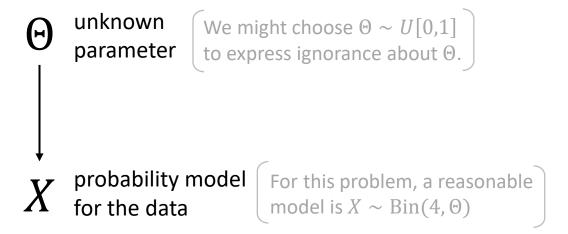
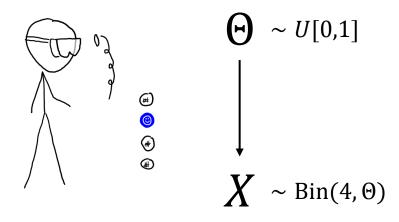


I tossed four coins and got one head. What is it reasonable to infer about the probability of heads (call it  $\theta$ )?

- "The maximum likelihood estimator is  $\hat{\theta} = 25\%$ , unjustified! thus the true probability of heads is 25%" (hence if I tossed millions more coins that's the fraction of heads I'd see)
- "All we know for certain is that  $0 < \theta < 1$ " logical, but useless!
- Let's use a random variable to express our beliefs about  $\Theta$ . Thus, we're proposing a joint model  $\Pr_{\Theta,X}$ .



Bayes's rule tells us how to update our beliefs in the light of data: simply use the conditional distribution,  $(\Theta|X=x)$ .

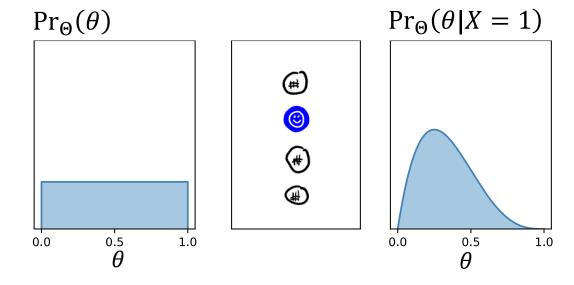


 $Pr_{\Theta}(\theta)$  is called the prior.

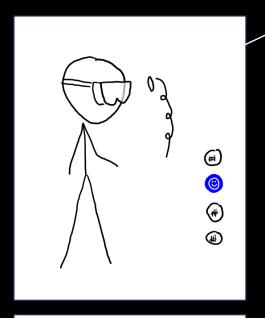
It expresses our beliefs prior to having seen the data.

 $Pr_{\Theta}(\theta|X=x)$  is called the posterior.

It expresses our beliefs about  $\Theta$  in the light of the data.

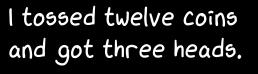


## A PARADOX ABOUT THE MEANING OF PROBABILITY



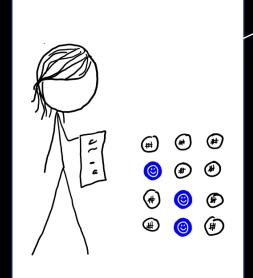
I tossed four coins and got one head.

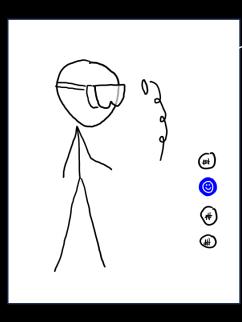
Using a  $Bin(n, \theta)$  model, I estimate the probability of heads is  $\hat{\theta} = 25\%$ 



Using a  $Bin(n,\theta)$  model, I estimate the probability of heads is  $\hat{\theta}=25\%$ 

I thought "probability" measured uncertainty.
Surely there's a difference in uncertainty between these two cases?!





I tossed four coins and got one head.

Using a  $Bin(n,\theta)$  model, I estimate the probability of heads is  $\hat{\theta}=25\%$ 

QUESTION. If we toss 100 more coins, how many heads do you predict we'd see?

