

GOV 2001/ 1002/ E-2001 Section 4

Maximum Likelihood Estimation

Mayya Komisarchik¹

Harvard University

February 24, 2016

¹Heartfelt thanks to Solé Prillaman, Stephen Pettigrew and all of the other Gov 2001 TFs of yesteryear. These slides draw heavily (but not i.i.d.) from their materials

LOGISTICS

Paper Selection- Comments coming soon!

Reading Assignment-

1. UPM Chapter 5
2. “Making the Most of Statistical Analyses: Improving Interpretation and Presentation” - King, Tomz, Wittenberg 2000
3. Sections 1 & 2 from “Toward a Common Framework for Statistical Analysis and Development” - Imai, King, Lau 2008
4. Recommended: “Election Forecasting: Principles and Practice” - Beck 2005

Problem Set 4- Due by 6pm next Wednesday (on Canvas).

Assessment Question- Due by 6pm next Wednesday on Canvas. You must work alone and only one attempt.

REPLICATION PAPER

Start getting data- journal website, author's website, ICPSR, email author.

Start replicating the paper- Literally like right now.

Replication- due March 23 by 6pm (draft with no text but relevant tables and figures replicated).

OUTLINE

MAXIMUM LIKELIHOOD ESTIMATORS

Steps to finding the MLE:

1. Write out the model.
2. Calculate the likelihood ($L(\theta|y)$) for all observations.
3. Take the log of the likelihood ($\ell(\theta|Y)$).
4. Plug in the systematic component for θ_i .
5. Bring in observed data.
6. Maximize $\ell(\theta|y)$ with respect to θ and confirm that this is a maximum.
7. Find the variance of your estimate.

UNIVARIATE EXAMPLE



Ex. Waiting for the Redline – How long will it take for the next T to get here?

1. WRITE THE MODEL

Y is a Exponential random variable with parameter λ .

$$f(y) = \lambda e^{-\lambda y}$$

Let's assume that Y is distributed Exponentially with some constant rate of seeing a train arrive across observations. We might also assume that observations are independent.

The model:

1. $Y_i \sim f_{\text{expo}}(y_i | \lambda_i)$.
2. $\lambda_i = \lambda$.
3. Y_i and Y_j are independent for all $i \neq j$.

2. CALCULATE $L(\lambda|y)$

$$L(\lambda_i|y_i) \propto f_{\text{expo}}(y_i|\lambda_i)$$

$$\propto \lambda_i e^{-\lambda_i y_i}$$

$$L(\lambda|y) \propto \lambda_1 e^{-\lambda_1 y_1} \cdot \lambda_2 e^{-\lambda_2 y_2} \cdot \dots \cdot \lambda_n e^{-\lambda_n y_n}$$

$$\propto \prod_{i=1}^n \lambda_i e^{-\lambda_i y_i}$$

3. CALCULATE $\ell(\lambda|y)$

$$L(\lambda|y) \propto \prod_{i=1}^n \lambda_i e^{-\lambda_i y_i}$$

$$\ell(\lambda|y) \propto \ln\left(\prod_{i=1}^n \lambda_i e^{-\lambda_i y_i}\right)$$

$$\propto \sum_{i=1}^n \ln(\lambda_i e^{-\lambda_i y_i})$$

$$\propto \sum_{i=1}^n (\ln(\lambda_i) + \ln(e^{-\lambda_i y_i}))$$

$$\propto \sum_{i=1}^n (\ln \lambda_i - \lambda_i y_i)$$

4. PLUG IN SYSTEMATIC COMPONENT

Remember in our model we assumed that the rate was constant across observations (i.e. $\lambda_i = \lambda$).

