## Lecture 1: outline and introduction to inference

Ben Lambert[1]
ben.lambert@some.ox.ac.uk

[1]Somerville College
University of Oxford

October 26, 2016

# Outline

1. Logistics

2. Course outline

3. The theory and practice of inference
   - A conceptual introduction to inference
   - Frequentist and Bayesian world views
   - Understanding probability distributions
   - A short introduction to Bayes' rule for inference

# Who am I?



L. B. & S. C. R. STATION, TUNBRIDGE WELLS.

# Who am I?

- Researcher in epidemiology in the department of Zoology.

# Who am I?

- Researcher in epidemiology in the department of Zoology.
- Used Bayesian statistics for the past 7 years.



L. B. & S. C. R. STATION, TUNBRIDGE WELLS.

# Who am I?

- Researcher in epidemiology in the department of Zoology.
- Used Bayesian statistics for the past 7 years.
- Born in the same town as Thomas Bayes (Tunbridge Wells.)



L. B. & S. C. R. STATION, TUNBRIDGE WELLS.

# Course outline

Target audience:

Prerequisites:

Target audience:

- Researchers who want to apply statistical inference in their work.

Prerequisites:

Target audience:

- Researchers who want to apply statistical inference in their work.

Prerequisites:

- A basic knowledge of mathematical programming in R, Matlab, Mathematica, Python, C++, or similar.

## Course outline

Target audience:

- Researchers who want to apply statistical inference in their work.

Prerequisites:

- A basic knowledge of mathematical programming in R, Matlab, Mathematica, Python, C++, or similar.
- If you're not comfortable with calculus, don't worry. However, it might be worth looking at an A-level textbook to brush up your skills.

# Lecture timetable

Every Wednesday at 2pm. Problem class starting at 3pm/3.15pm.

# Problem class format

- No class today!

## Problem class format

- No class today!
- (Usually) immediately follow a lecture.

## Problem class format

- No class today!
- (Usually) immediately follow a lecture.
- Happen upstairs in the 'IT suite'.

## Problem class format

- No class today!
- (Usually) immediately follow a lecture.
- Happen upstairs in the 'IT suite'.
- A mix of applied and theoretical problems (mostly applied.)

## Problem class format

- No class today!
- (Usually) immediately follow a lecture.
- Happen upstairs in the 'IT suite'.
- A mix of applied and theoretical problems (mostly applied.)
- A few demonstrators to help with specific problem issues.

## Problem class format

- No class today!
- (Usually) immediately follow a lecture.
- Happen upstairs in the 'IT suite'.
- A mix of applied and theoretical problems (mostly applied.)
- A few demonstrators to help with specific problem issues.
- Restricted to 56 people. Sign up for the these with sheet at front. If sign ups exceed places we will have a ballot.

## Problem class format

- No class today!
- (Usually) immediately follow a lecture.
- Happen upstairs in the 'IT suite'.
- A mix of applied and theoretical problems (mostly applied.)
- A few demonstrators to help with specific problem issues.
- Restricted to 56 people. Sign up for the these with sheet at front. If sign ups exceed places we will have a ballot.
- (However I will be putting the problem sets online...)

- Work in groups to reproduce and extend published results.

# Hackathon: 7th December

- Work in groups to reproduce and extend published results.
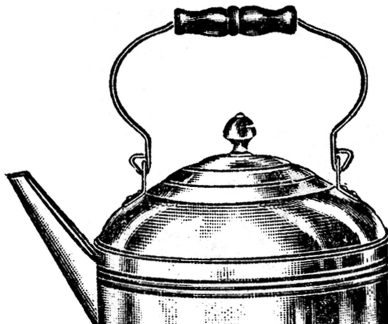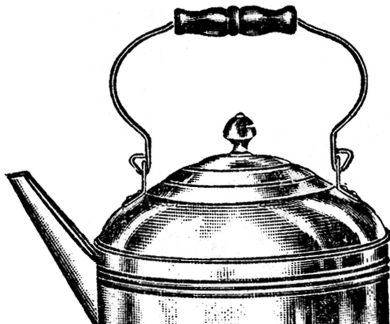- Alternatively, can work on your own research problem.

# A couple of things

- Lecture notes available from
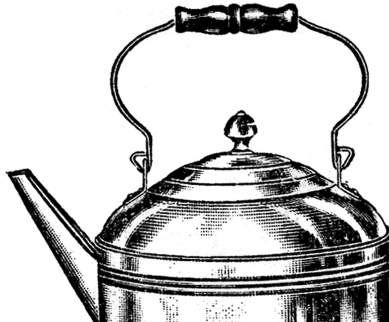  `www.ben-lambert.com/bayesian-lecture-slides/`.

# A couple of things

- Lecture notes available from
  www.ben-lambert.com/bayesian-lecture-slides/.
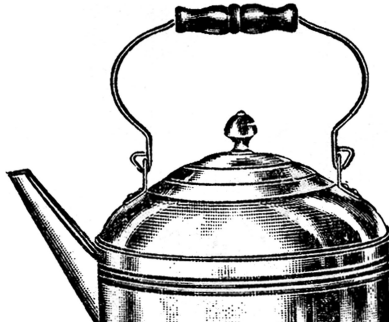- Printed lecture notes from weeks 2-7.

## A couple of things

- Lecture notes available from
  www.ben-lambert.com/bayesian-lecture-slides/.
- Printed lecture notes from weeks 2-7.
- No drink/food upstairs in IT suite.

## A couple of things

- Lecture notes available from
  www.ben-lambert.com/bayesian-lecture-slides/.
- Printed lecture notes from weeks 2-7.
- No drink/food upstairs in IT suite.
- We will sort logins for everyone at the beginning of class
  next week.

- Simple and intuitive model building (unlike frequentist statistics there is no need to remember lots of specific formulae).

## Tangible benefits of Bayesian inference

- Simple and intuitive model building (unlike frequentist statistics there is no need to remember lots of specific formulae).
- Exhaustive and creative model testing.
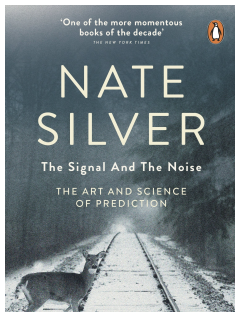
## Tangible benefits of Bayesian inference

- Simple and intuitive model building (unlike frequentist statistics there is no need to remember lots of specific formulae).
- Exhaustive and creative model testing.
- Straightforward interpretation of results.

- Simple and intuitive model building (unlike frequentist statistics there is no need to remember lots of specific formulae).
- Exhaustive and creative model testing.
- Straightforward interpretation of results.
- The best predictions (for example, Nate Silver's correct prediction of 2008 US Presidential election results.)

# Tangible benefits of Bayesian inference

- Simple and intuitive model building (unlike frequentist statistics there is no need to remember lots of specific formulae).
- Exhaustive and creative model testing.
- Straightforward interpretation of results.
- The best predictions (for example, Nate Silver's correct prediction of 2008 US Presidential election results.)

- Existing texts put a strong emphasis on its (seemingly) complex mathematical basis.

- Existing texts put a strong emphasis on its (seemingly) complex mathematical basis.
- Poor explanation of *why* we need MCMC algorithms.

# Why don't more people use Bayesian inference?

- Existing texts put a strong emphasis on its (seemingly) complex mathematical basis.
- Poor explanation of *why* we need MCMC algorithms.
- Poor explanation of *how* these MCMC algorithms work, and how to implement them in practice.

# Why don't more people use Bayesian inference?

- Existing texts put a strong emphasis on its (seemingly) complex mathematical basis.
- Poor explanation of *why* we need MCMC algorithms.
- Poor explanation of *how* these MCMC algorithms work, and how to implement them in practice.
- The view that Bayesian inference is more wishy-washy than frequentist inference.

# Why don't more people use Bayesian inference?

- Existing texts put a strong emphasis on its (seemingly) complex mathematical basis.
- Poor explanation of *why* we need MCMC algorithms.
- Poor explanation of *how* these MCMC algorithms work, and how to implement them in practice.
- The view that Bayesian inference is more wishy-washy than frequentist inference.

$\implies$ this course aims to avoid the first three issues.

# Why don't more people use Bayesian inference?

- Existing texts put a strong emphasis on its (seemingly) complex mathematical basis.
- Poor explanation of *why* we need MCMC algorithms.
- Poor explanation of *how* these MCMC algorithms work, and how to implement them in practice.
- The view that Bayesian inference is more wishy-washy than frequentist inference.

$\implies$ this course aims to avoid the first three issues. (We will cover the last point in this lecture.)

likelihood

# Syllabus

likelihood + prior

likelihood + prior  $\xrightarrow{\text{Bayes' rule}}$

# Syllabus
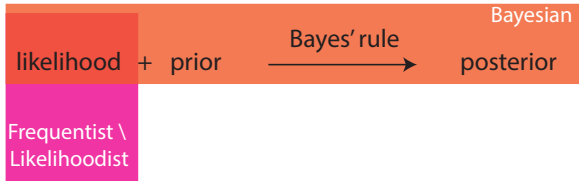
$$\text{likelihood} \; + \; \text{prior} \quad \xrightarrow{\text{Bayes' rule}} \quad \text{posterior}$$

# Syllabus

likelihood + prior $\xrightarrow{\text{Bayes' rule}}$ posterior

Bayesian

# Syllabus

# Syllabus

Lecture 1: The theory of inference, today



likelihood + prior →(Bayes' rule) Bayesian posterior

Frequentist \ Likelihoodist

# Syllabus

likelihood + prior  →  Bayes' rule  →  Bayesian posterior

# Syllabus

likelihood + prior $\xrightarrow{\text{Bayes' rule}}$ Bayesian posterior

Lecture 2: Analytic Bayesian inference
2nd November

# Syllabus

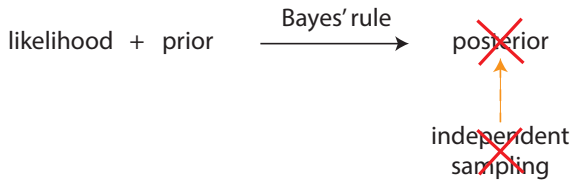likelihood + prior $\xrightarrow{\text{Bayes' rule}}$ posterior