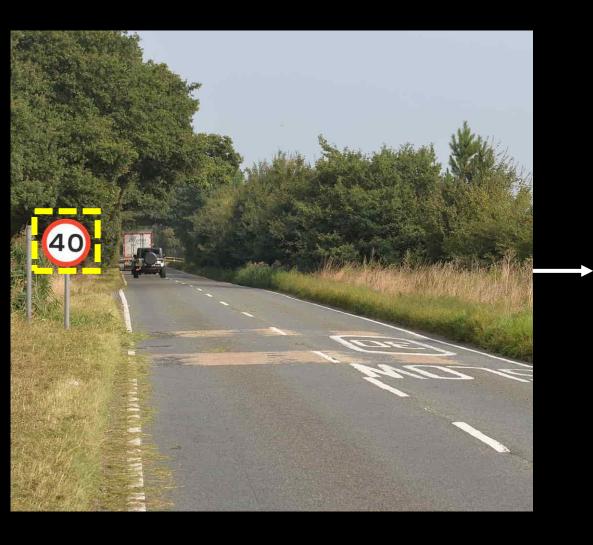
There are two types of uncertainty in our prediction:

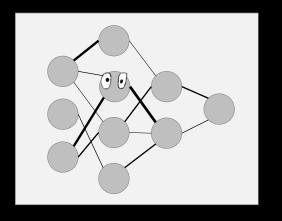
#### **EPISTEMIC UNCERTAINTY**

• We don't know  $\theta$  exactly

#### **ALEATORIC UNCERTAINTY**

 Even if we knew θ exactly, we still wouldn't know the exact number of heads we'll get





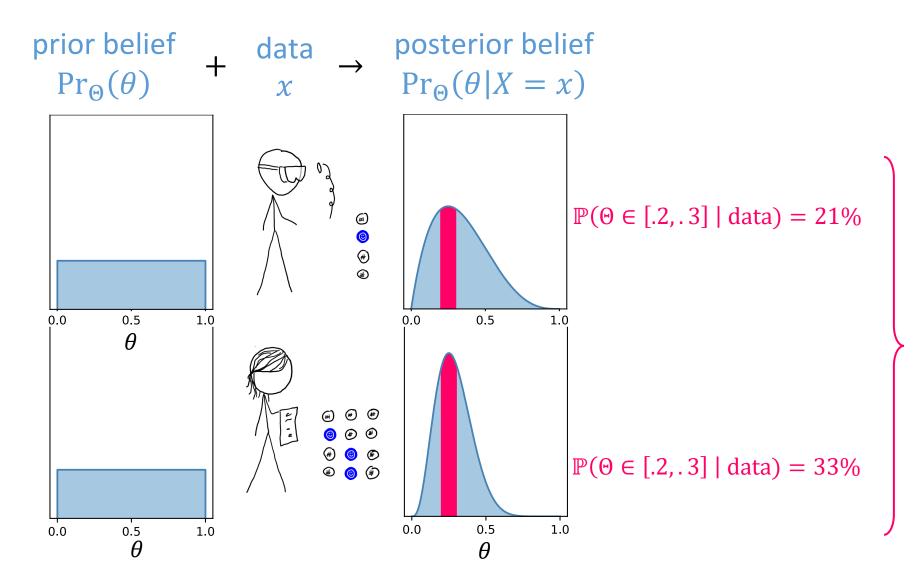
"This is a 40mph speed limit, with probability 98%."

Does it mean  $98\% \pm 0.01\%$  or  $98\% \pm 63\%$  ?

Neural networks classifiers report aleatoric *probabilities*, but they don't tell us their epistemic *confidence*.



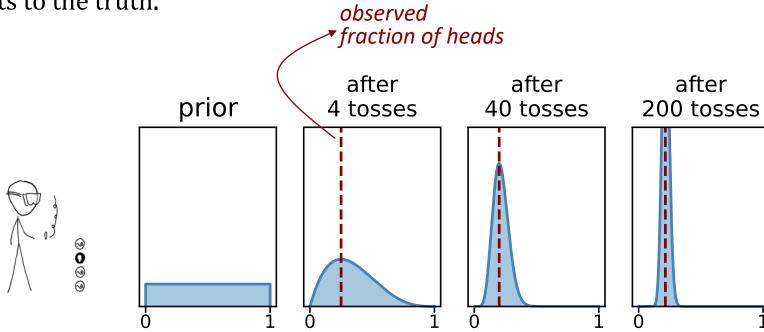
In Bayesianism, the posterior distribution *is exactly* our epistemic belief about the parameter. It depends on the amount of evidence we've seen.



