

# Likelihood

## Bayesian Modeling for Socio-Environmental Data

Chris Che-Castaldo, Mary B. Collins, N. Thompson Hobbs

May 2018



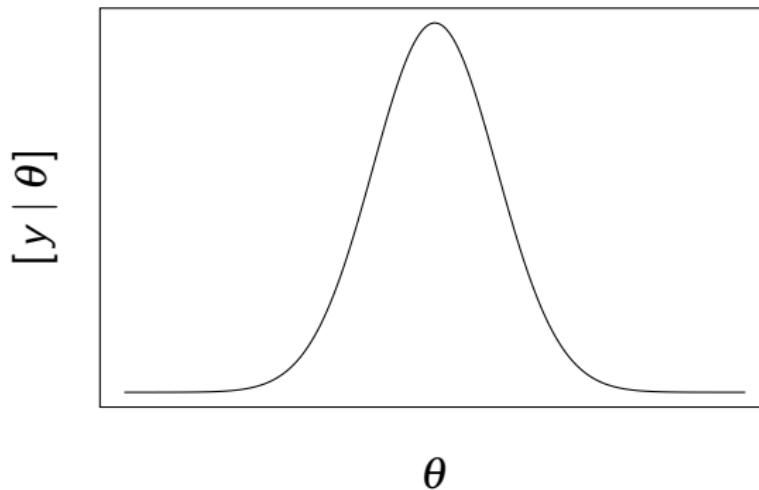
# Why Likelihood?

- Likelihood is a component of all Bayesian models.
- Maximum likelihood is an important statistical approach in its own right.

## Learning objectives for lecture

- Understand the concept of likelihood and its relationship to the probability of data conditional on parameters.
- Describe a likelihood profile and how it differs from a probability density function.
- Understand how to compose likelihoods for multiple parameters and multiple observations.

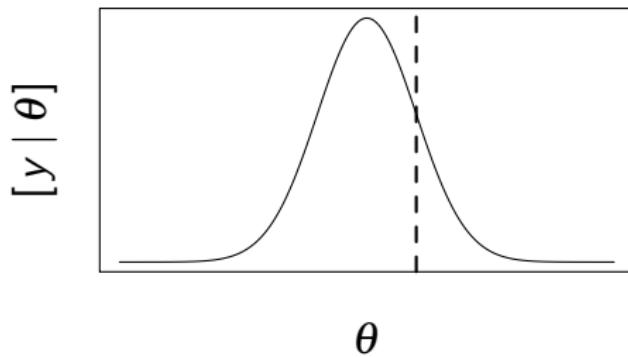
Inference from likelihood is based on  $[y | \theta]$



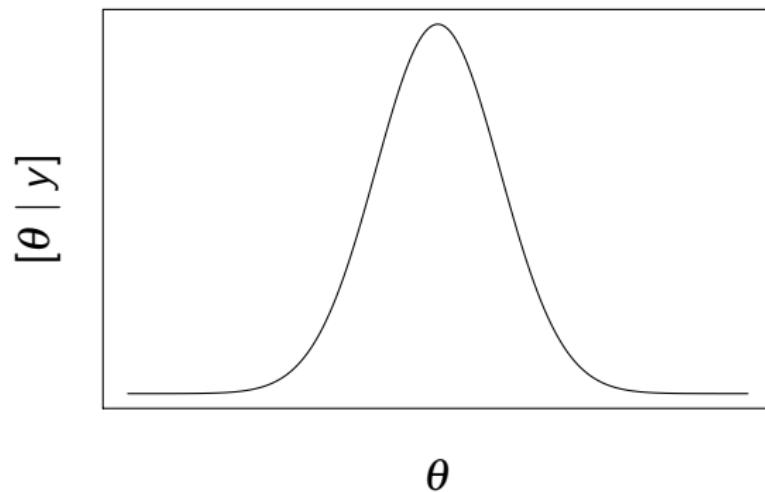
Likelihood allows us to compare alternative parameter values and models by calculating the probability of the data conditional on the parameters  $[y | \theta]$ . As you will see, all evidence based on likelihood is relative.