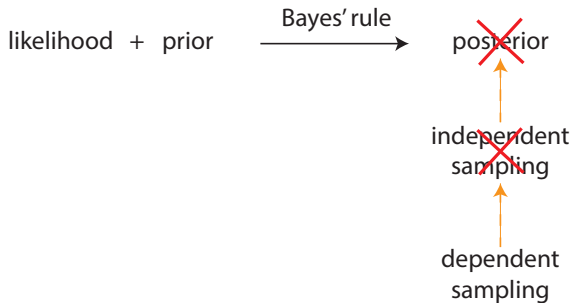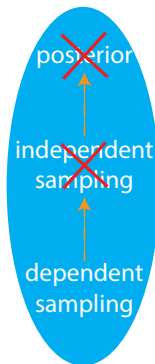likelihood + prior $\xrightarrow{\text{Bayes' rule}}$ pos~~t~~erior

likelihood + prior   —Bayes' rule→   posterior

independent sampling

likelihood + prior → Bayes' rule → posterior

independent sampling

dependent sampling, via MCMC:

a. Random Walk Metropolis

likelihood  +  prior  $\xrightarrow{\text{Bayes' rule}}$  posterior

independent
sampling

dependent
sampling, via MCMC:

a. Random Walk Metropolis

likelihood + prior  → Bayes' rule →  pos~~te~~rior

independent sam~~p~~ling

dependent sampling, via MCMC:

a. Random Walk Metropolis

increasing approximation quality

likelihood  +  prior  →  Bayes' rule  →  posterior

independent sampling

dependent sampling, via MCMC:
a. Random Walk Metropolis
b. Gibbs

increasing approximation quality

# Syllabus

likelihood  +  prior  $\xrightarrow{\text{Bayes' rule}}$  pos~~t~~erior

independent
sampling

dependent
sampling, via MCMC:

a. Random Walk Metropolis
b. Gibbs
c. Hamiltonian Monte Carlo

Lecture 5: Modern
MCMC methods
23rd November

increasing
approximation
quality

likelihood + prior →(Bayes' rule)→ posterior ✗

independent sampling ✗

dependent sampling, via MCMC:

a. Random Walk Metropolis
b. Gibbs
c. Hamiltonian Monte Carlo

increasing approximation quality

likelihood + prior $\xrightarrow{\text{Bayes' rule}}$ pos~~te~~rior

independent
sam~~pli~~ng

dependent
sampling, via MCMC:

a. Random Walk Metropolis
b. Gibbs

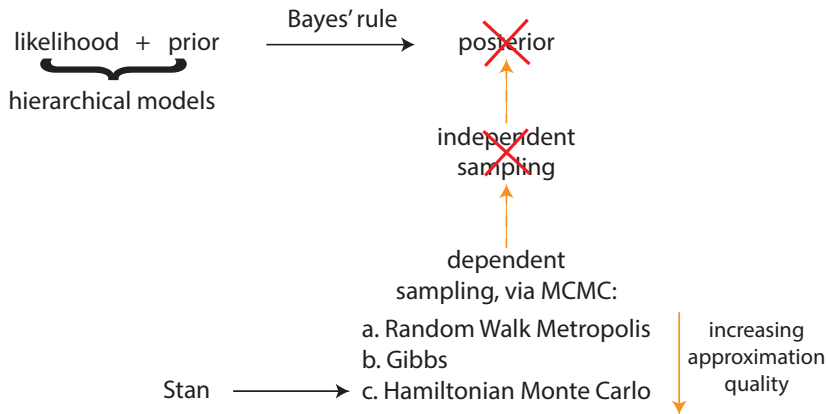Stan $\longrightarrow$ c. Hamiltonian Monte Carlo

increasing
approximation
quality

likelihood + prior $\xrightarrow{\text{Bayes' rule}}$ post~~erior~~

hierarchical models

independent sampling

dependent sampling, via MCMC:

a. Random Walk Metropolis
b. Gibbs
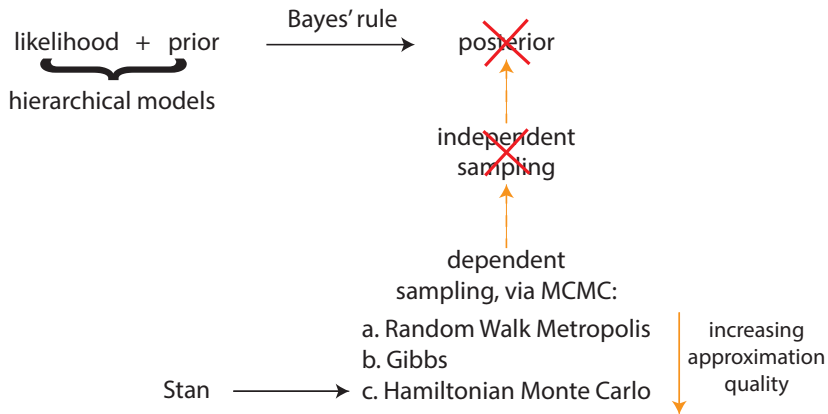Stan $\longrightarrow$ c. Hamiltonian Monte Carlo

increasing approximation quality
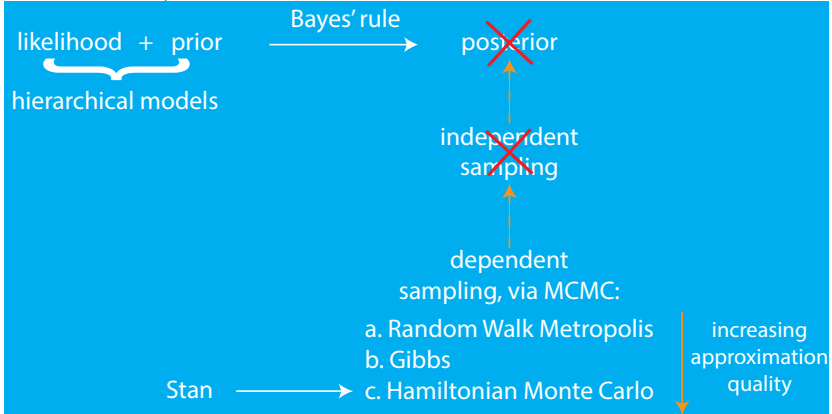
Lecture 7: further applied Bayesian inference
and hackathon, 7th December

By the end of this course you should:

By the end of this course you should:

- Understand the basic theory and motivation of Bayesian inference.

By the end of this course you should:

- Understand the basic theory and motivation of Bayesian inference.
- Know how to critically assess a statistical model.

## Course outcomes

By the end of this course you should:

- Understand the basic theory and motivation of Bayesian inference.
- Know how to critically assess a statistical model.
- Appreciate why we often need to use sampling in Bayesian inference.

## Course outcomes

By the end of this course you should:

- Understand the basic theory and motivation of Bayesian inference.
- Know how to critically assess a statistical model.
- Appreciate why we often need to use sampling in Bayesian inference.
- Grasp how modern MCMC algorithms work intuitively and how to implement these in practice.

## Course outcomes

By the end of this course you should:

- Understand the basic theory and motivation of Bayesian inference.
- Know how to critically assess a statistical model.
- Appreciate why we often need to use sampling in Bayesian inference.
- Grasp how modern MCMC algorithms work intuitively and how to implement these in practice.
- Know how to code up most models in *Stan*.

## Course outcomes

By the end of this course you should:

- Understand the basic theory and motivation of Bayesian inference.
- Know how to critically assess a statistical model.
- Appreciate why we often need to use sampling in Bayesian inference.
- Grasp how modern MCMC algorithms work intuitively and how to implement these in practice.
- Know how to code up most models in *Stan*.
- Recognise the benefits of hierarchical models and how these can be used to provide robust inferences.

## Lecture outcomes

By the end of this lecture you should:

By the end of this lecture you should:

1. Understand the motivation behind inference.

## Lecture outcomes

By the end of this lecture you should:

1. Understand the motivation behind inference.
2. Appreciate the similarities and differences between Frequentist and Bayesian approaches to inference.

## Lecture outcomes

By the end of this lecture you should:

1. Understand the motivation behind inference.
2. Appreciate the similarities and differences between Frequentist and Bayesian approaches to inference.
3. Know how to manipulate probability distributions.

1. Consider an observable characteristic we are trying to explain, for example the heights of 5 randomly chosen individuals.

## The Big world

1. Consider an observable characteristic we are trying to explain, for example the heights of 5 randomly chosen individuals.

2. Assume that there exists a true process $T$ that generates the heights of all individuals in our sample.

## The Big world

1. Consider an observable characteristic we are trying to explain, for example the heights of 5 randomly chosen individuals.

2. Assume that there exists a true process $T$ that generates the heights of all individuals in our sample.

3. There is variability in the observables outputted by $T$; this can either be ontological (for example due to the inherent variability in picking our random sample), or epistemological (for example, because we lack knowledge of the genetics and environmental factors that affect growth).

## The Big world

1. Consider an observable characteristic we are trying to explain, for example the heights of 5 randomly chosen individuals.

2. Assume that there exists a true process $T$ that generates the heights of all individuals in our sample.

3. There is variability in the observables outputted by $T$; this can either be ontological (for example due to the inherent variability in picking our random sample), or epistemological (for example, because we lack knowledge of the genetics and environmental factors that affect growth).

4. Imagine a set of all conceivable processes that could result in our sample of height observations, which we call the "Big World".
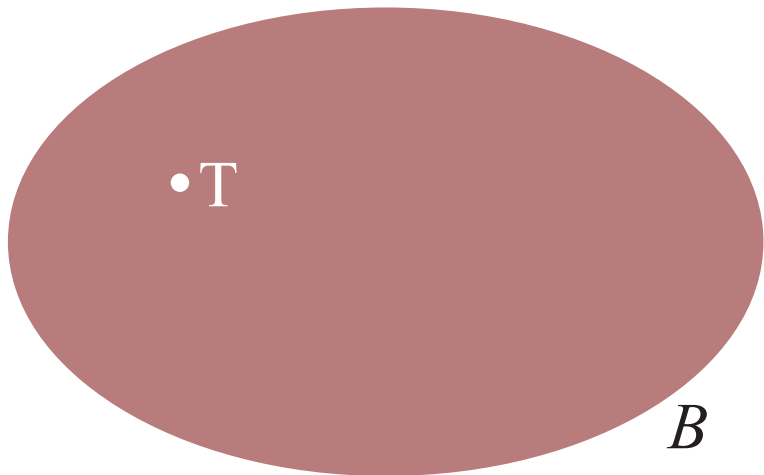
Figure: Images adapted from "A Technical Introduction to Probability and Bayesian Inference for Stan Users", *Stan Development Team*, 2016.

- **Motivation:** update our knowledge of $T$ in light of data, and use the updated knowledge to estimate quantities of interest.

- **Motivation:** update our knowledge of $T$ in light of data, and use the updated knowledge to estimate quantities of interest.
  - In our height example we might want to estimate the mean height of the entire population having witnessed our sample of 5 individuals.

- **Motivation:** update our knowledge of $T$ in light of data, and use the updated knowledge to estimate quantities of interest.
    - In our height example we might want to estimate the mean height of the entire population having witnessed our sample of 5 individuals.
- **Method:**

- **Motivation:** update our knowledge of $T$ in light of data, and use the updated knowledge to estimate quantities of interest.
  - In our height example we might want to estimate the mean height of the entire population having witnessed our sample of 5 individuals.
- **Method:**
  1. Find areas of the Big World that are closest to $T$; ideally we would find $T$ itself!

# What is inference?

- **Motivation:** update our knowledge of $T$ in light of data, and use the updated knowledge to estimate quantities of interest.
  - In our height example we might want to estimate the mean height of the entire population having witnessed our sample of 5 individuals.

- **Method:**
  1. Find areas of the Big World that are closest to $T$; ideally we would find $T$ itself!
  2. Estimate quantities of interest using these subsets of the Small World.