A tighter posterior distribution for  $\Theta$  means we are more confident about its value.

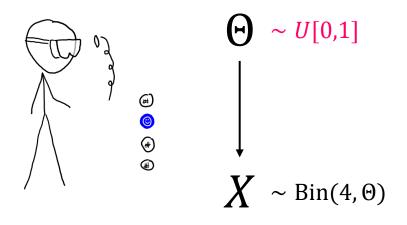


Typically, the more data we have, the closer our posterior gets to the truth.





We *must* have a prior belief about every unknown parameter. We *must* choose it before seeing the dataset in question.



But where does the prior come from?

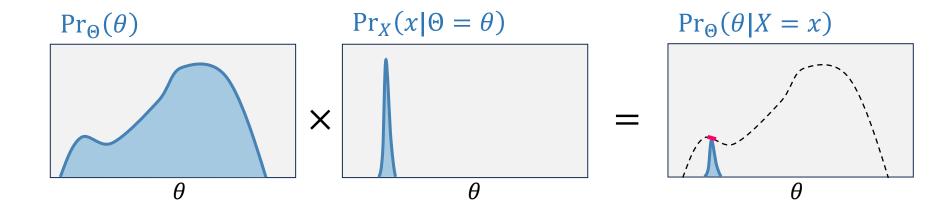
It comes from what you know already — it's how you can integrate your existing knowledge into your modelling.

If you don't have a prior to start off with, you have no business even thinking about the experiment!



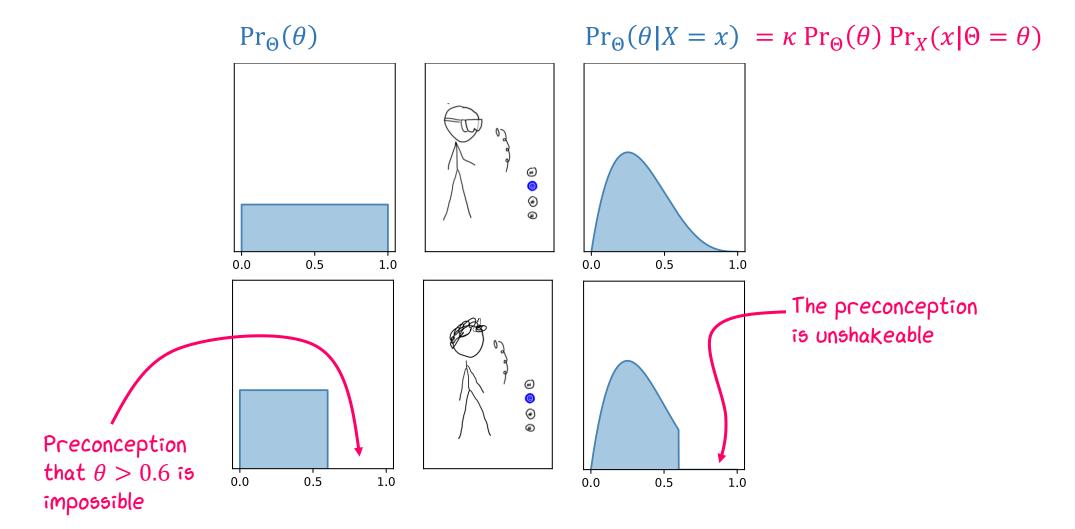
Often, with lots of data, the prior doesn't make much difference.