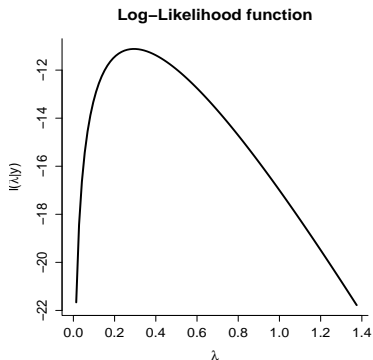
$$\begin{aligned}\ell(\lambda|y) &\propto \sum_{i=1}^n (\ln \lambda_i - \lambda_i y_i) \\ &\propto \sum_{i=1}^n (\ln \lambda - \lambda y_i) \\ &\propto n \ln \lambda - \lambda \sum_{i=1}^n y_i\end{aligned}$$

5. BRING IN OBSERVED DATA

Let's say that we take the train 5 times this week and observe the following times until arrival $Y : \{1, 5, 8, 2, 1\}$.

$$\begin{aligned}\ell(\lambda|y) &\propto n \ln \lambda - \lambda \sum_{i=1}^n y_i \\ &\propto 5 \ln \lambda - \lambda(1 + 5 + 8 + 2 + 1) \\ &\propto 5 \ln \lambda - 17\lambda\end{aligned}$$

6. MAXIMIZE $\ell(\theta|Y)$



Where is the maximum?

How could we find it analytically using

$$\ell(\lambda|y) = 5 \ln \lambda - 17\lambda$$

OPTIMIZATION

Optimization is the process of minimizing or maximizing a function by systematically choosing the values of variables from within an allowed set. For example:

$$\min_{x \in [-\infty, \infty]} f(x) = -\frac{1}{2}(3 - x)^2$$

- $f(x)$ is called the objective function
- x is the parameter (for us λ or β , π , σ , etc.)
- $x \in [-\infty, \infty]$ is the allowed set or the parameter space

Two ways to solve:

1. Analytically
2. Numerically

ANALYTIC OPTIMIZATION

Step One: Take the first derivative of the function with respect to the parameter of interest.

- ▶ The derivative of a function at a value x_0 , denoted by $f'(x_0)$ or $\frac{\partial f}{\partial x}(x_0)$, is the instantaneous rate of change in $f(x)$ at x_0 .
- ▶ Define the **Score** as:

$$S(\theta) = \frac{\partial \ell(\theta | \mathbf{x})}{\partial \theta}$$

ANALYTIC OPTIMIZATION

Step One: Take the first derivative of the function with respect to the parameter of interest.

$$\begin{aligned}\ell(\lambda|y) &= 5 \ln \lambda - 17\lambda \\ S(\lambda) &= \frac{\partial \ell(\lambda|y)}{\partial \lambda} = \frac{\partial (5 \ln \lambda - 17\lambda)}{\partial \lambda} \\ &= \frac{5}{\lambda} - 17\end{aligned}$$

ANALYTIC OPTIMIZATION

Step Two: Set the first derivative of the function equal to 0 and identify the critical value(s) of our parameter.

- ▶ $f'(x_0)$ describes the behavior of a function on an interval $[a,b]$
 - If $f'(x) > 0$ for all $x \in [a, b]$, then f is increasing on the interval $[a, b]$
 - If $f'(x) < 0$ for all $x \in [a, b]$, then f is decreasing on the interval $[a, b]$
 - If $f'(x) = 0$ at some $x \in [a, b]$ then we say x is a critical value of f . Critical values may be maxima, minima, or saddle points.

ANALYTIC OPTIMIZATION

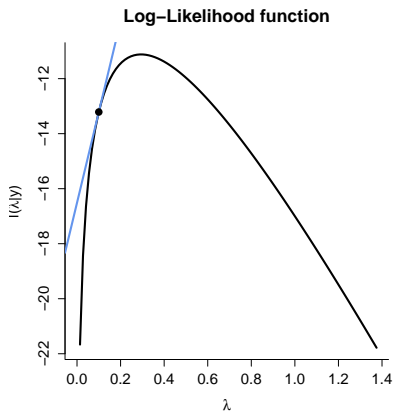
Step Two: Set the first derivative of the function equal to 0 and identify the critical value(s) of our parameter.

$$\begin{aligned}\frac{\partial \ell(\lambda|y)}{\partial \lambda} &= 0 \\ \frac{5}{\lambda} - 17 &= 0 \\ \hat{\lambda} &= \frac{5}{17}\end{aligned}$$

$\hat{\lambda}$ is the maximum likelihood estimate given our observed data!