1. The infinity of the Big World is too large to be useful.

## The Small world

1. The infinity of the Big World is too large to be useful.
2. Instead we first consider a subset of possible data generating processes which we call the "Small World", or $\Theta$.
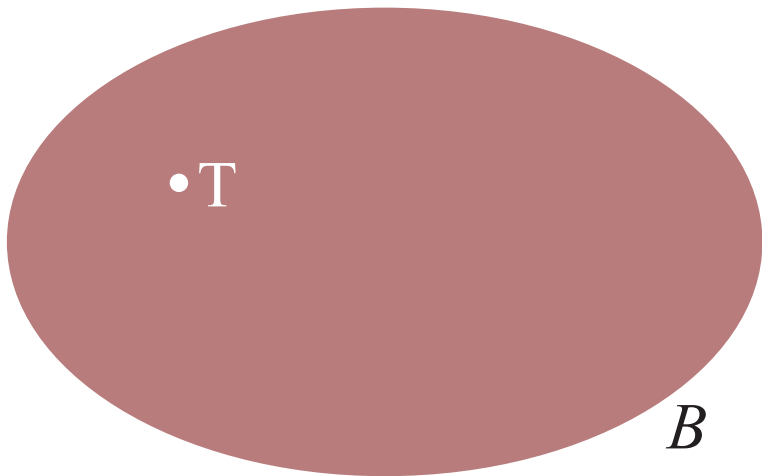
## The Small world

1. The infinity of the Big World is too large to be useful.
2. Instead we first consider a subset of possible data generating processes which we call the "Small World", or $\Theta$.
3. The Small World corresponds to a single probability model framework; in our height example we might suppose that $H \sim N(\mu, \sigma)$, where $\mu$ is the mean height, and $\sigma$ is their standard deviation.
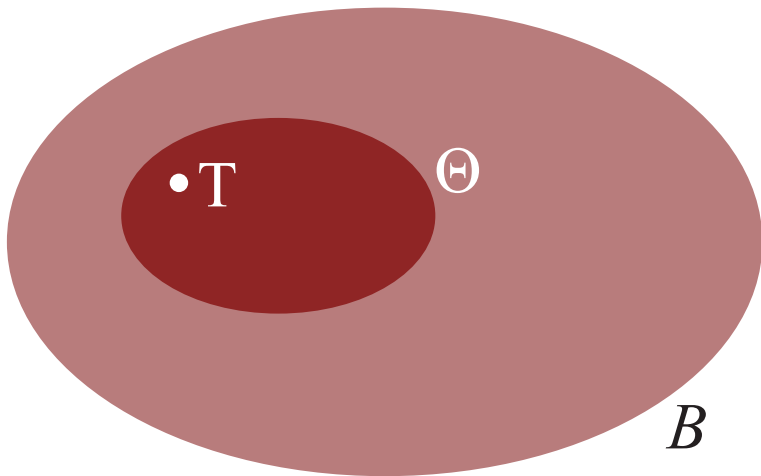
## The Small world

1. The infinity of the Big World is too large to be useful.
2. Instead we first consider a subset of possible data generating processes which we call the "Small World", or $\Theta$.
3. The Small World corresponds to a single probability model framework; in our height example we might suppose that $H \sim N(\mu, \sigma)$, where $\mu$ is the mean height, and $\sigma$ is their standard deviation.
4. By varying our parameters $\theta = (\mu, \sigma)$ we get different data generating processes.

## The Small world

1. The infinity of the Big World is too large to be useful.
2. Instead we first consider a subset of possible data generating processes which we call the "Small World", or $\Theta$.
3. The Small World corresponds to a single probability model framework; in our height example we might suppose that $H \sim N(\mu, \sigma)$, where $\mu$ is the mean height, and $\sigma$ is their standard deviation.
4. By varying our parameters $\theta = (\mu, \sigma)$ we get different data generating processes.
5. The collection of probability distributions we get by varying $\theta \subset \Theta$ in the Small World is known as the *Likelihood*.
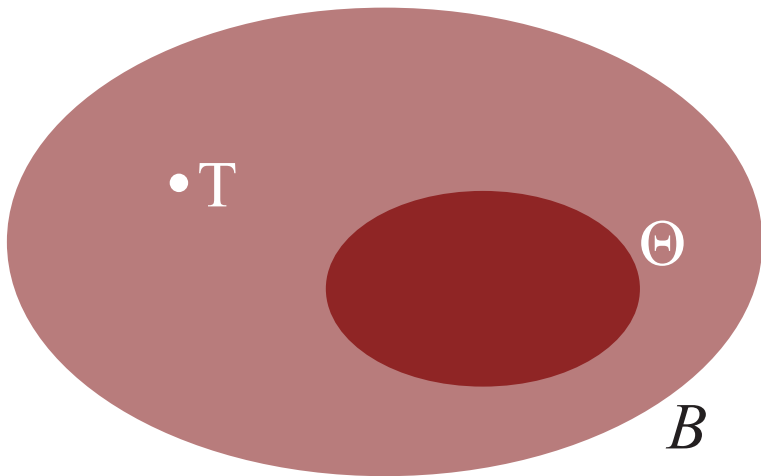
1. The Small World is still too big for our purposes.

## The prior

1. The Small World is still too big for our purposes.
2. We usually have *some* knowledge about which areas of the Small World are nearest to $T$. For example we don't believe that $\mu = 100m$ and $\mu = 1.5m$ are equally probable.
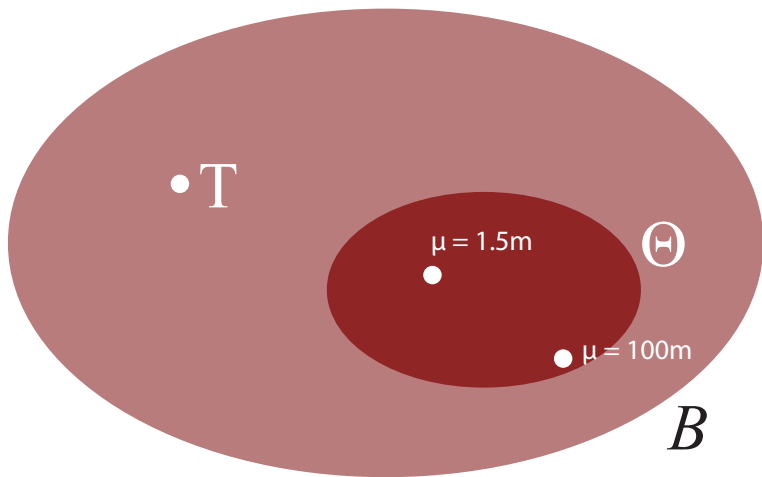
## The prior

1. The Small World is still too big for our purposes.

2. We usually have *some* knowledge about which areas of the Small World are nearest to $T$. For example we don't believe that $\mu = 100m$ and $\mu = 1.5m$ are equally probable.

3. As such, in Bayesian inference we define a *prior* probability density that gives a weighting to all $\theta \in \Theta$ reflecting our beliefs.
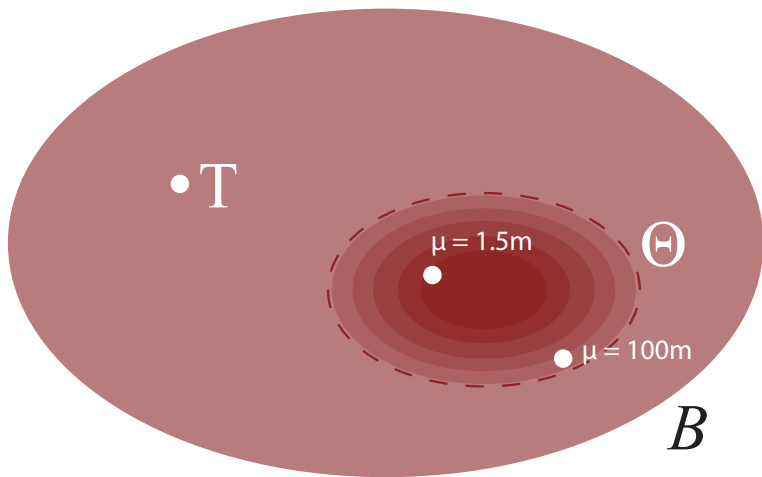
## The prior

1. The Small World is still too big for our purposes.

2. We usually have *some* knowledge about which areas of the Small World are nearest to $T$. For example we don't believe that $\mu = 100m$ and $\mu = 1.5m$ are equally probable.

3. As such, in Bayesian inference we define a *prior* probability density that gives a weighting to all $\theta \in \Theta$ reflecting our beliefs.

4. Frequentist inference does not require us to specify a prior (this causes issues later on that we will discuss).

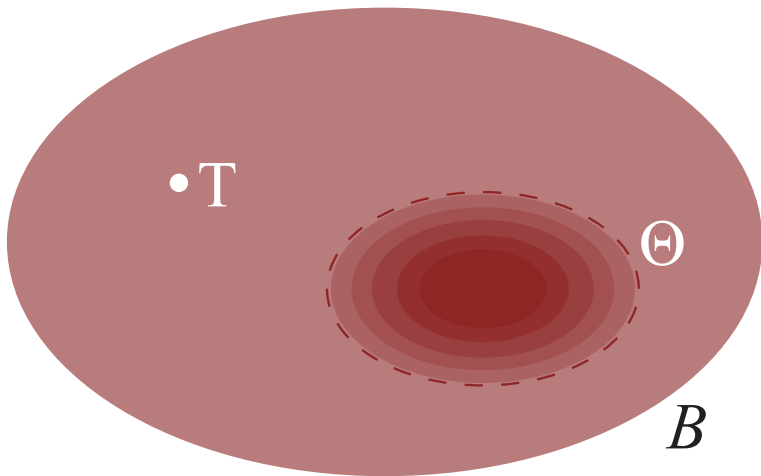1. Inference is the process of updating our prior knowledge in light of data.

1. Inference is the process of updating our prior knowledge in light of data.

2. In Bayesian inference with a likelihood and our prior knowledge explicitly stated we use Bayes' rule to find our posterior probability density over $\theta \in \Theta$.
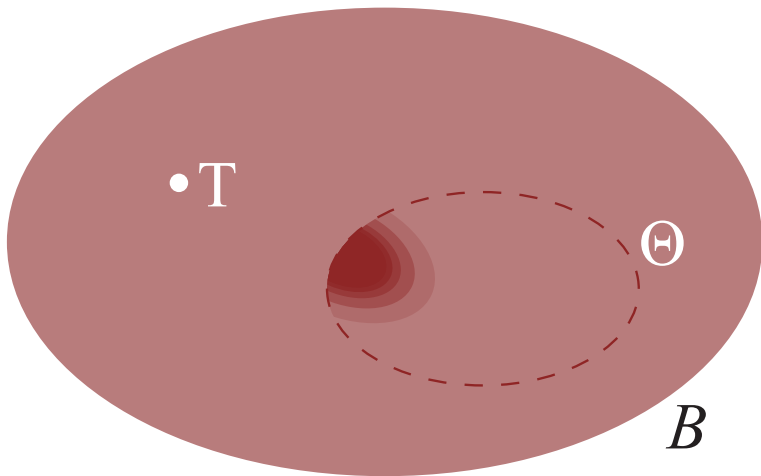
## The data

1. Inference is the process of updating our prior knowledge in light of data.
2. In Bayesian inference with a likelihood and our prior knowledge explicitly stated we use Bayes' rule to find our posterior probability density over $\theta \in \Theta$.
3. The lack of a prior means that in Frequentist inference we generate posterior weightings approximately using rules of thumb (more on this in a minute).

What is the whole (Bayesian) inference process?

Define the observables: The Big World

Define the observables: The Big World

Specify a likelihood

Specify a prior

Input the data

# Example likelihood: frequency of lift malfunctioning[1]

[1]Inspired by Prof. Philip Maini.

# Example likelihood: frequency of lift malfunctioning[1]

- Imagine we want to create a model for the frequency a lift (elevator) breaks down in a given year, $X$.