## Data conditional on a parameter: $[y | \theta]$

Prevalence is a term used in disease ecology to indicate the proportion of a population that is infected. The true prevalence,  $\theta$ , of chronic wasting disease in male mule deer on the winter range north of Fort Collins is 12%. A sample of 24 male deer includes 4 infected individuals. What is the probability of obtaining these data if the estimate of prevalence is true?



Bayesian inference is based on  $[\theta | y]$



## Parameter conditional the data: $[\theta | y]$

We obtain a new sample of 24 male mule deer on the winter range north of Fort Collins that includes 4 infected individuals. In light of these data, what is the probability that the true value of prevalence,  $\theta$ , is found in  $q_L \leq \theta \leq q_U$ ?

## The key idea in likelihood

- In a probability mass or probability density function, the parameter  $\theta$  is constant (known) and the data  $y$  are random variables. The function sums or integrates to 1 over its support.
- In a likelihood function, the data are constant (known) and the parameter is unknown but fixed. We use  $[y | \theta]$  to assess the likelihood of different values of  $\theta$  in light of the data. In this case, the function does not sum or integrate to one over all possible values of the parameter.

$$\underbrace{L(\theta | y)}_{\text{likelihood function}} \propto \underbrace{[y | \theta]}_{\text{PDF or PMF}}$$

Likelihood is *proportional* to probability or probability density

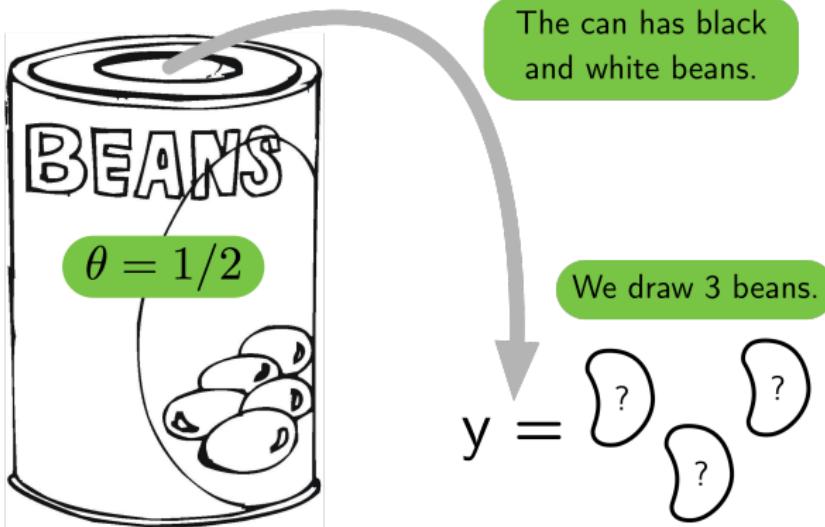
## Discuss notation

$$L(\theta | y) \propto [y | \theta] \quad (1)$$

$$L(\theta | y) = c[y | \theta] \quad (2)$$

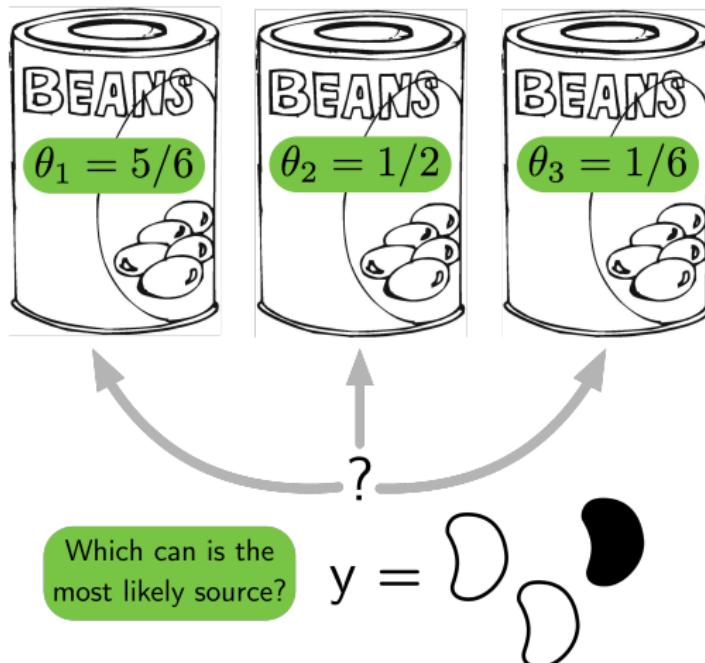
$$L(\theta | y) = [y | \theta] \quad (3)$$

