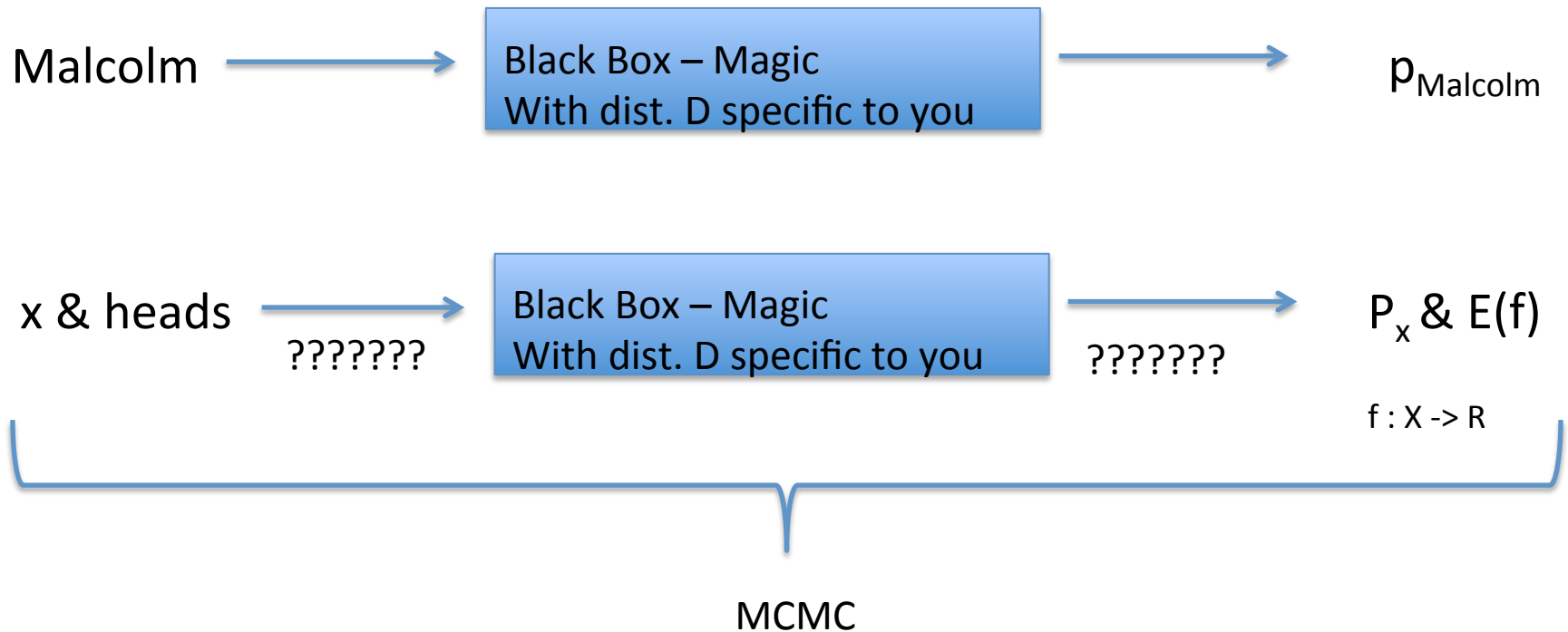
Landscape formed by Uniform priors.