---

[1]Inspired by Prof. Philip Maini.

## Example likelihood: frequency of lift malfunctioning[1]

- Imagine we want to create a model for the frequency a lift (elevator) breaks down in a given year, $X$.
- This model will be used to plan expenditure on lift repairs over the following few years.

---

[1]Inspired by Prof. Philip Maini.

# Example likelihood: frequency of lift malfunctioning[1]

- Imagine we want to create a model for the frequency a lift (elevator) breaks down in a given year, $X$.
- This model will be used to plan expenditure on lift repairs over the following few years.
- An aside: how to survive a falling lift

---

[1]Inspired by Prof. Philip Maini.

# Example likelihood: frequency of lift malfunctioning[1]

- Imagine we want to create a model for the frequency a lift (elevator) breaks down in a given year, $X$.
- This model will be used to plan expenditure on lift repairs over the following few years.
- An aside: how to survive a falling lift



Figure: Taken from www.npr.org

---

[1]Inspired by Prof. Philip Maini.

- Assume a range of unpredictable and uncorrelated factors (temperature, lift usage, etc.) affect the functioning of the lift.

- Assume a range of unpredictable and uncorrelated factors (temperature, lift usage, etc.) affect the functioning of the lift.
- $\implies X \sim Poisson(\theta)$, where $\theta$ is the mean number of times the lift breaks in one year.

- Assume a range of unpredictable and uncorrelated factors (temperature, lift usage, etc.) affect the functioning of the lift.
- $\implies X \sim Poisson(\theta)$, where $\theta$ is the mean number of times the lift breaks in one year.
- By specifying that $X$ is Poisson-distributed we define the boundaries of the Small World.

- Assume a range of unpredictable and uncorrelated factors (temperature, lift usage, etc.) affect the functioning of the lift.
- $\implies X \sim Poisson(\theta)$, where $\theta$ is the mean number of times the lift breaks in one year.
- By specifying that $X$ is Poisson-distributed we define the boundaries of the Small World.
- **Important:** we don't *a priori* know the *true* value of $\theta$

# Example likelihood: frequency of lift malfunctioning

- Assume a range of unpredictable and uncorrelated factors (temperature, lift usage, etc.) affect the functioning of the lift.
- $\implies X \sim Poisson(\theta)$, where $\theta$ is the mean number of times the lift breaks in one year.
- By specifying that $X$ is Poisson-distributed we define the boundaries of the Small World.
- **Important:** we don't *a priori* know the *true* value of $\theta$ $\implies$ our model defines collection of probability models; one for each value of $\theta$.
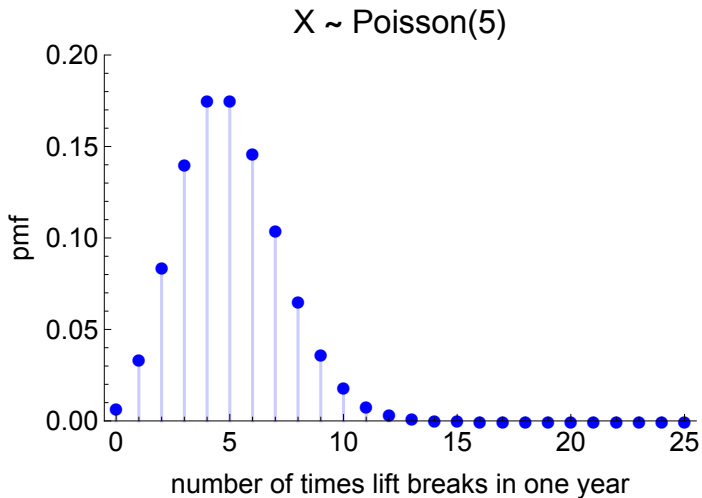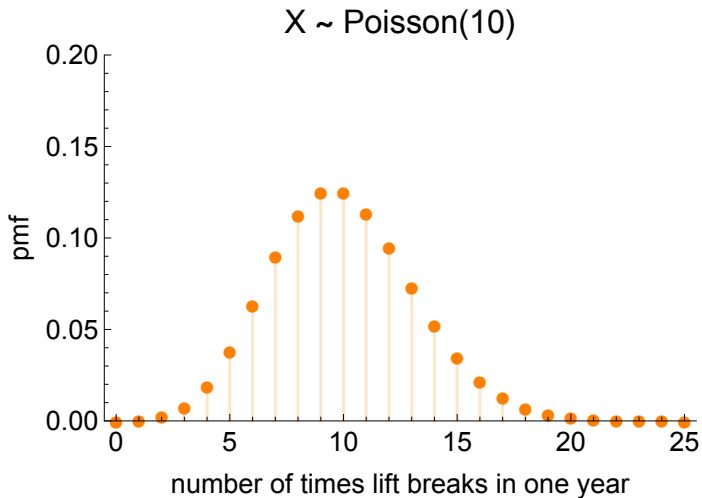
# Example likelihood: frequency of lift malfunctioning

- Assume a range of unpredictable and uncorrelated factors (temperature, lift usage, etc.) affect the functioning of the lift.
- $\implies X \sim Poisson(\theta)$, where $\theta$ is the mean number of times the lift breaks in one year.
- By specifying that $X$ is Poisson-distributed we define the boundaries of the Small World.
- **Important:** we don't *a priori* know the *true* value of $\theta$ $\implies$ our model defines collection of probability models; one for each value of $\theta$.
- We call this collection of models the *Likelihood*.

X ~ Poisson(5)

X ~ Poisson(10)

number of times lift breaks in one year

X ~ Poisson(15)

pmf

number of times lift breaks in one year

X ~ Poisson($\theta$)

pmf

number of times lift breaks in one year

$\theta = 5$
$\theta = 10$
$\theta = 15$

In summary:

In summary:

- By specifying a model framework $X \sim Poisson(\theta)$ we defined the boundaries of the "Small World".

In summary:

- By specifying a model framework $X \sim Poisson(\theta)$ we defined the boundaries of the "Small World".
- The Small World contains a collection of probability distributions known as the *Likelihood*.

- Assume we find that the lift broke down 8 times in the past year.

- Assume we find that the lift broke down 8 times in the past year.
- Our likelihood gives us an *infinite* number of possible ways in which this could have come about.

- Assume we find that the lift broke down 8 times in the past year.
- Our likelihood gives us an *infinite* number of possible ways in which this could have come about.
- Each of these ways corresponds to a unique value of $\theta$.

# The aim of inference: inverting the likelihood

X ~ Poisson(5)  X ~ Poisson(10)  X ~ Poisson(15)

X = 8

- We know that any of these models, each corresponding to different values of $\theta$, could generate the data.

- We know that any of these models, each corresponding to different values of $\theta$, could generate the data.
- In inference we want to use our prior knowledge and data to help us choose which of these models make most sense.

- We know that any of these models, each corresponding to different values of $\theta$, could generate the data.
- In inference we want to use our prior knowledge and data to help us choose which of these models make most sense.
- Essentially we want to run the process in reverse.

# The aim of inference: inverting the likelihood

Start with data



$X = 8$

# The aim of inference: inverting the likelihood

Infer the data generating process



$X = 8$

# The aim of inference: inverting the likelihood

- Both Frequentists and Bayesians essentially invert:
  $p(X|\theta) \rightarrow p(\theta|X)$.

# The aim of inference: inverting the likelihood

- Both Frequentists and Bayesians essentially invert:
  $p(X|\theta) \rightarrow p(\theta|X)$.
- This amounts to going from an 'effect' back to a 'cause'.

# The aim of inference: inverting the likelihood

- Both Frequentists and Bayesians essentially invert:
  $p(X|\theta) \rightarrow p(\theta|X)$.
- This amounts to going from an 'effect' back to a 'cause'.
- Their methods of inversion are *different*.

Frequentist inference considers a single hypothesis $\theta$ about data generating process at a time.

(1)

(2)

Frequentist inference considers a single hypothesis $\theta$ about data generating process at a time.

$$H_0 : \text{A hypothesis } \theta \text{ is true} \tag{1}$$

$$\tag{2}$$

Frequentist inference considers a single hypothesis $\theta$ about data generating process at a time.

$$H_0 : \text{A hypothesis } \theta \text{ is true} \tag{1}$$
$$H_1 : \text{A hypothesis } \theta \text{ is false} \tag{2}$$

Frequentist inference considers a single hypothesis $\theta$ about data generating process at a time.

$$H_0 : \text{A hypothesis } \theta \text{ is true} \tag{1}$$
$$H_1 : \text{A hypothesis } \theta \text{ is false} \tag{2}$$

Frequentists use a rule of thumb:

Frequentist inference considers a single hypothesis $\theta$ about data generating process at a time.

$$H_0 : \text{A hypothesis } \theta \text{ is true} \tag{1}$$

$$H_1 : \text{A hypothesis } \theta \text{ is false} \tag{2}$$

Frequentists use a rule of thumb:

- If $Pr(\text{data as or more extreme than } X | \theta) < 0.05$, then $\theta$ is false, $\implies p(\theta | X) = 0$

Frequentist inference considers a single hypothesis $\theta$ about data generating process at a time.

$$H_0 : \text{A hypothesis } \theta \text{ is true} \tag{1}$$
$$H_1 : \text{A hypothesis } \theta \text{ is false} \tag{2}$$

Frequentists use a rule of thumb:

- If $Pr(\text{data as or more extreme than } X|\theta) < 0.05$, then $\theta$ is false, $\implies p(\theta|X) = 0$
- If $Pr(\text{data as or more extreme than } X|\theta) \geq 0.05$, then $\theta$ could be true, $\implies p(\theta|X) =?$

- For $X = 8$ we can carry out a series of these hypothesis tests across a range of $\theta$.

- For $X = 8$ we can carry out a series of these hypothesis tests across a range of $\theta$.
- For example, assume $\theta = 15$:

# Frequentist inversion: null hypothesis testing

- For $X = 8$ we can carry out a series of these hypothesis tests across a range of $\theta$.
- For example, assume $\theta = 15$:

- For $X = 8$ we can carry out a series of these hypothesis tests across a range of $\theta$.
- For example, assume $\theta = 15$:



$\Pr(X \leq 8 | \theta = 15) \simeq 0.037$

- For $X = 8$ we can carry out a series of these hypothesis tests across a range of $\theta$.
- For example, assume $\theta = 15$:



$Pr(X \leq 8 | \theta = 15) \simeq 0.037 < 0.05 \therefore$ reject !

pmf

number of times lift breaks in one year

- If we carry out a series of similar hypothesis tests over the range of $\theta$ we find the 90% confidence intervals (90% because we have used two one sided 5% test sizes):

- If we carry out a series of similar hypothesis tests over the range of $\theta$ we find the 90% confidence intervals (90% because we have used two one sided 5% test sizes):

$$4.0 \leq \theta \leq 14.4 \tag{3}$$

- If we carry out a series of similar hypothesis tests over the range of $\theta$ we find the 90% confidence intervals (90% because we have used two one sided 5% test sizes):

$$4.0 \leq \theta \leq 14.4 \qquad (3)$$

- **Question:** what does this interval mean?

- If we carry out a series of similar hypothesis tests over the range of $\theta$ we find the 90% confidence intervals (90% because we have used two one sided 5% test sizes):

$$4.0 \leq \theta \leq 14.4 \qquad (3)$$

- **Question:** what does this interval mean?
  - **Answer:** invoke fictitious samples; if we collected an infinite number of data samples and for each one constructed a 90% confidence interval, then 90% of these intervals would contain the true value.