$$Pr_{\Theta}(\theta|X=x) = \kappa Pr_{\Theta}(\theta) Pr_X(x|\Theta=\theta)$$



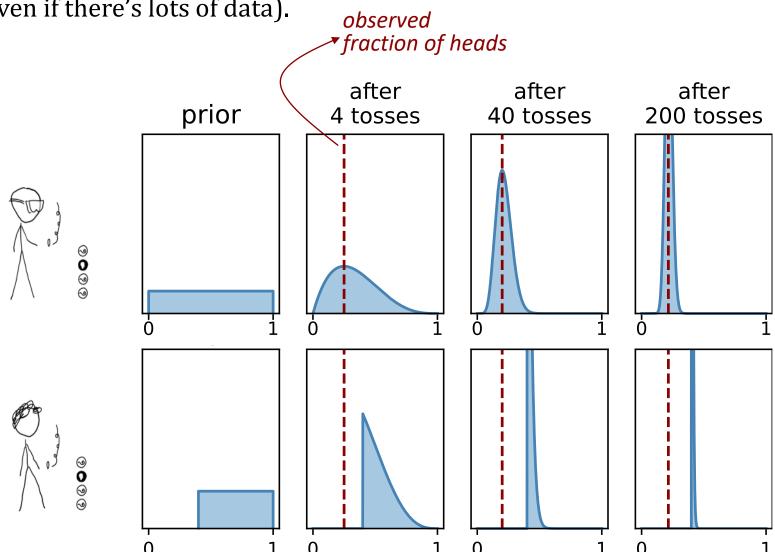


You are entitled to your own personal prior beliefs. They are entirely your choice.



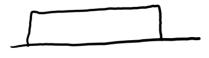


If your prior is extreme, it will be reflected in your posterior (even if there's lots of data).



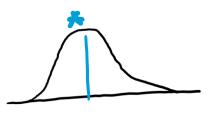
# How should we report the posterior distribution?

Prior distribution for  $\Theta$ 



Posterior distribution for  $\Theta$ 





We could report the *posterior mean*.

We could report the point with highest likelihood, the *MAP* or *maximum a-posteriori* estimate.

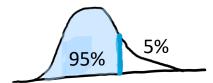
#### Example (Laplace smoothing).

We counted x successful outcomes from n trials. Using the model  $X \sim \text{Bin}(n, \Theta)$ , and the prior  $\Theta \sim U[0,1]$ , the posterior mean of  $\Theta$  is (x+1)/(n+2).



We could report a 95% confidence interval [lo,hi] such that

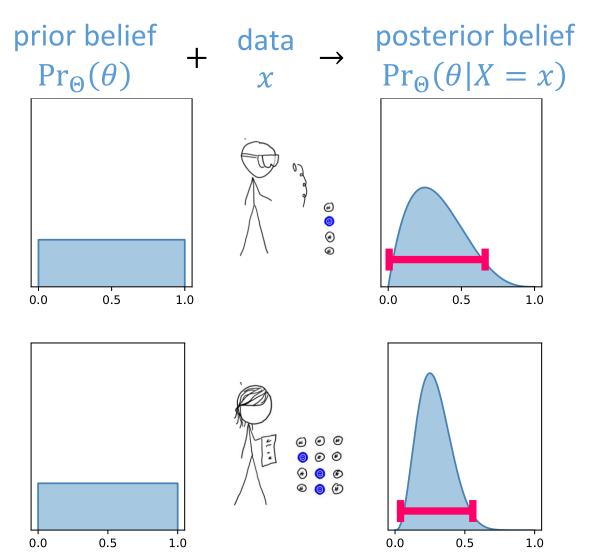
$$\mathbb{P}(\Theta < \text{lo} \mid \text{data}) = 2.5\%$$
  
 $\mathbb{P}(\Theta > \text{hi} \mid \text{data}) = 2.5\%$ 



or indeed any other 95% confidence interval e.g.

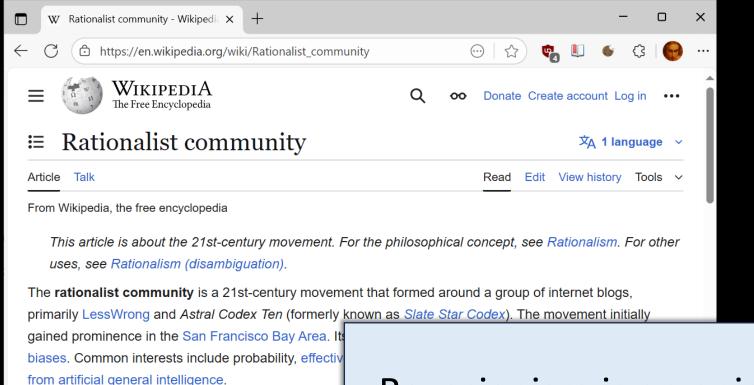
$$lo = -\infty$$
  
 $\mathbb{P}(\Theta > hi \mid data) = 5\%$ 

(though these only really work well for continuous  $\Theta$ , as for discrete  $\Theta$  we might not be able to hit the probabilities exactly)



I estimate the probability of heads is 25%, and my 95% confidence interval is [3%, 72%]

I estimate the probability of heads is 25%, and my 95% confidence interval is [12%, 51%]



Bayesianism is an epistemology, i.e. a theory of knowledge and evidence, that passes all these "smell tests".