Uniform prior landscape; alternate view



Problem Example

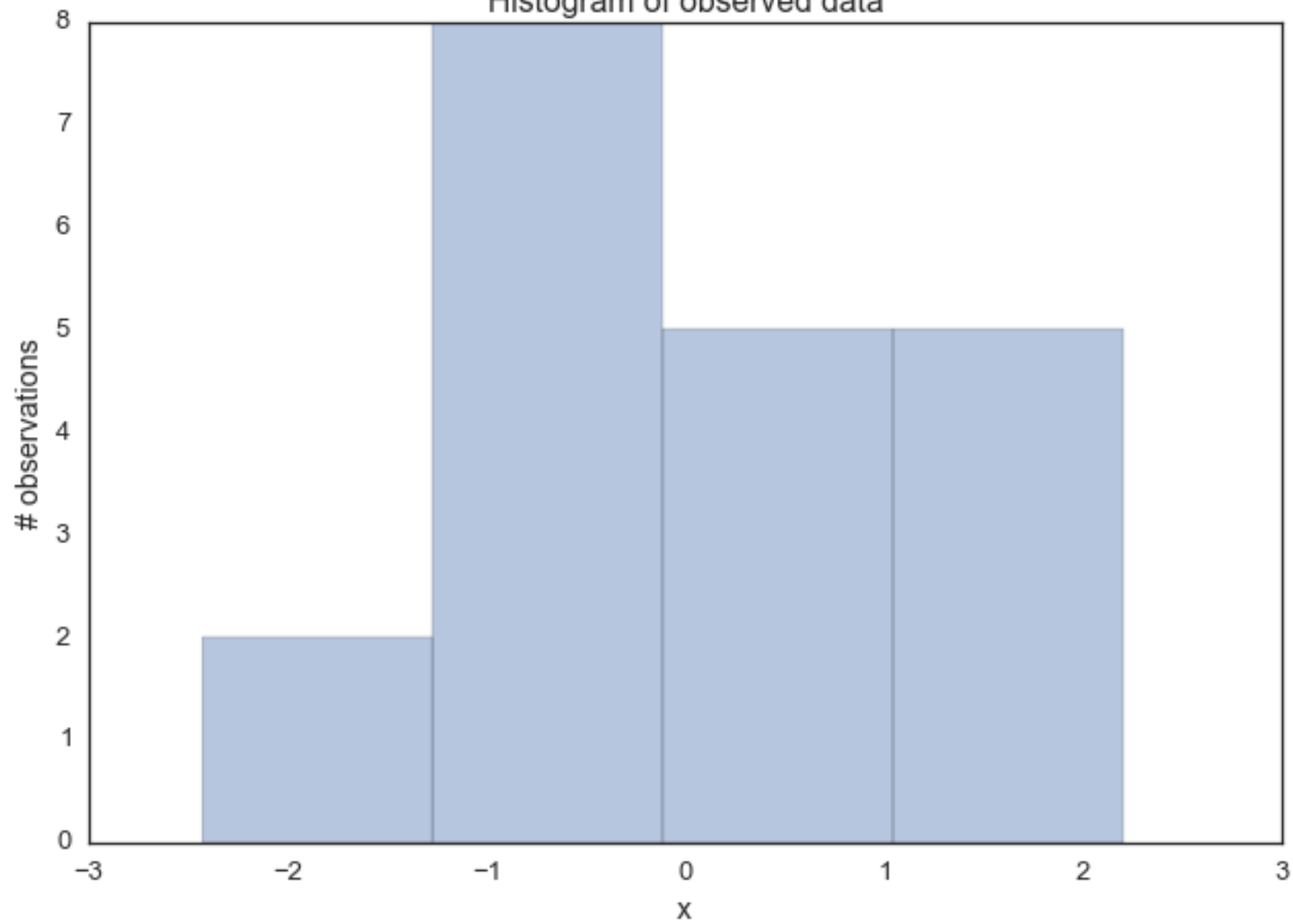