- If we carry out a series of similar hypothesis tests over the range of $\theta$ we find the 90% confidence intervals (90% because we have used two one sided 5% test sizes):

$$4.0 \leq \theta \leq 14.4 \qquad (3)$$

- **Question:** what does this interval mean?
  - **Answer:** invoke fictitious samples; if we collected an infinite number of data samples and for each one constructed a 90% confidence interval, then 90% of these intervals would contain the true value.
- **Second question:** what is the logic for assuming that we obtained data "more extreme than X" in the hypothesis test?

# Frequentist inversion: null hypothesis testing

- If we carry out a series of similar hypothesis tests over the range of $\theta$ we find the 90% confidence intervals (90% because we have used two one sided 5% test sizes):

$$4.0 \leq \theta \leq 14.4 \qquad (3)$$

- **Question:** what does this interval mean?
  - **Answer:** invoke fictitious samples; if we collected an infinite number of data samples and for each one constructed a 90% confidence interval, then 90% of these intervals would contain the true value.
- **Second question:** what is the logic for assuming that we obtained data "more extreme than X" in the hypothesis test?
  - **Answer:** invoke fictitious samples; we needed it because otherwise $Pr(X = 8|\theta)$ would be small for all values of $\theta$.

Bayesians instead use a rule consistent with the rules of probability known as *Bayes' rule*:

Bayesians instead use a rule consistent with the rules of probability known as *Bayes' rule*:

$$p(\theta|X) = \frac{p(X|\theta) \times p(\theta)}{p(X)} \tag{4}$$

Bayesians instead use a rule consistent with the rules of probability known as *Bayes' rule*:

$$p(\theta|X) = \frac{p(X|\theta) \times p(\theta)}{p(X)} \qquad (4)$$

Resulting in an accumulation of evidence (not binary decision) across *all* potential hypotheses $\theta$.

# Bayesian inversion

# Bayesian inversion



mean number of times lift breaks in one year ($\theta$)

- Often we are required to give summary intervals for estimated parameters.

- Often we are required to give summary intervals for estimated parameters.
- There are a number of choices here, for example, the central posterior interval or the highest posterior density interval (although here these are the same).

- Often we are required to give summary intervals for estimated parameters.
- There are a number of choices here, for example, the central posterior interval or the highest posterior density interval (although here these are the same).
- These intervals are known as *credible* intervals, in contrast to the *confidence* intervals of Frequentism.

# Bayesian inversion: finding summary intervals

- Often we are required to give summary intervals for estimated parameters.
- There are a number of choices here, for example, the central posterior interval or the highest posterior density interval (although here these are the same).
- These intervals are known as *credible* intervals, in contrast to the *confidence* intervals of Frequentism.
- These are found by finding an interval such that X% of the area under the pdf (probability mass) is contained within it.

number of times lift breaks in one year ($\theta$)

# Bayesian inversion

- $\implies$ find a 90% central posterior interval of $3.6 \leq \theta \leq 12.4$.

# Bayesian inversion

- $\implies$ find a 90% central posterior interval of $3.6 \le \theta \le 12.4$.
- **Question:** what does this interval mean?

## Bayesian inversion

- $\implies$ find a 90% central posterior interval of $3.6 \leq \theta \leq 12.4$.
- **Question:** what does this interval mean?
  - **Answer:** conditional on our prior knowledge and the data we estimate a 90% probability that this interval contains the true value of $\theta$.

# Frequentist and Bayesian perspectives on probability

The two paradigms differ in their definition of probability:

# Frequentist and Bayesian perspectives on probability

The two paradigms differ in their definition of probability:

- **Frequentists** assume that probabilities represent frequencies of an event's occurrence across an infinite number of exact repetitions of a given experiment

# Frequentist and Bayesian perspectives on probability

The two paradigms differ in their definition of probability:

- **Frequentists** assume that probabilities represent frequencies of an event's occurrence across an infinite number of exact repetitions of a given experiment $\implies$ does not make sense to "update" probabilities.

# Frequentist and Bayesian perspectives on probability

The two paradigms differ in their definition of probability:

- **Frequentists** assume that probabilities represent frequencies of an event's occurrence across an infinite number of exact repetitions of a given experiment $\implies$ does not make sense to "update" probabilities.
- **Bayesians** instead assume that probabilities measure the strength of our underlying beliefs.

The two paradigms differ in their definition of probability:

- **Frequentists** assume that probabilities represent frequencies of an event's occurrence across an infinite number of exact repetitions of a given experiment $\implies$ does not make sense to "update" probabilities.

- **Bayesians** instead assume that probabilities measure the strength of our underlying beliefs. $\implies$ free to update our beliefs using Bayes' rule!

# Frequentist and Bayesian perspectives on probability

## Frequentist

| | | | | | Probability |
|---|---|---|---|---|---|
| Heads | yes | no | ... | no | 0.49 |
| Tails | no | yes | ... | yes | 0.51 |

# Frequentist and Bayesian perspectives on probability

## Frequentist

|        |     |     |       |     | Probability |
|--------|-----|-----|-------|-----|-------------|
| Heads  | yes | no  | …     | no  | 0.49        |
| Tails  | no  | yes | …     | yes | 0.51        |

## Bayesian

Impossible      Probability      Certain

0      1

"The next US President will be the queen"

"The next US President will build a wall"

"The next US President will be American"

All methods of inference attempt to invert the likelihood to make it a valid probability distribution.

All methods of inference attempt to invert the likelihood to make it a valid probability distribution.

**Frequentists:**

- Use a heuristic to do this: if the probability of obtaining data as or more extreme than the actual observation is low when conditioned on $\theta$, then we reject $\theta$.

# Frequentist versus Bayesians: summary

All methods of inference attempt to invert the likelihood to make it a valid probability distribution.

**Frequentists:**

- Use a heuristic to do this: if the probability of obtaining data as or more extreme than the actual observation is low when conditioned on $\theta$, then we reject $\theta$.
    - $\implies$ need to invoke fictitious "more-extreme" data.

All methods of inference attempt to invert the likelihood to make it a valid probability distribution.

**Frequentists:**

- Use a heuristic to do this: if the probability of obtaining data as or more extreme than the actual observation is low when conditioned on $\theta$, then we reject $\theta$.
  - $\implies$ need to invoke fictitious "more-extreme" data.
- Confidence intervals are constructed by repeating this process over a range of $\theta$.

All methods of inference attempt to invert the likelihood to make it a valid probability distribution.

**Frequentists:**

- Use a heuristic to do this: if the probability of obtaining data as or more extreme than the actual observation is low when conditioned on $\theta$, then we reject $\theta$.
    - $\implies$ need to invoke fictitious "more-extreme" data.
- Confidence intervals are constructed by repeating this process over a range of $\theta$.
    - $\implies$ to interpret these intervals we again need to invoke fictitious samples!

# A problem with inverse reasoning[2]

---
[2]Adapted from Gill 1999.

# A problem with inverse reasoning[2]

Even if we disregard the "fictitious samples" criticism there are still problems.

---
[2]Adapted from Gill 1999.

# A problem with inverse reasoning[2]

Even if we disregard the "fictitious samples" criticism there are still problems. Consider the following statement:

---

[2]Adapted from Gill 1999.

# A problem with inverse reasoning[2]

Even if we disregard the "fictitious samples" criticism there are still problems. Consider the following statement:

"Very few Americans are Senators."

[2] Adapted from Gill 1999.

# A problem with inverse reasoning[2]

Even if we disregard the "fictitious samples" criticism there are still problems. Consider the following statement:

"Very few Americans are Senators."

$\implies Pr(S|A)$ is small.

---

[2]Adapted from Gill 1999.

# A problem with inverse reasoning[2]

Even if we disregard the "fictitious samples" criticism there are still problems. Consider the following statement:

"Very few Americans are Senators."

$\implies Pr(S|A)$ is small.
We then suppose that an event occurs:

---

[2]Adapted from Gill 1999.

# A problem with inverse reasoning[2]

Even if we disregard the "fictitious samples" criticism there are still problems. Consider the following statement:

"Very few Americans are Senators."

$\implies Pr(S|A)$ is small.
We then suppose that an event occurs:

We meet a Senator.

---

[2]Adapted from Gill 1999.

# A problem with inverse reasoning[2]

Even if we disregard the "fictitious samples" criticism there are still problems. Consider the following statement:

"Very few Americans are Senators."

$\implies Pr(S|A)$ is small.
We then suppose that an event occurs:

We meet a Senator.

Since $Pr(S|A)$ is small we conclude that $Pr(A|S)$ is too.
Therefore we make the statement:

---

[2]Adapted from Gill 1999.

# A problem with inverse reasoning[2]

Even if we disregard the "fictitious samples" criticism there are still problems. Consider the following statement:

"Very few Americans are Senators."

$\implies Pr(S|A)$ is small.
We then suppose that an event occurs:

We meet a Senator.

Since $Pr(S|A)$ is small we conclude that $Pr(A|S)$ is too.
Therefore we make the statement:

The person is not American.

---

[2]Adapted from Gill 1999.

# Frequentist versus Bayesians: summary

**Bayesians:**

**Bayesians:**

- Use Bayes' law for inversion, which requires we specify a prior distribution.

**Bayesians:**

- Use Bayes' law for inversion, which requires we specify a prior distribution.
- Credible intervals can be constructed by finding the relevant area under the posterior curve.

Intervals we found are similar:

- Frequentist: $4.0 \leq \theta \leq 14.4$

Intervals we found are similar:

- Frequentist: $4.0 \leq \theta \leq 14.4$
- Bayesian: $3.5 \leq \theta \leq 12.4$

Intervals we found are similar:

- Frequentist: $4.0 \leq \theta \leq 14.4$
- Bayesian: $3.5 \leq \theta \leq 12.4$
- This is often the case for the two approaches.

Different views on probability:

Different views on probability:

- **Frequentists** view probabilities as frequencies.

# Frequentist versus Bayesians: summary

Different views on probability:

- **Frequentists** view probabilities as frequencies.
- **Bayesians** view probabilities as subjective measures of uncertainty.

# The dependence of Bayesian inference on probability distributions

Bayesian inference quantifies uncertainty through *probability distributions*.

Bayesian inference quantifies uncertainty through *probability distributions*. Central to Bayesian inference is Bayes' rule:

Bayesian inference quantifies uncertainty through *probability distributions*. Central to Bayesian inference is Bayes' rule:

$$p(\theta|X) = \frac{p(X|\theta) \times p(\theta)}{p(X)} \tag{5}$$

Bayesian inference quantifies uncertainty through *probability distributions*. Central to Bayesian inference is Bayes' rule:

$$p(\theta|X) = \frac{p(X|\theta) \times p(\theta)}{p(X)} \tag{5}$$

$\implies$ we need to be *very* comfortable with probability distributions, to avoid freaking out!

Discrete case:

- Only defined at integer values of $Y$.

## What is a probability distribution?

Discrete case:

- Only defined at integer values of $Y$.
- The value of the distribution at a particular value of $Y$ corresponds to a probability, or probability mass; accordingly we call this distribution a *probability mass function* or pmf.