Surely it's how any rational person should think!

groups.<sup>[1]</sup> Members who diverge from typical rationali "post-rationalist" (also known as "ingroup" and "TPOT adjacent".<sup>[3]</sup>

The borders of the rationalist community are blurry ar

#### Description [edit]

#### Rationality [edit]

Rationalists define rationality to include epistemic rationstrumental rationality — acting in a way to achieve of

The rationalists are concerned with applying science to Bayesian inference.<sup>[6]</sup> According to Ellen Huet, the unbiased, even when the conclusions are scary".<sup>[7]</sup>

To apply Bayes's rule, first write out your probability model

 $\Theta \longrightarrow X$  and state the distributions of  $\Theta$  and of  $(X; \Theta)$ . Then,

## **ALGEBRAIC BAYES**

- 1. Write out  $Pr_{\Theta}(\theta)$
- 2. Use the formula

$$\Pr_{\Theta}(\theta|X \neq x) = \kappa \Pr_{\Theta}(\theta) \Pr_{\kappa}(x|\Theta \neq \theta)$$
  
then find  $\kappa$  to make this integrate to 1

... but these are usually intractable

This lets us calculate confidence intervals:

$$\mathbb{P}(\Theta \in \text{range}|X \neq x) = \int_{\Theta \in \text{range}} \mathbb{P}(\Theta(\theta|X = x)) d\theta$$

One way to do

## COMPUTATIONAL BAYES

- 1. Generate a sample  $(\theta_1, ..., \theta_n)$  from  $\Theta$
- 2. Compute weights

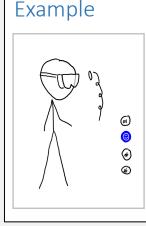
$$w_i = \Pr_X(x|\Theta = \theta_i),$$

then rescale weights to sum to one

$$\mathbb{P}(\Theta \in \text{range} | X \neq x) \approx \sum_{i=1}^{n} w_i \mathcal{I}_{\theta_i} \in \text{range}$$

It's more elegant to use the generalized version

$$\mathbb{E}[h(\Theta)|X=x] \approx \Sigma_i w_i h(\theta_i)$$



I got x=1 head out of n=4 coin tosses. I propose the probability model  $X \sim \text{Bin}(n,\Theta)$ . I don't know  $\Theta$ , so I'll treat it as a random variable,  $\Theta \sim U[0,1]$ .

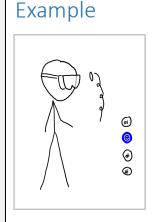
Plot the distribution of  $(\Theta|X=x)$ .

Likelihood of the data:  $\times \sim Bin(n,8)$   $Pr_{\times}(\times \mid 0=0) = \binom{n}{x} \theta^{\times} (1-\theta)^{n-x} = 4 \theta (1-\theta)^{3}$ 

- 1. Generate a sample  $(\theta_1, ..., \theta_n)$  from  $\Theta$ :

  Osamp = np. random. uniform (0, 1, size = (00000))
- 2. Compute weights  $w_i = \Pr_X(x|\Theta = \theta_i)$ ,  $w = 4 \times (0 \text{ samp *x 1}) \times ((1 0 \text{ samp}) \text{ *x 3})$  then rescale weights to sum to one:  $w = w \setminus \text{Asum}(w)$

Reason about  $(\Theta|X=x)$  indirectly, using  $\mathbb{E}[h(\Theta)|X=x] \approx \Sigma_i w_i h(\theta_i)$ 



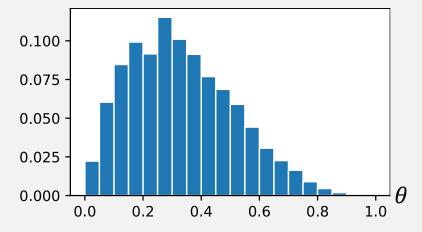
I got x=1 head out of n=4 coin tosses. I propose the probability model  $X \sim \text{Bin}(n,\Theta)$ . I don't know  $\Theta$ , so I'll treat it as a random variable,  $\Theta \sim U[0,1]$ .