Putting These Methods into Practice

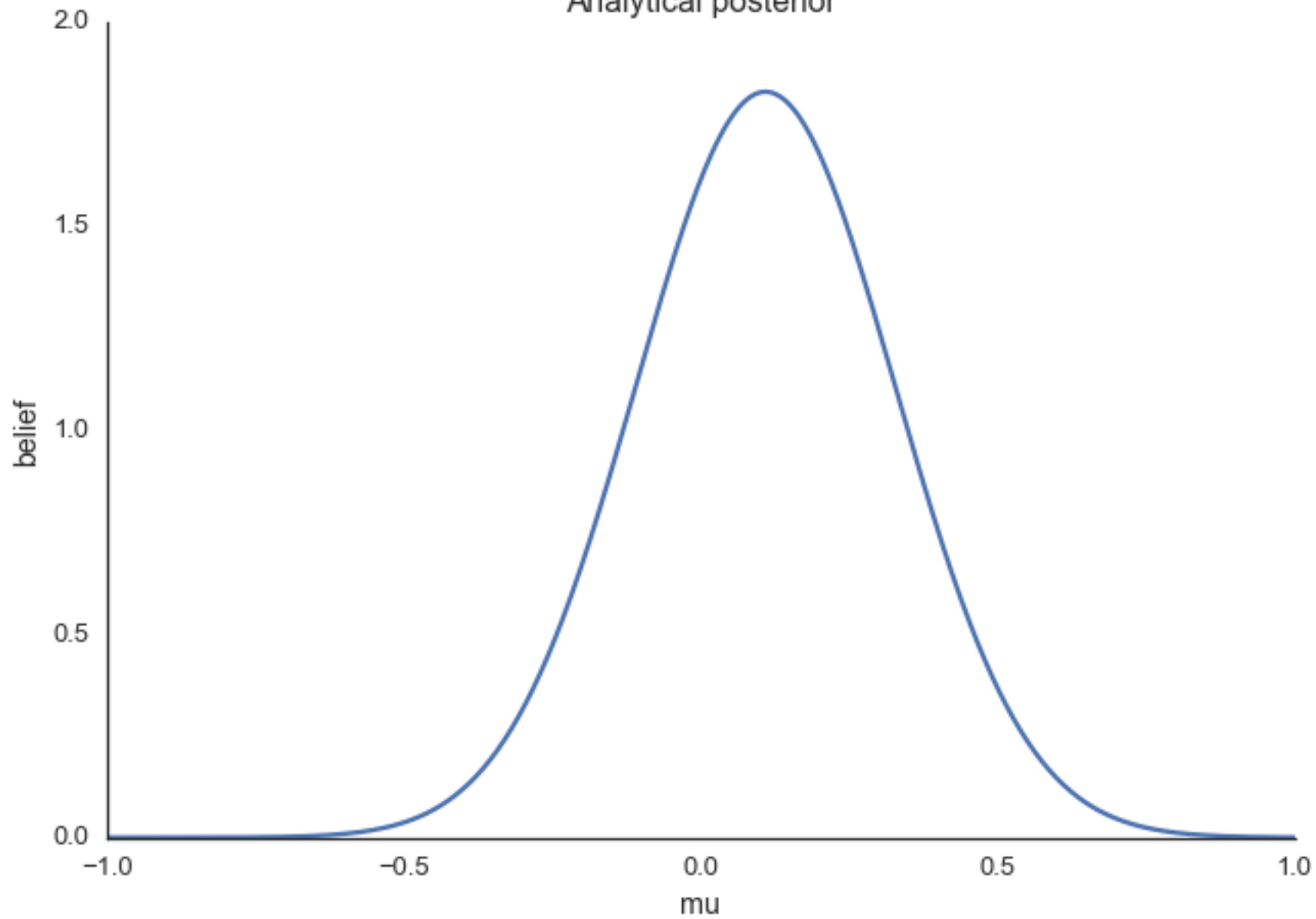
Markov Chain Monte Carlo (MCMC)