## What is a probability distribution?

Discrete case:

- Only defined at integer values of $Y$.
- The value of the distribution at a particular value of $Y$ corresponds to a probability, or probability mass; accordingly we call this distribution a *probability mass function* or pmf.
- Probabilities must be non-negative.

## What is a probability distribution?

Discrete case:

- Only defined at integer values of $Y$.
- The value of the distribution at a particular value of $Y$ corresponds to a probability, or probability mass; accordingly we call this distribution a *probability mass function* or pmf.
- Probabilities must be non-negative.
- Sum of probabilities across all allowed values of $Y$ is 1.

# Example discrete distribution: coin flips



Flip a coin ten times and record the number of heads, $Y$

# Example discrete distribution: coin flips



Flip a coin ten times and record the number of heads, $Y$

$$\implies Y \sim Binomial(10, \theta) \tag{6}$$

Flip a coin ten times and record the number of heads, $Y$

$$\implies Y \sim Binomial(10, \theta) \qquad (6)$$

where $\theta$ is the probability of obtaining "heads" on one throw.

# Example discrete distribution: coin flips

**Question:** how do we calculate the probability of any $Y$?

**Question:** how do we calculate the probability of any $Y$?
**Answer:** use the Binomial distribution where we **hold the parameter constant**.

**Question:** how do we calculate the probability of any $Y$?
**Answer:** use the Binomial distribution where we **hold the parameter constant**. For example, if the coin is fair $\theta = \frac{1}{2}$ and the probability is given by,

**Question:** how do we calculate the probability of any $Y$?
**Answer:** use the Binomial distribution where we **hold the parameter constant**. For example, if the coin is fair $\theta = \frac{1}{2}$ and the probability is given by,

$$Pr(Y = y | \theta = \frac{1}{2}) = \binom{10}{5} \left(\frac{1}{2}\right)^y \left(1 - \frac{1}{2}\right)^{10-y} \qquad (7)$$

## Example discrete distribution: coin flips

**Question:** how do we calculate the probability of any $Y$?
**Answer:** use the Binomial distribution where we **hold the parameter constant**. For example, if the coin is fair $\theta = \frac{1}{2}$ and the probability is given by,

$$Pr(Y = y | \theta = \frac{1}{2}) = \binom{10}{5} \left(\frac{1}{2}\right)^y \left(1 - \frac{1}{2}\right)^{10-y} \qquad (7)$$

where $\binom{10}{5}$ is the number of ways of obtaining $5/10$ heads.

# Example discrete distribution: coin flips

**Question:** what does the graph of this (probability) distribution look like?
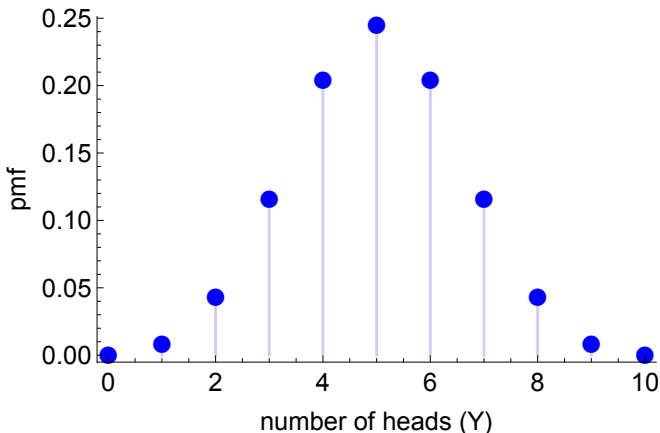
## Example discrete distribution: coin flips

**Question:** what does the graph of this (probability) distribution look like?
**Answer:** it is discrete!

# Example discrete distribution: coin flips

**Question:** what does the graph of this (probability) distribution look like?
**Answer:** it is discrete!

# Example discrete distribution: coin flips

# Example discrete distribution: coin flips

**Question**: how do we calculate the **likelihood**?

# Example discrete distribution: coin flips

**Question**: how do we calculate the **likelihood**?
**Answer:** use the Binomial distribution where we **hold the data constant**.

**Question**: how do we calculate the **likelihood**?
**Answer:** use the Binomial distribution where we **hold the data constant**. For example, $Y = 5$:

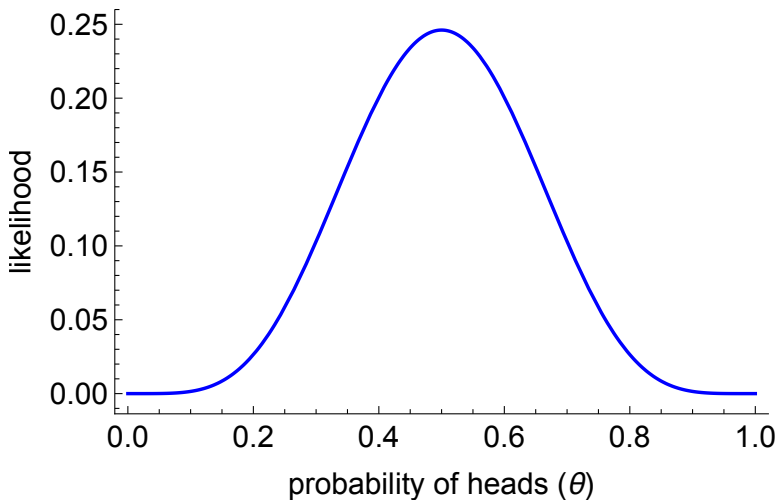# Example discrete distribution: coin flips

**Question**: how do we calculate the **likelihood**?
**Answer:** use the Binomial distribution where we **hold the data constant**. For example, $Y = 5$:

$$\mathcal{L}(\theta|Y = 5) = Pr(Y = 5|\theta)$$

**Question**: how do we calculate the **likelihood**?
**Answer:** use the Binomial distribution where we **hold the data constant**. For example, $Y = 5$:
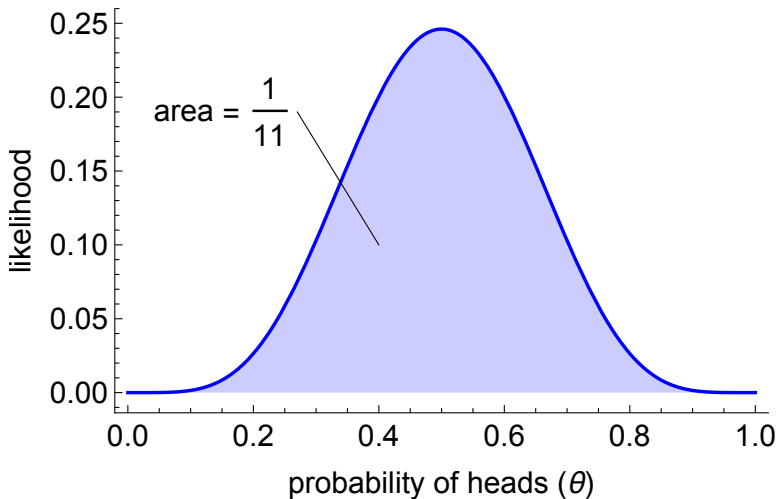
$$\mathcal{L}(\theta | Y = 5) = Pr(Y = 5 | \theta)$$
$$= \binom{10}{5} \theta^5 (1 - \theta)^5$$

## Example discrete distribution: coin flips

**Question**: how do we calculate the **likelihood**?
**Answer:** use the Binomial distribution where we **hold the data constant**. For example, $Y = 5$:

$$\mathcal{L}(\theta | Y = 5) = Pr(Y = 5 | \theta)$$
$$= \binom{10}{5} \theta^5 (1 - \theta)^5$$

**Note:** this is a continuous function of $\theta$, unlike the probability distribution! (Which is a discrete distribution of $Y$.)

# Example discrete distribution: coin flips

area = $\dfrac{1}{11}$

# Continuous distributions

Imagine a continuous variable $W$, for example the wingspan of a seagull.
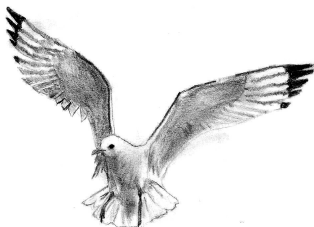
## Continuous distributions

Imagine a continuous variable $W$, for example the wingspan of a seagull.



At infinite precision we wouldn't bet on any one wingspan, for example, 1m or 1.000000001m

# Continuous distributions

Imagine a continuous variable $W$, for example the wingspan of a seagull.



At infinite precision we wouldn't bet on any one wingspan, for example, 1m or 1.000000001m $\implies$ The probability for **any** one value is zero.

# Example continuous distribution: seagull wingspan

Instead of specifying probabilities for individual values we
specify probabilities for **intervals**.

Instead of specifying probabilities for individual values we specify probabilities for **intervals**. For a narrow interval we obtain probability **mass** by multiplying a **probability density** by the **volume** of that interval:

Instead of specifying probabilities for individual values we specify probabilities for **intervals**. For a narrow interval we obtain probability **mass** by multiplying a **probability density** by the **volume** of that interval:

$$\text{probability mass} = \text{probability density} \times \text{volume} \qquad (8)$$

Instead of specifying probabilities for individual values we specify probabilities for **intervals**. For a narrow interval we obtain probability **mass** by multiplying a **probability density** by the **volume** of that interval:

$$\text{probability mass} = \text{probability density} \times \text{volume} \qquad (8)$$

We call the graph of probability densities $p(W)$ a **probability density function** or pdf.

Instead of specifying probabilities for individual values we specify probabilities for **intervals**. For a narrow interval we obtain probability **mass** by multiplying a **probability density** by the **volume** of that interval:

$$\text{probability mass} = \text{probability density} \times \text{volume} \qquad (8)$$

We call the graph of probability densities $p(W)$ a **probability density function** or pdf. To ensure $p(W)$ is a **valid** probability distribution we assume,

Instead of specifying probabilities for individual values we specify probabilities for **intervals**. For a narrow interval we obtain probability **mass** by multiplying a **probability density** by the **volume** of that interval:

$$\text{probability mass} = \text{probability density} \times \text{volume} \qquad (8)$$

We call the graph of probability densities $p(W)$ a **probability density function** or pdf. To ensure $p(W)$ is a **valid** probability distribution we assume,

- $p(W) \geq 0$.