# The parameter is known and the data are random



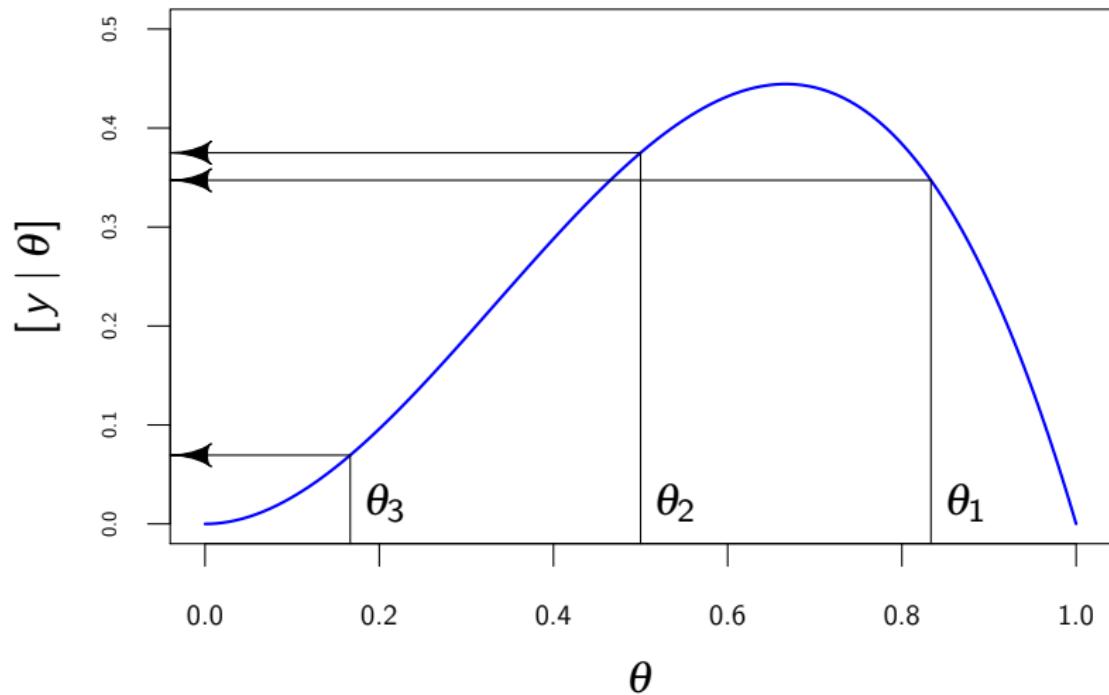
- What are the possible outcomes?
- What probability mass function would you use to model these data?
- What is the probability of each outcome?
- What is the sum of the individual probabilities?

# The data are known and the parameter is random



- What is the likelihood of each parameter value?