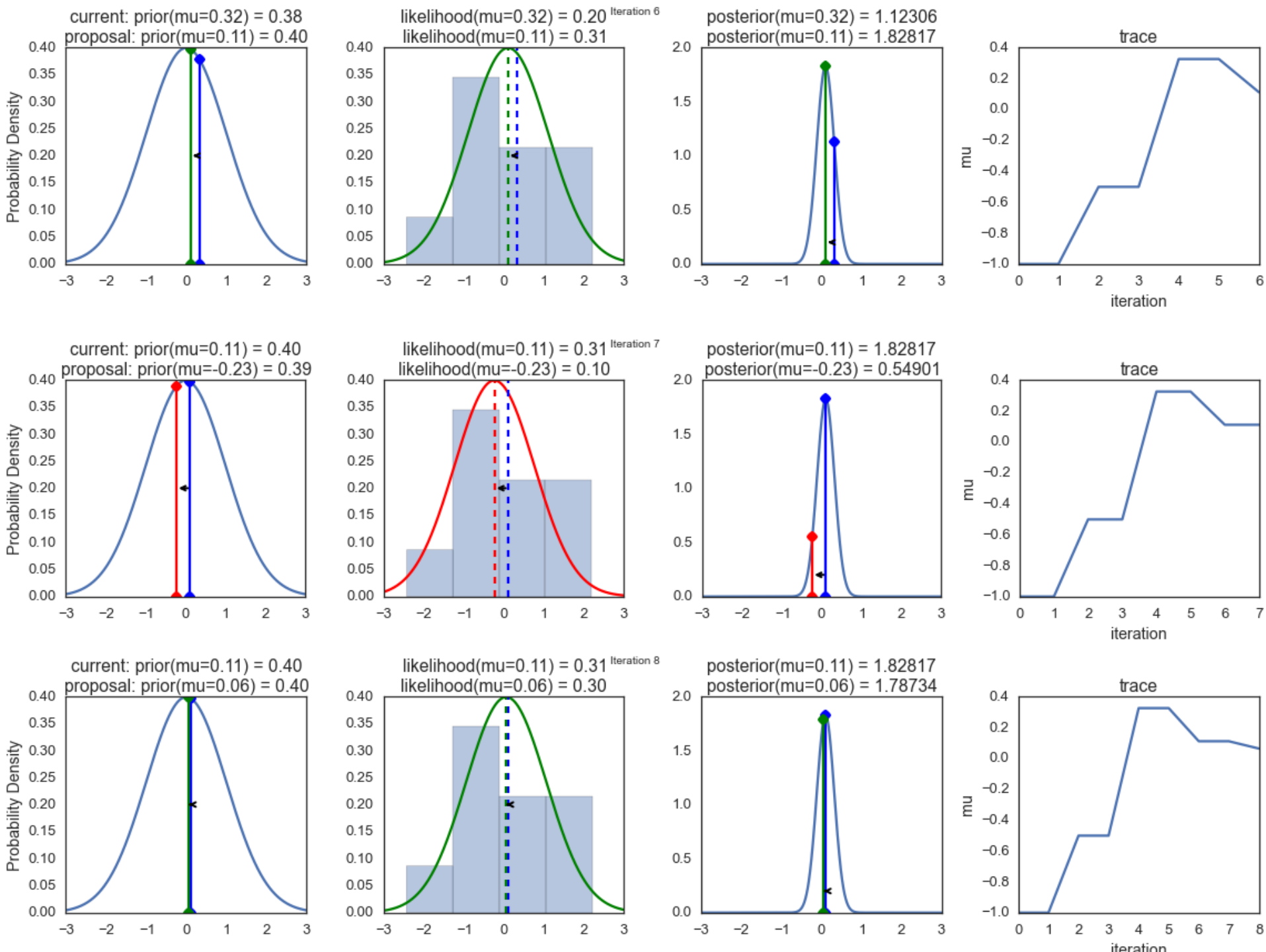
- Bayesian inference problem with N unknowns, we are implicitly creating an N dimensional space for the prior distributions to exist in. Associated with the space is an additional dimension, which we can describe as the surface, or curve, that sits on top of the space, that reflects the prior probability of a particular point.