Instead of specifying probabilities for individual values we specify probabilities for **intervals**. For a narrow interval we obtain probability **mass** by multiplying a **probability density** by the **volume** of that interval:

$$\text{probability mass} = \text{probability density} \times \text{volume} \qquad (8)$$

We call the graph of probability densities $p(W)$ a **probability density function** or pdf. To ensure $p(W)$ is a **valid** probability distribution we assume,

- $p(W) \geq 0$.
- The total area under the graph is 1,
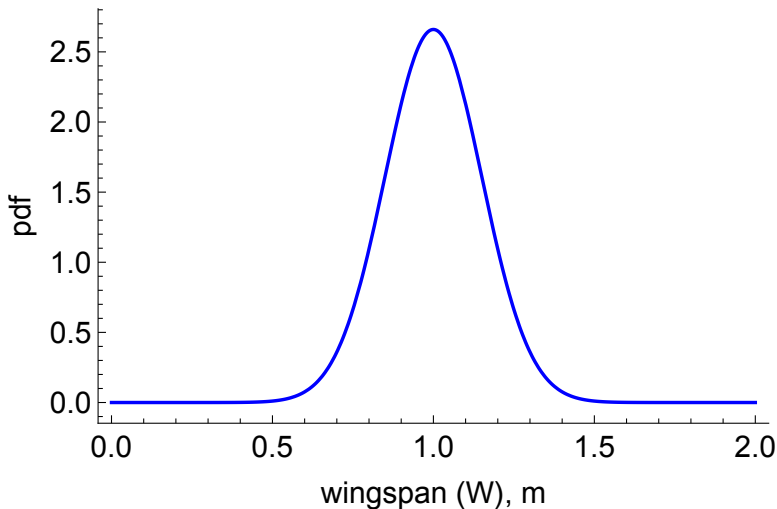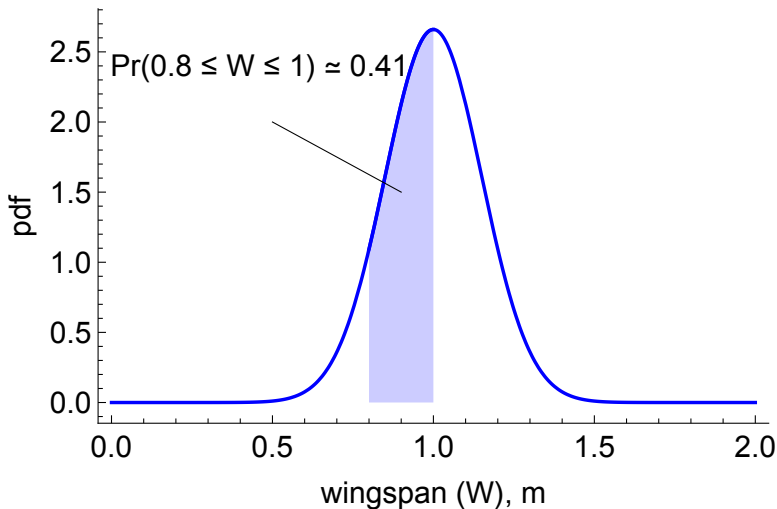
$$Pr(0 \leq W \leq \infty) = 1 \qquad (9)$$

Based on past experience with seagulls we assume
$W \sim N(1, 0.15)$; so they have a mean wingspan of $1m$ and a
standard deviation of $0.15m$.

# Example continuous distribution: seagull wingspan

Based on past experience with seagulls we assume
$W \sim N(1, 0.15)$; so they have a mean wingspan of $1m$ and a
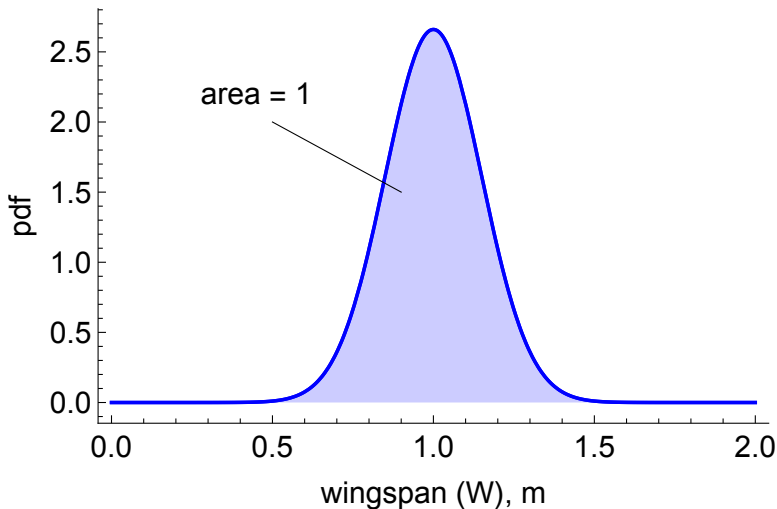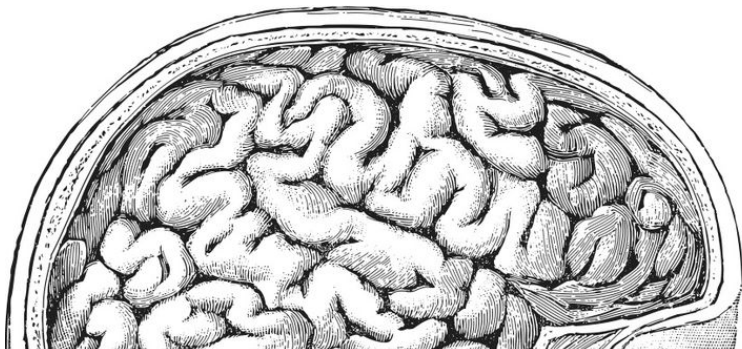standard deviation of $0.15m$.

# Two-dimensional probability distributions

- Imagine you are interested in the interrelation between the circumference of a person's head ($H$) and the volume of their brain ($B$).
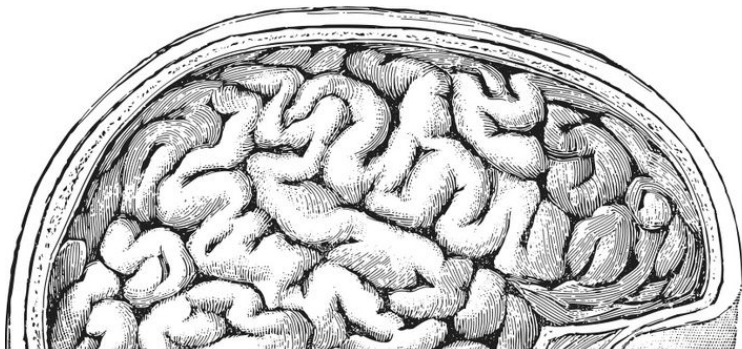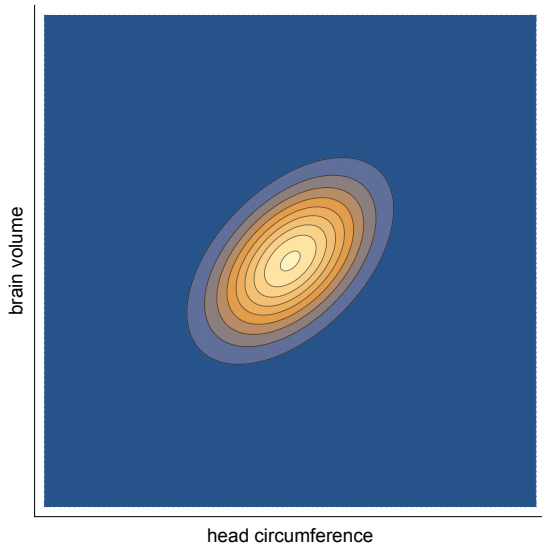
# Two-dimensional probability distributions

- Imagine you are interested in the interrelation between the circumference of a person's head ($H$) and the volume of their brain ($B$).
- Based on data we find there is a positive correlation between these two variables, which we represent in a distribution $p(H, B)$.

brain volume

head circumference

- **Question:** what does the distribution of head circumference look like **irrespective** of brain volume?

- **Question:** what does the distribution of head circumference look like **irrespective** of brain volume?
- **Answer:** sample from the distribution and average-over/remove all possible brain volumes.

# Marginal distributions

- **Question:** what does the distribution of head circumference look like **irrespective** of brain volume?
- **Answer:** sample from the distribution and average-over/remove all possible brain volumes. But what does this mean exactly?

Draw a sample of head circumference and brain volumes from their respective joint distribution:

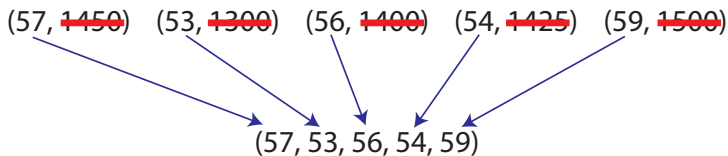(57, 1450)   (53, 1300)   (56, 1400)   (54, 1425)   (59, 1500)

Remove/forget-about the observations of brain volume.

(57, ~~1450~~)  (53, ~~1300~~)  (56, ~~1400~~)  (54, ~~1425~~)  (59, ~~1500~~)

Examine the distribution of the remaining observations.

(57, ~~1450~~)  (53, ~~1300~~)  (56, ~~1400~~)  (54, ~~1425~~)  (59, ~~1500~~)

(57, 53, 56, 54, 59)

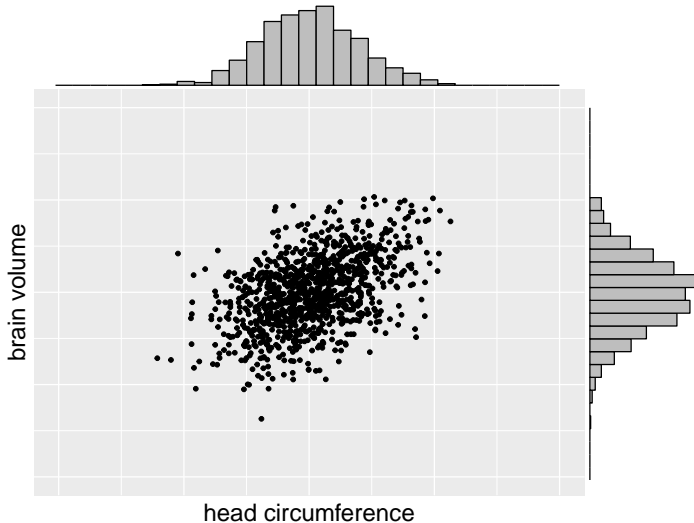# Marginal distribution by sampling

Taking a larger sample:

# Marginal distribution by sampling

Looking at the marginal distributions for each variable:

- The marginal distribution of head circumference is written $p(H)$.

- The marginal distribution of head circumference is written $p(H)$.
- We obtained this distribution by **sampling** from the joint distribution $p(H, B)$.

# Marginal distributions: relationship between sampling and integration

- The marginal distribution of head circumference is written $p(H)$.
- We obtained this distribution by **sampling** from the joint distribution $p(H, B)$.
- Sampling is **approximate** but for large enough samples the approximation is very good.

## Marginal distributions: relationship between sampling and integration

- The marginal distribution of head circumference is written $p(H)$.
- We obtained this distribution by **sampling** from the joint distribution $p(H, B)$.
- Sampling is **approximate** but for large enough samples the approximation is very good.
- The exact way to find the marginals for continuous variables is to use integration, although for most applied problems sampling is far easier.

- The marginal distribution of head circumference is written $p(H)$.
- We obtained this distribution by **sampling** from the joint distribution $p(H, B)$.
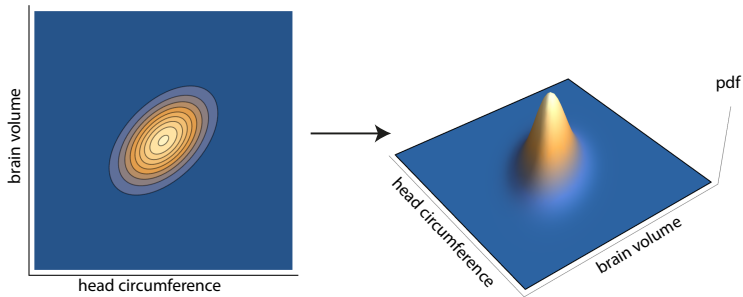- Sampling is **approximate** but for large enough samples the approximation is very good.
- The exact way to find the marginals for continuous variables is to use integration, although for most applied problems sampling is far easier.
- In the next few lectures we will come back to this link between sampling and integration.