Plot the distribution of  $(\Theta|X=x)$ .

Reason about  $(\Theta|X=x)$  indirectly, using  $\mathbb{E}[h(\Theta)|X=x] \approx \Sigma_i w_i h(\theta_i)$ 

## Let's plot a histogram:

we'll split the range of  $\theta$  into bins, and for each  $\theta$ -bin, we'll draw a bar of height  $\mathbb{P}(\theta \in \text{bin } | X = x)$ 



P ( O e bin (duta)

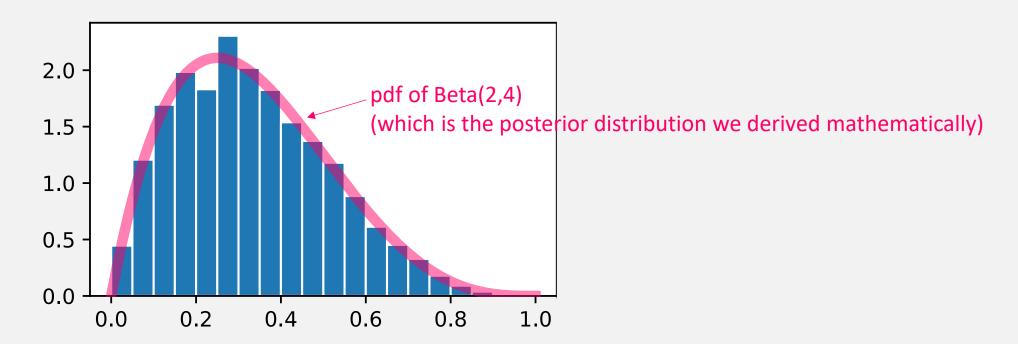
In code, to plot all bins: plt.hist( $\theta$ samp, weights=w)

bar height = 
$$\frac{\mathbb{P}(\Theta \in \text{bin}|\text{data})}{\text{bin width}}$$

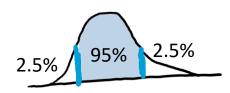
For continuous random variables, I prefer to scale the bar heights so that the total *area* is 1. This is known as a *density histogram*.

It means the plot is directly comparable to a pdf.

plt.hist(θsamp, weights=w) density=True)



## How do we compute confidence intervals from a weighted sample?



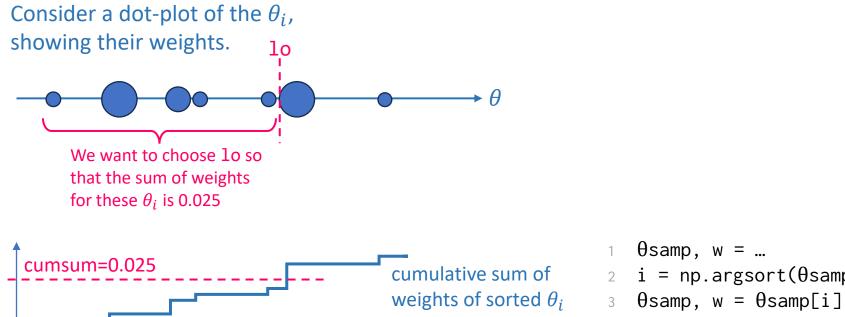
We'd like to report a 95% confidence interval [lo, hi] such that

$$\mathbb{P}(\Theta < 1o \mid data) = 2.5\%$$

$$\mathbb{P}(\Theta > \text{hi} \mid \text{data}) = 2.5\%$$

We can find lo via a computational Bayes estimate (and hi similarly):

$$\mathbb{P}(\Theta < \text{lo} \mid \text{data}) \approx \sum_{i} w_i \, 1_{\theta_i < \text{lo}}$$



- $i = np.argsort(\theta samp)$
- $\theta$ samp, w =  $\theta$ samp[i], w[i]
- F = np.cumsum(w)
- $lo = \theta samp[F<0.025][-1]$