Histogram of observed data



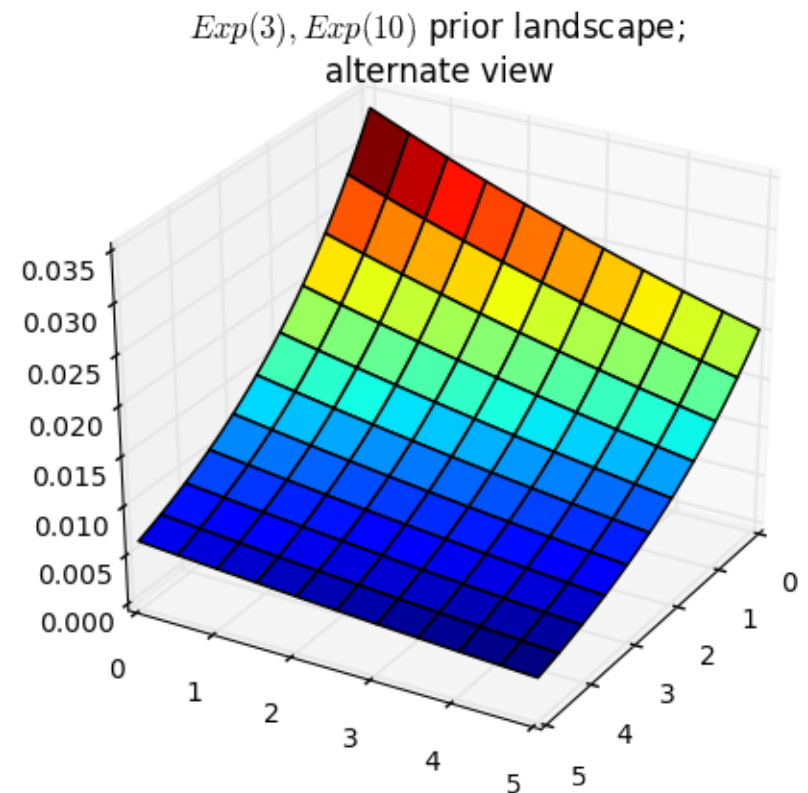
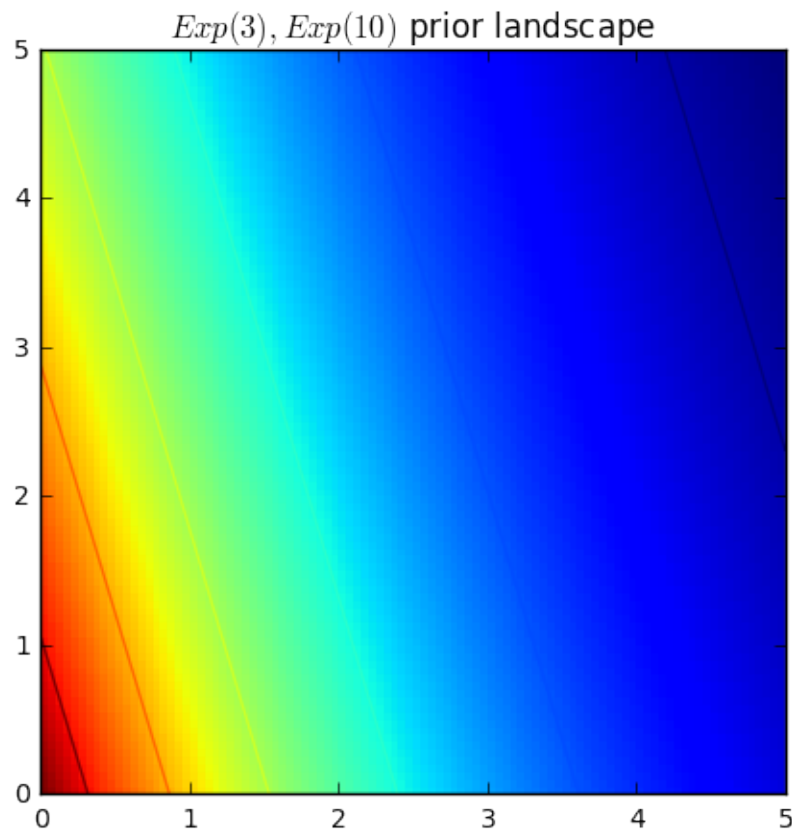
Analytical posterior