# Conditional distributions

- **Question:** If an individual has a brain volume of $1450cm^3$, then what does the distribution for their head circumference look like?

- **Question:** If an individual has a brain volume of $1450cm^3$, then what does the distribution for their head circumference look like?
- **Answer:** Use law of conditional probability:

$$p(H|B = 1450) = \frac{p(B = 1450, H)}{p(B = 1450)} \qquad (10)$$

## Conditional distributions

- **Question:** If an individual has a brain volume of $1450 cm^3$, then what does the distribution for their head circumference look like?

- **Answer:** Use law of conditional probability:

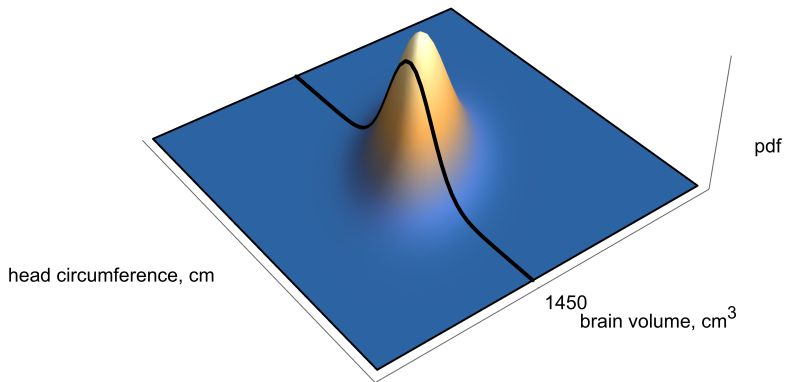$$p(H|B = 1450) = \frac{p(B = 1450, H)}{p(B = 1450)} \qquad (10)$$

- Analogy: imagine walking over the probability distribution along a line of $B = 1450 cm^3$, and recording your height as you go.

# Conditional distributions

## Bayes' rule

- By repeated application of the law of conditional probability we can obtain **Bayes' rule for probabilities**.

## Bayes' rule

- By repeated application of the law of conditional probability we can obtain **Bayes' rule for probabilities**.
- $\implies$ another way to calculate the distribution of head circumferences for an individual with a brain volume of $1450\ cm^3$:

## Bayes' rule

- By repeated application of the law of conditional probability we can obtain **Bayes' rule for probabilities**.
- $\implies$ another way to calculate the distribution of head circumferences for an individual with a brain volume of 1450 $cm^3$:

$$p(H|B = 1450) = \frac{p(B = 1450|H) \times p(H)}{p(B = 1450)} \qquad (11)$$
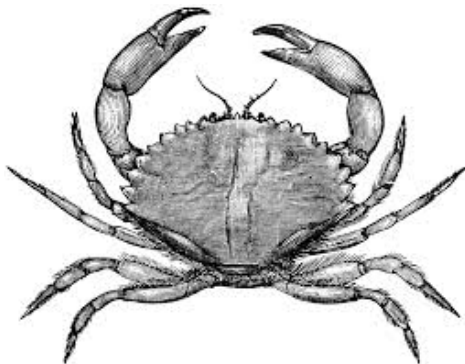
## Bayes' rule

- By repeated application of the law of conditional probability we can obtain **Bayes' rule for probabilities**.
- $\implies$ another way to calculate the distribution of head circumferences for an individual with a brain volume of $1450\ cm^3$:
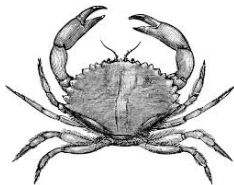
$$p(H|B = 1450) = \frac{p(B = 1450|H) \times p(H)}{p(B = 1450)} \qquad (11)$$

- where $p(H)$ and $p(B)$ are the marginal probability distributions for the head circumferences and brain volumes respectively.

Suppose:

Suppose:

- The probability that a randomly chosen 40 year old woman has breast cancer is approximately $\frac{1}{100}$.

# Bayes' rule in action: breast cancer screening



Suppose:

- The probability that a randomly chosen 40 year old woman has breast cancer is approximately $\frac{1}{100}$.
- If a woman has breast cancer the probability they will test positive in a mammography is about 90%.

Suppose:

- The probability that a randomly chosen 40 year old woman has breast cancer is approximately $\frac{1}{100}$.
- If a woman has breast cancer the probability they will test positive in a mammography is about 90%.
- However there is a risk of about 8% of a false positive result of the test.
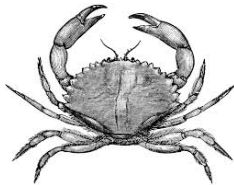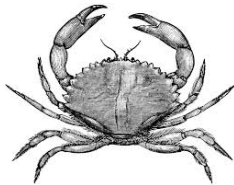
# Bayes' rule in action: breast cancer screening



Suppose:

- The probability that a randomly chosen 40 year old woman has breast cancer is approximately $\frac{1}{100}$.
- If a woman has breast cancer the probability they will test positive in a mammography is about 90%.
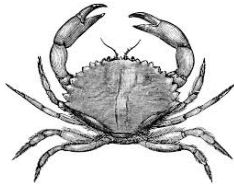- However there is a risk of about 8% of a false positive result of the test.

**Question:** given that a woman tests positive, what is the probability that they have breast cancer?

# Bayes' rule in action: breast cancer screening

**Answer:** we want to find the probability the woman has cancer *given* she has tested positive, which we can do via Bayes' rule (it's the same for pmfs as it was for pdfs):

# Bayes' rule in action: breast cancer screening

**Answer:** we want to find the probability the woman has cancer *given* she has tested positive, which we can do via Bayes' rule (it's the same for pmfs as it was for pdfs):

$$\Pr(\text{🦀} \mid +) = \frac{\Pr(+ \mid \text{🦀}) \times \Pr(\text{🦀})}{\Pr(+)}$$

$$\Pr(\text{🦀} \mid +) \quad = \quad \frac{\overbrace{\Pr(+ \mid \text{🦀})}^{0.9} \times \overbrace{\Pr(\text{🦀})}^{0.01}}{\underbrace{\Pr(+)}_{?}}$$

# Bayes' rule in action: breast cancer screening

$$\Pr(\text{🦀}\,|\,+) \quad = \quad \frac{\overbrace{\Pr(+\,|\,\text{🦀})}^{0.9} \times \overbrace{\Pr(\text{🦀})}^{0.01}}{\underbrace{\Pr(+)}_{?}}$$

- Marginalise out any cancer dependence via summation (discrete equivalent of integration):

# Bayes' rule in action: breast cancer screening

$$\Pr(\text{🦀} \mid +) \quad = \quad \frac{\overbrace{\Pr(+ \mid \text{🦀})}^{0.9} \times \overbrace{\Pr(\text{🦀})}^{0.01}}{\underbrace{\Pr(+)}_{?}}$$

- Marginalise out any cancer dependence via summation (discrete equivalent of integration):

$$\Pr(+) = \underbrace{\Pr(+ \mid \text{🦀}) \times \Pr(\text{🦀})}_{0.9 \quad \times \quad 0.01} + \underbrace{\Pr(+ \mid \text{🕷}) \times \Pr(\text{🕷})}_{0.08 \quad \times \quad 0.99}$$

$$\approx 0.09$$

Putting this into Bayes' rule:

Putting this into Bayes' rule:

$$\Pr(\text{🦀} \mid +) \ = \ \frac{0.9 \quad \times \quad 0.01}{0.09}$$

$$\approx \ 0.1$$

Putting this into Bayes' rule:

$$\Pr(\text{🦀} \mid +) \quad = \quad \frac{0.9 \quad \times \quad 0.01}{0.09}$$

$$\approx \quad 0.1$$

Intuitively, the number of false positives dwarfs the number of true positives.

Take Bayes' rule for probability density of $A$ given $B$:

Take Bayes' rule for probability density of $A$ given $B$:

$$p(A|B) = \frac{p(B|A) \times p(A)}{p(B)} \tag{12}$$

Using a sleight of hand replace: $A \to \theta$ and $B \to X$, where $\theta$ is a parameter vector, and $X$ is a data sample.

# Bayes' rule for inference

Using a sleight of hand replace: $A \to \theta$ and $B \to X$, where $\theta$ is a parameter vector, and $X$ is a data sample.

$$p(\theta|X) = \frac{p(X|\theta) \times p(\theta)}{p(X)} \qquad (13)$$

Using a sleight of hand replace: $A \to \theta$ and $B \to X$, where $\theta$ is a parameter vector, and $X$ is a data sample.

$$p(\theta|X) = \frac{p(X|\theta) \times p(\theta)}{p(X)} \qquad (13)$$

But what do these terms mean? Next lecture.

## Summary

- All methods of inference involve the subjective decision of defining the boundaries of the Small World (likelihood).

## Summary

- All methods of inference involve the subjective decision of defining the boundaries of the Small World (likelihood).
- Small World inference involves inversion of the likelihood.

## Summary

- All methods of inference involve the subjective decision of defining the boundaries of the Small World (likelihood).
- Small World inference involves inversion of the likelihood.
- Frequentists and Bayesians differ in their approach to carry out this inversion:

## Summary

- All methods of inference involve the subjective decision of defining the boundaries of the Small World (likelihood).
- Small World inference involves inversion of the likelihood.
- Frequentists and Bayesians differ in their approach to carry out this inversion:
    - Frequentists use null hypothesis tests.

## Summary

- All methods of inference involve the subjective decision of defining the boundaries of the Small World (likelihood).
- Small World inference involves inversion of the likelihood.
- Frequentists and Bayesians differ in their approach to carry out this inversion:
    - Frequentists use null hypothesis tests.
    - Bayesians use Bayes' rule, which requires us to specify a prior.

# Summary

- All methods of inference involve the subjective decision of defining the boundaries of the Small World (likelihood).
- Small World inference involves inversion of the likelihood.
- Frequentists and Bayesians differ in their approach to carry out this inversion:
    - Frequentists use null hypothesis tests.
    - Bayesians use Bayes' rule, which requires us to specify a prior.
- Bayesian statistics requires us to be able to manipulate probability distributions.

## Light reading

- "New engineering applications of Information Theory", Jaynes 1963.
- "Lifetime earnings and the Vietnam era draft lottery: evidence from social security administration records", Angrist 1990.
- "The insignificance of Null Hypothesis significance testing", Gill 1999.
- "The difference Between 'Significant' and 'Not Significant' is not itself statistically significant", Gelman and Stern 2006.
- "Why most published research findings are false", Ioannidis 2005.
- "Publication and related bias in meta-analysis: power of statistical tests and prevalence in the literature", Sterne, Gavaghan and Egger 2000.

- No problem class this time, but will be next time.

- No problem class this time, but will be next time.
- See you next week on Wednesday at 2pm for "Analytic Bayesian inference".

Thanks!

- David Gavaghan.
- Sam Miles, Francesca Wright.
- Simon Ellis.
- Lab group in Zoology.
- Stan development team.

Bayesian statistics:

$$p(\theta|\boldsymbol{D}) = \frac{p(\boldsymbol{D}|\theta) \times p(\theta)}{p(\boldsymbol{D})} \qquad (14)$$

Beigeian statistics:

$$p(\theta|\boldsymbol{D}) = \frac{p(\boldsymbol{D}|\theta) \times p(\theta)}{p(\boldsymbol{D})} \qquad (15)$$