## Exercise 8.3.2 (Multiple unknowns)

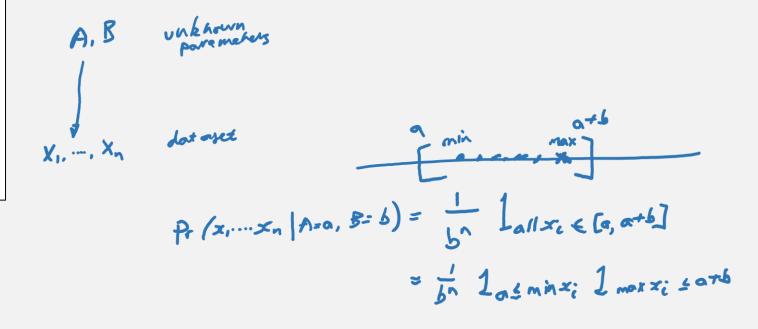
We have a dataset  $[x_1, ..., x_n]$ . We propose to model it as independent samples from U[A, A + B], where A and B are unknown parameters.

Using  $A \sim \text{Exp}(0.5)$  and  $B \sim \text{Exp}(1.0)$  as prior distributions for the unknown parameters, find the distribution of (B|data).

1. Generate a sample  $(\theta_1, \dots, \theta_n)$  from  $\theta$ 

2. Compute weights  $w_i = \Pr_{X}(x|\theta - \theta_i)$ . then rescale weights to sum to one

Reason about  $(\Theta|X=x)$  indirectly



```
# The dataset
x = [2, 3, 2.1, 2.4, 3.14, 1.8]
# Step 1
asamp = np.random.exponential(scale=1/0.5, size=100000)
bsamp = np.random.exponential(scale=1/1.0, size=100000)
#absamp = zip(asamp, bsamp)
# Step 2
w = 1/bsamp**len(x) * np.where((asamp <= min(x)) & (max(x) <= asamp+bsamp), 1, 0)
w = w / np.sum(w)
# Plot the posterior distribution of B
plt.hist(bsamp, weights=w, density=True)
```

## Exercise 8.3.2 (Multiple unknowns)

We have a dataset  $[x_1, ..., x_n]$ . We propose to model it as independent samples from U[A, A + B], where A and B are unknown parameters.

Using  $A \sim \text{Exp}(0.5)$  and  $B \sim \text{Exp}(1.0)$  as prior distributions for the unknown parameters, find the distribution of (B|data).

**TIP.** If n is large, we can run into underflow problems when we compute  $\Pr(x_1, ..., x_n | \text{params})$  directly.

Fix: be clever about rescaling the weights, using the log-sum-exp trick (exercise 8.3.4).

**TIP.** First find the joint posterior distribution for *all* the unknown parameters. Then, just report on the readout we're interested in.

In maths,

$$Pr_B(b|data) = \int_a Pr_{A,B}(a,b|data) da$$

### Computationally,

- 1. Generate samples  $(a_i, b_i)$  from the joint prior
- 2. Compute a weight  $w_i$  for each pair
- 3. Plot a weighted histogram of just the  $b_i$

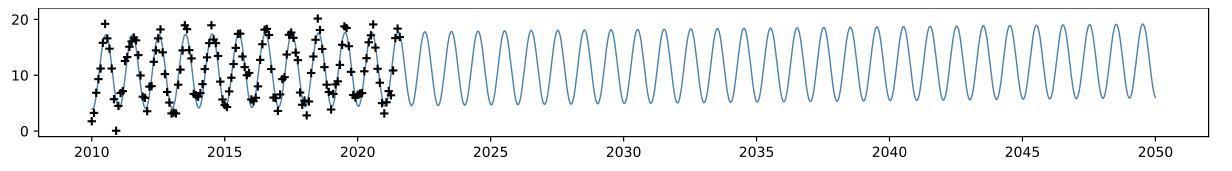
This is called *marginalization*.

# Why does computational Bayes work?

It's a simple argument — see §6.2 of lecture notes.

Consider the dataset of monthly average temperatures in Cambridge.

Proposed model: Temp  $\sim \alpha + \beta \sin(2\pi(t+\phi)) + \gamma(t-2000) + N(0,\sigma^2)$ 



If we fit this model we get the maximum likelihood estimate  $\hat{\gamma} = 0.027$  °C/year.

How **confident** are we about this value?

## Climate confidence challenge.

Find a 95% confidence interval for  $\gamma$ , for Cambridge from 1985 to the present. (Use your own priors for the unknowns.)

Please submit your answer on Moodle.