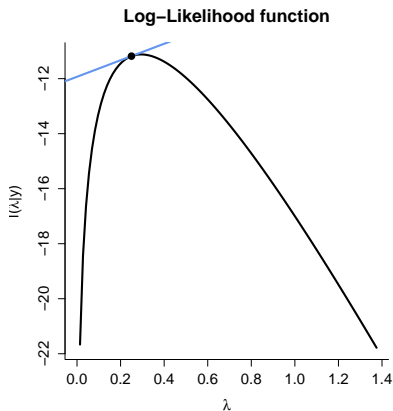
ANALYTIC OPTIMIZATION

Why do we think $\hat{\lambda}$ is the MLE? Put another way, why do we think $\hat{\lambda}$ maximizes the log-likelihood function?



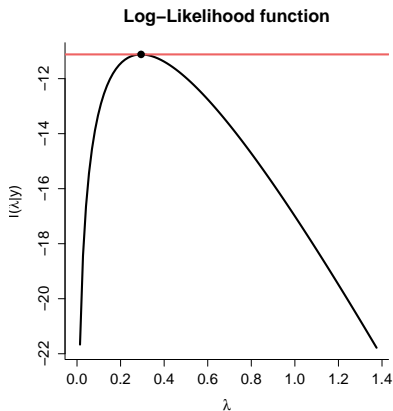
ANALYTIC OPTIMIZATION

Why do we think $\hat{\lambda}$ is the MLE? Put another way, why do we think $\hat{\lambda}$ maximizes the log-likelihood function?



ANALYTIC OPTIMIZATION

Why do we think $\hat{\lambda}$ is the MLE? Put another way, why do we think $\hat{\lambda}$ maximizes the log-likelihood function?



ANALYTIC OPTIMIZATION

Step Three: Compute the second derivative of the function at the critical value(s) and evaluate.

- ▶ The second derivative of a function $f''(x)$ or $\frac{\partial^2 f}{\partial x \partial x}(x)$ is the derivative of the derivative, or the rate of change of the rate of change.
- ▶ Use the following to evaluate your critical value(s):
 - If $f'(x_0) = 0$, and $f''(x_0) < 0$, then x_0 is a maximum
 - If $f'(x_0) = 0$ and $f''(x_0) > 0$, then x_0 is a minimum
 - If $f'(x_0) = 0$ and $f''(x_0) = 0$, then x_0 may be a minimum, a maximum, or neither.
- ▶ Define the **Hessian** for the univariate case as:

$$H(\theta) = \frac{\partial^2 \ell(\theta|\mathbf{x})}{\partial \theta^2}$$

ANALYTIC OPTIMIZATION

Step Three: Compute the second derivative of the function at the critical value(s) and evaluate.

$$\begin{aligned} S(\lambda) &= \frac{\partial \ell(\lambda|y)}{\partial \lambda} = \frac{5}{\lambda} - 17 \\ H(\lambda) &= \frac{\partial^2 \ell(\lambda|y)}{\partial \lambda^2} = -\frac{5}{\lambda^2} \\ \frac{\partial^2 \ell(\lambda|y)}{\partial \lambda^2} \Big|_{\lambda=\hat{\lambda}} &= -\frac{5}{(5/17)^2} \\ &= -\frac{289}{5} \end{aligned}$$

This confirms that $\hat{\lambda}$ is a **maximum** because the second derivative is negative.

NUMERICAL OPTIMIZATION

For problems of only slightly more complexity, using derivatives and solving for parameters in order to maximize may be not only impractical but impossible.