## A likelihood profile: 2 white beans on 3 draws

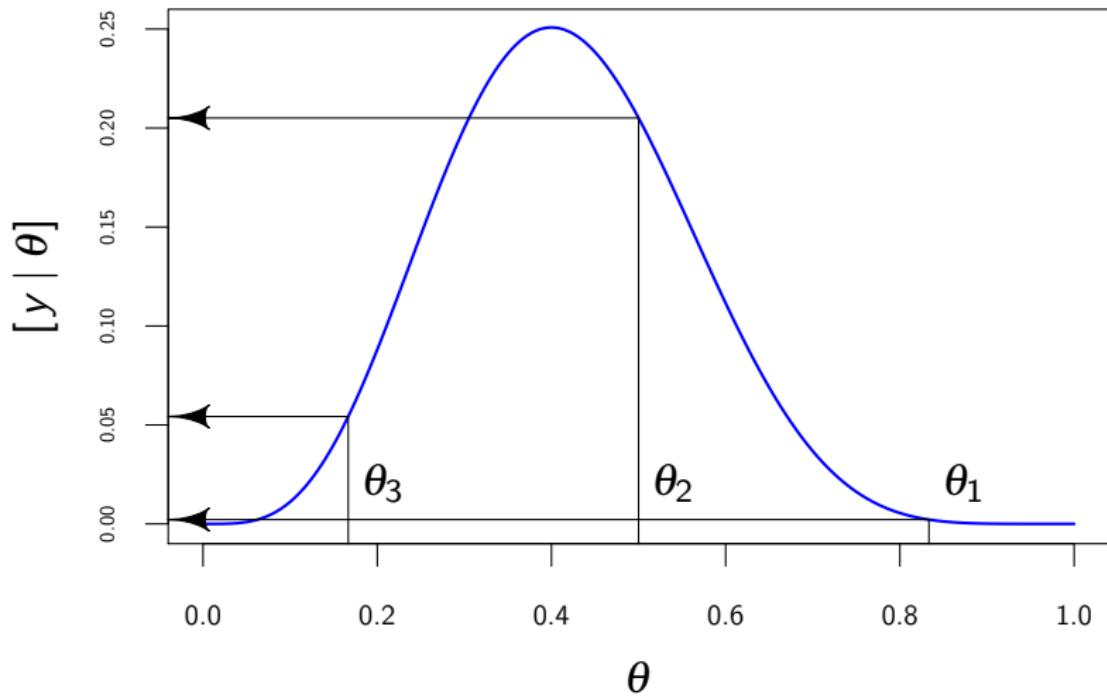


## Class exercise: Can of beans

Have someone take a draw of 10 beans from one of the cans where the identity of the can is unknown.

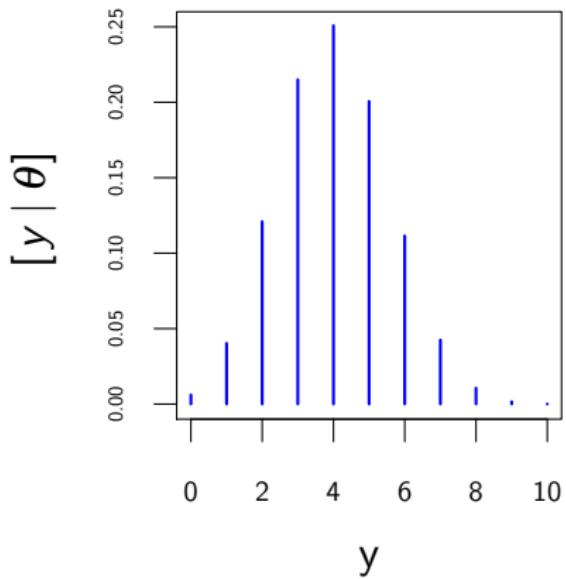
- Which of the three cans is the most likely to have produced this draw?
- How much more likely is this can than the other two?