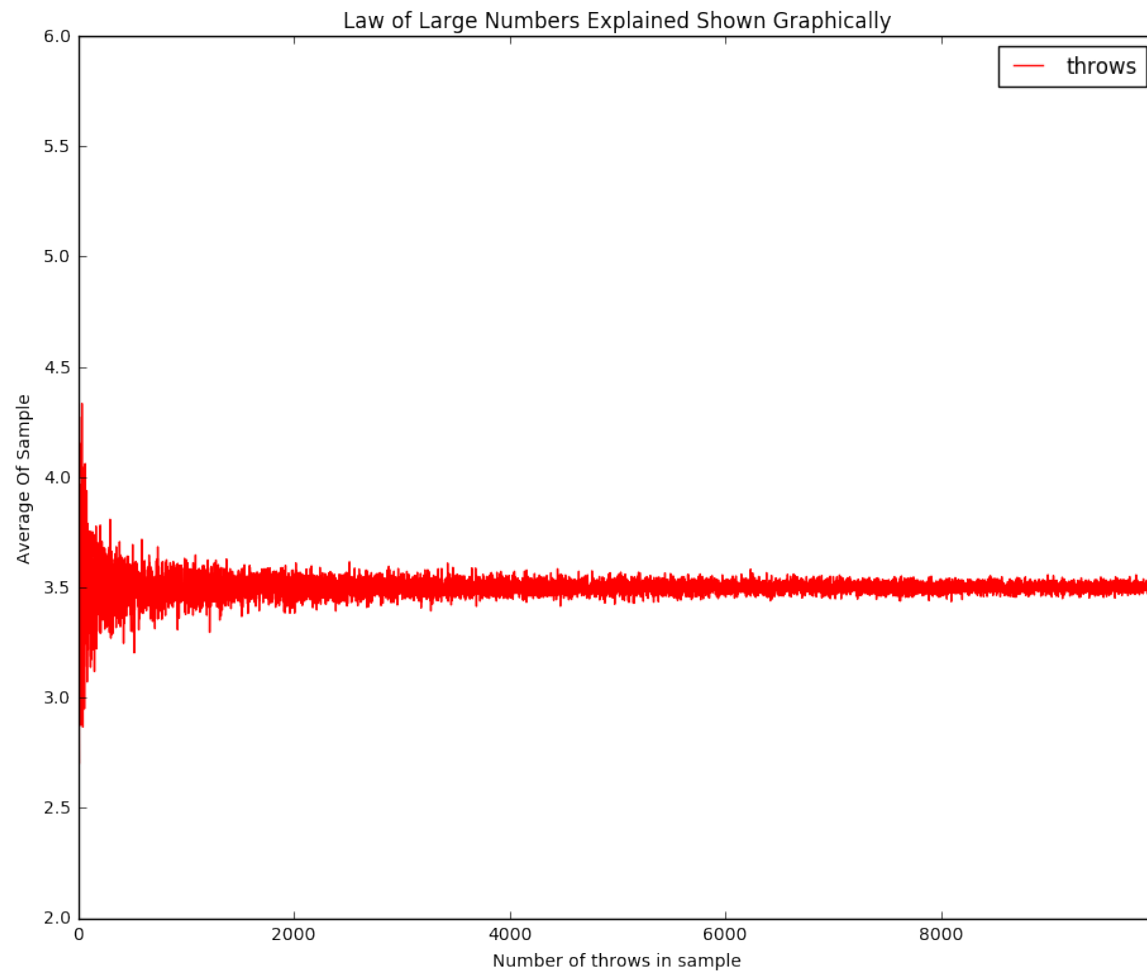
The Curse of Dimensionality

- We cannot naively search the space: any computer scientist will tell you that traversing N-dimensional space is exponentially difficult in N: the size of the space quickly blows-up as we increase N.



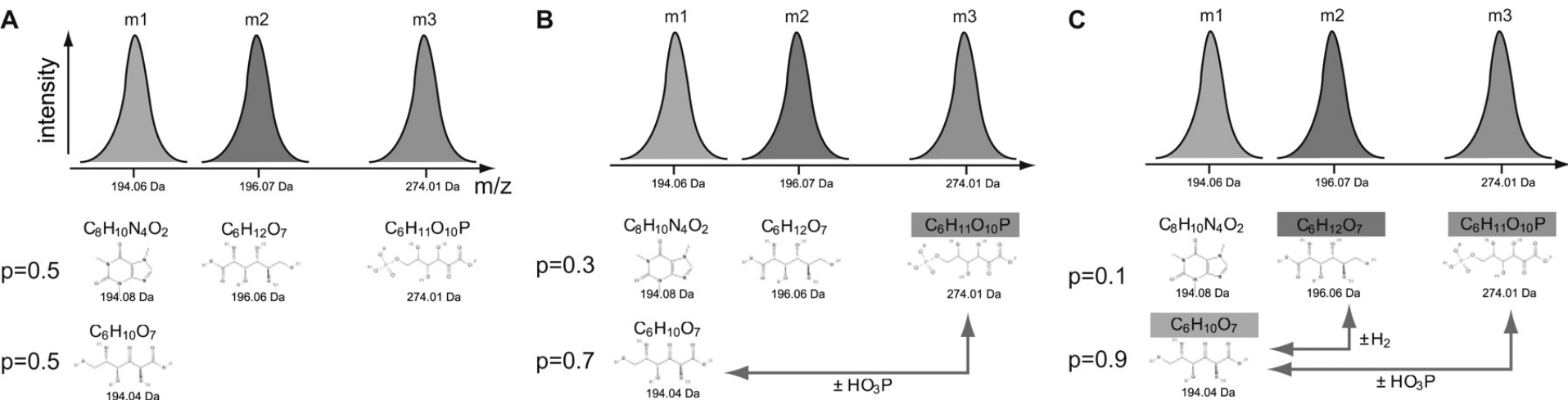
Building Intuition

- Law of Large Numbers



Example Putting It All Together

Probabilistic assignment of formulas to mass peaks in metabolomics experiments



Cartoon depicting the principle of our approach. Three peaks are observed in the mass spectrum. For mass m1 two empirical formulas are initially equally likely ($P=0.5$), as they both differ by 0.02 Da from the observed mass. Assignment of mass m3 as $C_6H_{11}O_{10}P$ provides support for the identification of m1 as $C_6H_{10}O_7$ from which it differs by one phosphorylation reaction. Assignment of m2 further increases the confidence in the assignment of m1 to $C_6H_{10}O_7$, with the posterior probability increasing to $P=0.9$.