There are a number of ways to do numerical optimization in R, including (listed from most to least DIY:)

1. Write a function for $\ell(\theta|y)$ and program your own iterative optimization algorithm (on your next problem set!)
2. Write a function for $\ell(\theta|y)$ and use R's built in `optim()` function to find your MLE(s) (also on your next problem set!)
3. Use R's `glm()` function in much the same way that you would use `lm()`. This is a great way to check your work.

NUMERICAL OPTIMIZATION

Step One: Write a function in R for your log-likelihood.

```
ll.expo <- function(par, y){  
  length(y)*log(par) - par*sum(y)  
}
```

par is our parameter lambda.

NUMERICAL OPTIMIZATION

Step Two: Optimize this function using `optim()`.

`optim()` takes a starting value (`par`) and a function (`fn`) as its main arguments.

```
# Create our data
y <- c(1,5,8,2,1)

# Optimize
opt.expo <- optim(par = 0.01, fn = ll.expo, y = y,
  method = "BFGS", control = list(fnscale = -1),
  hessian = TRUE)
```

Note that (`par`) takes a *vector* of starting values for *all* the parameters you will be estimating! We only have one in this case, but you'll be working with models that have multiple parameters!

(`fn`) will be your log-likelihood function.

NUMERICAL OPTIMIZATION

Step Two

```
# Create our data
y <- c(1,5,8,2,1)

# Optimize
opt.expo <- optim(par = 0.01, fn = ll.expo, y = y,
  method = "BFGS", control = list(fnscale = -1),
  hessian = TRUE)
```

What are these three extra arguments?

1. `method` is the algorithm used to find the maximum.
2. `fnscale` multiplies the function by the given constant. As a default `optim()` finds the minimum so multiplying our function by -1 fools `optim()` into finding the maximum.
3. `hessian = TRUE` requests that `optim` return a matrix of second derivatives which in the univariate case will be 1×1 .

NUMERICAL OPTIMIZATION

```
## Output from Optim
opt.expo

$par
[1] 0.2941228

$value
[1] -11.11888
...
$hessian
      [,1]
[1,] -57.79931

# If we want to pull out our MLE
mle <- opt.expo$par
# This is the same as 5/17, our analytic MLE

# If we want to pull out the matrix of second derivatives
# evaluated at the MLE hessian <- opt.expo$hessian
# This is the same as -289/5, our analytic hessian
```

WHAT IS `OPTIM` DOING?

It depends on your choice in the `method` argument.

- ▶ Nelder-Mead: this is the default; it is slow but somewhat robust to non-differentiable functions.
- ▶ BFGS: a quasi-Newton Method; it is fast but needs a well behaved objective function.
- ▶ L-BFGS-B: similar to BFGS but allows box-constraints (i.e. upper and lower bounds on variables).
- ▶ CG: conjugate gradient method, may work for really large problems (we won't really use this).
- ▶ SANN: uses simulated annealing – a stochastic global optimization method; it is very robust but *very* slow.

WHAT IF YOU WANTED TO DO THE SAME NUMERICAL OPTIMIZATION “BY HAND”?

One approach to this is to use the Newton-Raphson algorithm (indeed, `optim`’s BFGS is a variation on this)

The Newton-Raphson algorithm iteratively finds an MLE for θ (in our case λ) by applying the following formula (see the appendix to these slides for the derivation):

$$\theta_{n+1} = \theta_n - \frac{f'(\theta_n)}{f''(\theta_n)}$$

Where $f(\theta_n)$ generally denotes your log-likelihood function.

We can implement this in R!

WHAT IF YOU WANTED TO DO THE SAME NUMERICAL OPTIMIZATION “BY HAND”?

To iterate over possible values of λ in our wait times example, we need to calculate:

$$\lambda_{n+1} = \lambda_n - \frac{\ell'(\lambda|y)}{\ell''(\lambda|y)}$$

We found $\ell'(\lambda|y)$ and $\ell''(\lambda|y)$ analytically earlier, so we can plug in the results:

$$\lambda_{n+1} = \lambda_n - \frac{n/\lambda_n - \sum_{i=1}^n y_i}{-n/\lambda_n^2}$$

Note that we'd previously plugged in 5 for n because we have 5 observations and 17 for $\sum_{i=1}^n y_i$ since our y observations add to 17, but this is the general form.