## A likelihood profile: 4 white beans on 10 draws

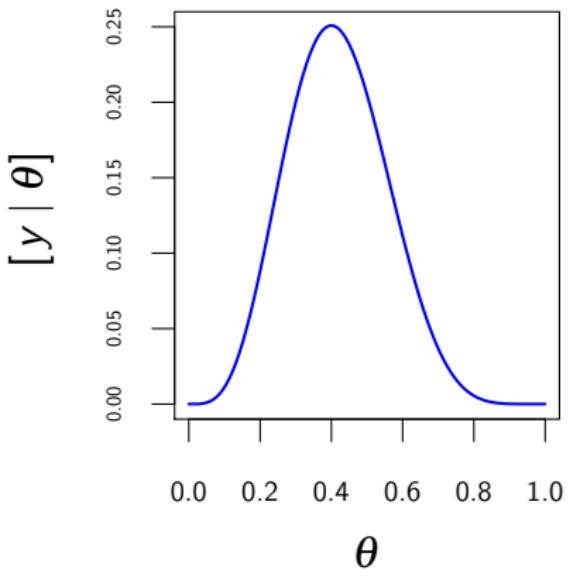


# Likelihood vs Probability:

Probability mass function  
Parameter is fixed



Likelihood profile  
Data are fixed



## How do we fit models with multiple parameters?

In the example we had a single parameter,  $\theta$ , one observation of 4 successes on 10 draws, and a binomial likelihood. However, we could have made the likelihood a function of the *predictions* of a model, and used any probability mass function or probability density function as a “wrapper” for the predictions, i.e.,

$$\begin{aligned}\mu_i &= g(\theta, x_i) \\ \mathcal{L}(\mu_i, \sigma^2 | y_i) &\propto \underbrace{[y_i | \mu_i, \sigma^2]}_{\text{PDF or PMF}}\end{aligned}$$

# Likelihood Surfaces

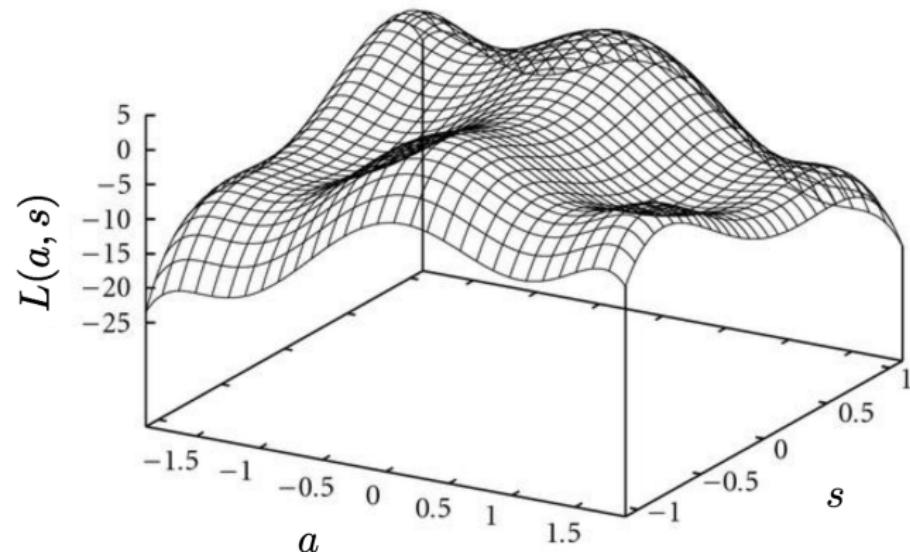


Figure courtesy of [Sergiy Nesterko](#).

## How do we fit models with multiple parameters and multiple data points?

The total likelihood is the product of the individual likelihoods, assuming the data are *conditionally independent*:

$$\mathcal{L}(\boldsymbol{\theta}, \sigma^2 | \mathbf{y}) = c \prod_{i=1}^n [y_i | g(\boldsymbol{\theta}, x_i), \sigma^2]$$

## What does conditionally independent mean?

Independence is an assumption! Remember from the chain rule:

$$Pr(y_1, \dots, y_n) = Pr(y_1 | y_2 \dots y_n) Pr(y_2 | y_3 \dots y_n) \dots Pr(y_n).$$

However, by assuming that these random variables are independent, you can simplify the joint probability into:

$$Pr(y_1, \dots, y_n) = Pr(y_1) Pr(y_2) \dots Pr(y_n),$$

such that the total likelihood a product of the individual likelihoods.

## What does conditionally independent mean?

We evaluate the independence assumption by examining the residuals ( $\varepsilon$ ) from a model, where ( $\varepsilon_i = y_i - g(\theta, x_i)$ ).

The independence assumption holds if knowing a residual tells you nothing about the other residuals.

We assess this by ensuring that the residuals:

- do not show a trend, meaning they should be centered on 0 throughout the range of fitted values,
- and are not autocorrelated.