Building Deeper Intuition, If we have time

Bayesian Point Estimate

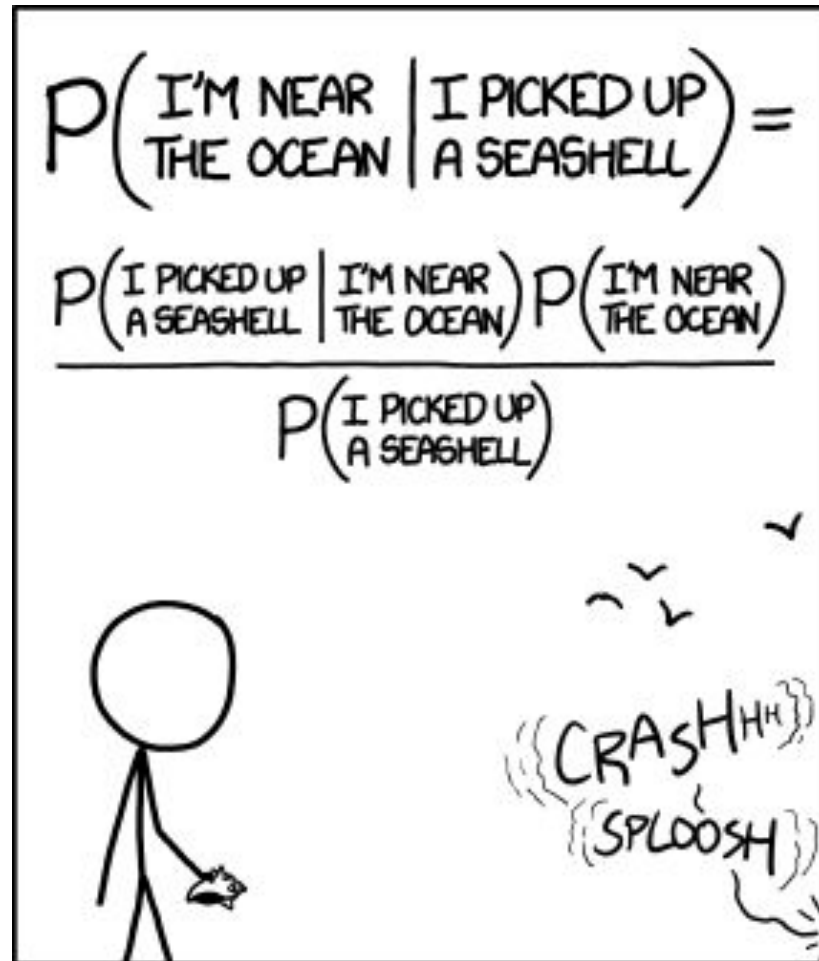
- We need to distill our posterior distribution down to a single value (or vector in the multivariate case). If the value is chosen intelligently, we can avoid the flaw of frequentist methodologies that mask the uncertainty and provide a more informative result. The value chosen, if from a Bayesian posterior, is a Bayesian point estimate.

Building Deeper Intuition, If we have time

Loss Function

- Measuring how bad our current estimation is
- So far we have been under the unrealistic assumption that we know the true parameter. Of course if we knew the true parameter, bothering to guess an estimate is pointless.
- As we have a whole distribution of what the unknown parameter could be (the posterior), we should be more interested in computing the expected loss given an estimate. This expected loss is a better estimate of the true loss than comparing the given loss from only a single sample from the posterior. (frequentist method)

Thanks for listening and Questions?



STATISTICALLY SPEAKING, IF YOU PICK UP A SEASHELL AND *DON'T* HOLD IT TO YOUR EAR, YOU CAN PROBABLY HEAR THE OCEAN.