## Log likelihoods:

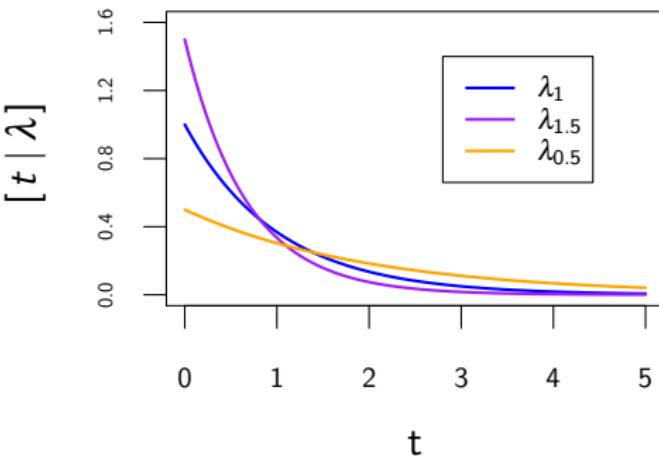
We often use the sum of the log likelihoods to get the total log likelihood as a basis for fitting models:

$$\log(\mathcal{L}(\boldsymbol{\theta}, \sigma^2 | y)) = L(\boldsymbol{\theta}, \sigma^2 | y) = \log(c) + \sum_{i=1}^n \log([y_i | g(\boldsymbol{\theta}, x_i), \sigma^2])$$

# The exponential distribution

$$t_i \sim \text{exponential}(\lambda)$$
$$P(t_i | \lambda) = \begin{cases} \lambda e^{-\lambda t_i} & t_i \geq 0 \\ 0 & t_i < 0 \end{cases}$$

if  $\lambda$  is the average number of events time $^{-1}$ , then  $1/\lambda$  is the average time between events and mean of the exponential distribution.



- The data,  $t_i$ , represent “waiting times” between events happening in a Poisson process. If the number of events per unit time is provided by the Poisson distribution, then the length of time between events is provided by the exponential distribution.
- See this [link](#) by John C.B. Cooper for a clear and concise explanation of the connection between these two distributions.

## Likelihood exercise for time to tweet

Generate a datum quantifying how long we have to wait for a new tweet from the POTUS. Write out the likelihood function assuming wait times for POTUS tweets are governed by an exponential distribution. Determine the maximum likelihood estimate (MLE) for  $\lambda$ .

## Generate the datum to estimate $\lambda$ for POTUS

- $\frac{991}{151} \sim 6.6$  tweets per day
- $t_1 = \frac{1}{6.6} \sim 0.15$  days between each tweet
- Fun fact: 3.7 hours respite between POTUS tweets!

## Analytically solve for $\lambda_{MLE}$ for POTUS

$$\mathcal{L}(\lambda) = [t_1 | \lambda] \quad \text{The likelihood equals the PDF.}$$

$$\mathcal{L}(\lambda) = \lambda e^{-\lambda t_1} \quad \text{Substitute the exponential PDF.}$$

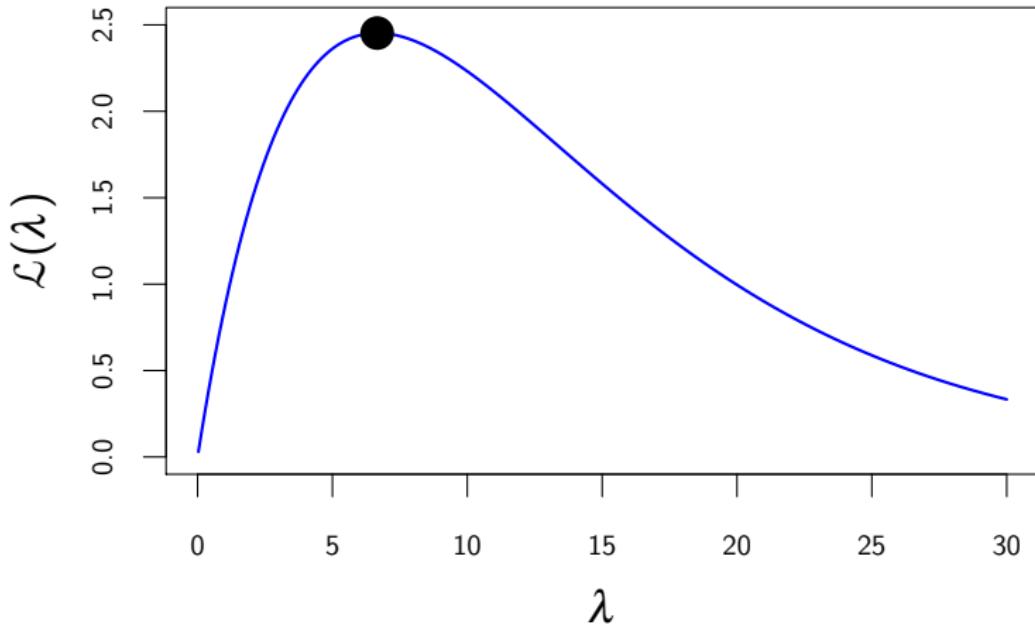
$$L(\lambda) = \log(\lambda) - \lambda t_1 \quad \text{Take logs to ease differentiation.}$$

$$\frac{dL(\lambda)}{d\lambda} = \frac{1}{\lambda} - t_1 \quad \text{Differentiate.}$$

$$0 = \frac{1}{\lambda_{MLE}} - t_1 \quad \text{Set derivative} = 0 \text{ and solve for } \lambda_{MLE}.$$

$$\lambda_{MLE} = \frac{1}{0.15} = 6.6 \quad \text{Note that } \lambda_{MLE} = \text{Poisson rate parameter.}$$

## Likelihood profile of $\lambda$ for POTUS



## Tweeting by heads of state<sup>1</sup>



Assuming tweet wait times for heads of state arise from an exponential distribution whose parameter is  $\lambda$ :

- Write the total likelihood and log-likelihood function.
- Advanced Find the maximum likelihood estimate for  $\lambda$ .
- Advanced Plot the likelihood profile for  $\lambda$  using R.

---

<sup>1</sup>Wait times between tweets collected with R package `twitteR`.

## Main points

- Likelihood allows us to evaluate the relative strength of evidence for one parameter or model relative to another.
- The data are fixed and the parameters are variable in likelihood functions. These functions do not integrate or sum to one over the range of values of the parameter.
- The data are variable and the parameter are fixed in probability mass functions and probability density functions. These functions sum or integrate to one over the support of the random variable,  $y$ .

## Likelihood ratio confidence intervals

Find the upper and lower bounds of an interval where all  $\lambda$  values within that interval are as consistent with the data as  $\lambda_{MLE}$ .

We compute the likelihood ratio statistic:

$$\mathcal{R} = 2 \log \left( \frac{\mathcal{L}(\lambda_{MLE} | t_1)}{\mathcal{L}(\lambda_0 | t_1)} \right) \sim \chi^2_{k=1}$$

which is distributed  $\chi^2$  with 1 degree of freedom. Note that we fail to reject  $H_0$  that  $\lambda = \lambda_0$  if  $\mathcal{R} < \chi^2_{k=1}(1 - \alpha)$ .

## Likelihood ratio confidence intervals

We determine the  $(1 - \alpha = 0.95)$  likelihood ratio confidence interval by finding the upper and lower bounds for all values of  $\lambda_0$  where we would fail to reject  $H_0$ .

$$\begin{aligned} 2 \log \left( \frac{\mathcal{L}(\lambda_{MLE} | t_1)}{\mathcal{L}(\lambda | t_1)} \right) &< \chi^2_{k=1}(0.95) \\ L(\lambda_{MLE}) - \frac{3.84}{2} &< L(\lambda | t_1) \\ L(\lambda | t_1) &> L(\lambda_{MLE}) - 1.92 \end{aligned}$$

## Likelihood ratio confidence intervals

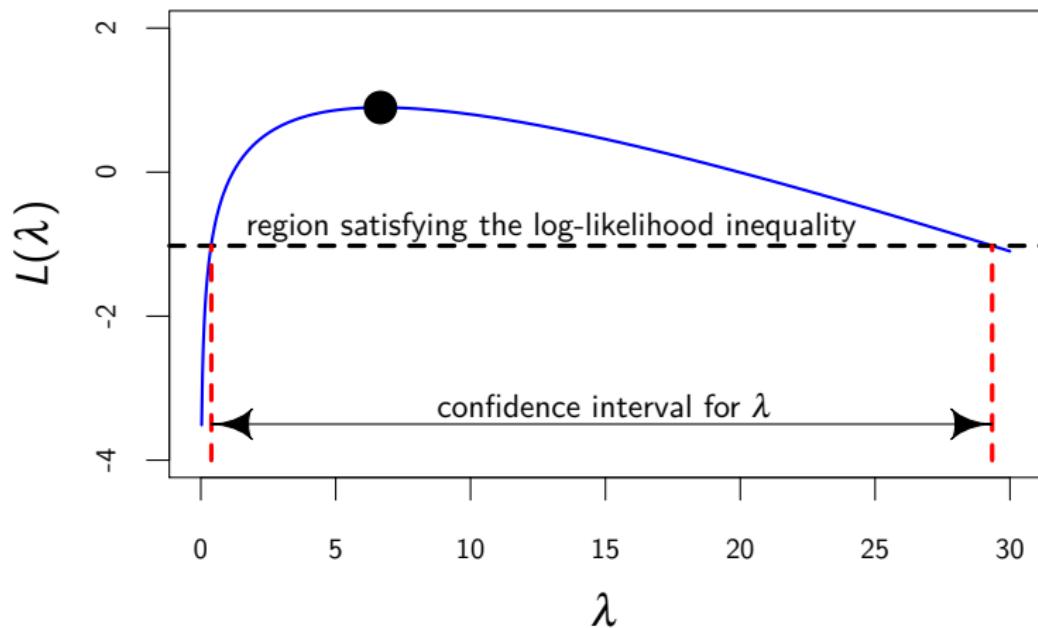
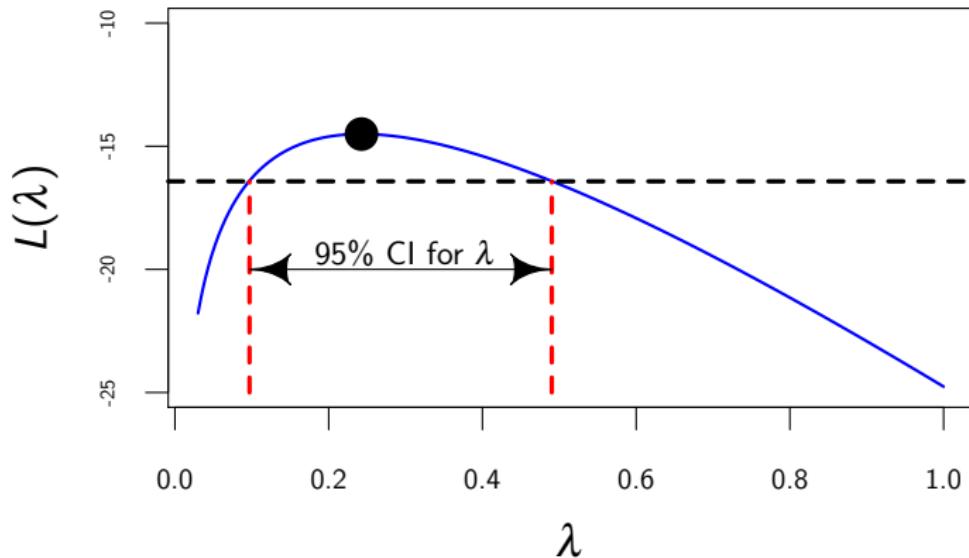


Figure courtesy of the UNC Biology Department.

# Likelihood profile of $\lambda$ for heads of state



```
lambda <- seq(.03, 1, length = 1000)
y <- NA
for(i in 1:length(lambda)) {y[i] <- log(prod(dexp(c(.15, 1.5, 22, .077, .0053, 1.1), lambda[i])))}
```