WHAT IF YOU WANTED TO DO THE SAME NUMERICAL OPTIMIZATION “BY HAND”?

Now we can write a function in R for our Newton-Raphson algorithm:

```
nr <- function(lambda,y){  
  return(lambda - (((length(y) / lambda) -  
    sum(y)) / (-length(y) / lambda^2)))  
}
```

WHAT IF YOU WANTED TO DO THE SAME NUMERICAL OPTIMIZATION “BY HAND”?

And loop over iterations of λ to see what our MLE for λ ($\hat{\lambda}$) converges to:

```
Y <- c(1,5,8,2,1) # observed Y
diff <- 10 # holder for diff. btw. lambda_n and lambda_n+1
init.lambda <- 0.01 # initialize lambda
while(diff > 0.0001){ # set to arbitrarily small gap,
  temp <- init.lambda # store prior lambda
  init.lambda <- nr(init.lambda, Y) # lambda_n+1
  print(init.lambda) # print lambda_n+1
  diff <- init.lambda - temp}
[1] 0.01966
[1] 0.03800585
[1] 0.07110058
[1] 0.1250132
[1] 0.1968901
[1] 0.2619768
[1] 0.2906053
[1] 0.2940757
[1] 0.2941176
```


WHAT IF YOU WANTED TO DO THE SAME NUMERICAL OPTIMIZATION “BY HAND”?

Boom. Converged to our analytical MLE for λ in 9 iterations.

7. FIND THE VARIANCE OF $\hat{\lambda}$

We are interested in calculating a measure of uncertainty of our MLE. That is, we are after the following:

$$\text{Var}(\hat{\theta}_{MLE})$$

Conceptually, we want to know how much information the MLE contains about the underlying parameter.

ASYMPTOTIC DISTRIBUTION OF THE MLE

It can be shown by the **Central Limit Theorem** that under certain regularity conditions, the MLE is distributed normally with a mean equal to the true parameter (θ_0) and the variance equal to the inverse of the expected sample Fisher information at the true parameter (denoted as $\mathcal{I}_n(\theta_0)$):

$$\hat{\theta}_{MLE} \sim \mathcal{N}\left(\theta, \underbrace{\left(-E\left[\frac{\partial^2 \ell(\theta|\mathbf{x})}{\partial \theta^2} \middle| \theta=\theta\right]\right)^{-1}}_{\mathcal{I}_n(\theta)}\right)$$

ASYMPTOTIC DISTRIBUTION OF THE MLE

Let's convince ourselves of this. We're going to generate 1000 datasets, of 10 observations each from an Exponential with $\lambda = 5/17 = 0.2941$.

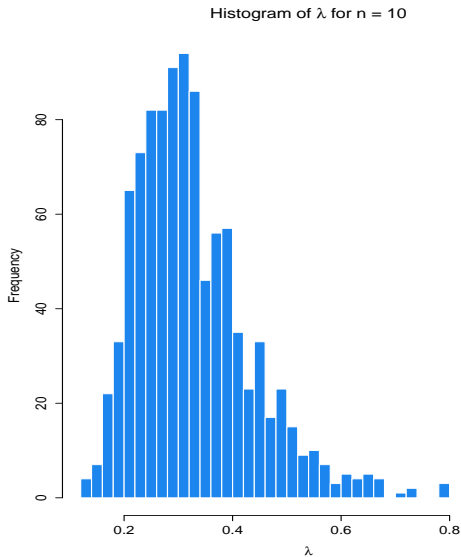
```
# We calculated analytically that the MLE is
true <- 5/17
n <- 10
set.seed(1234)
data <- sapply(seq(1,1000),function(x) rexp(n, rate=true))
dim(data)
[1] 10 1000
```

ASYMPTOTIC DISTRIBUTION OF THE MLE

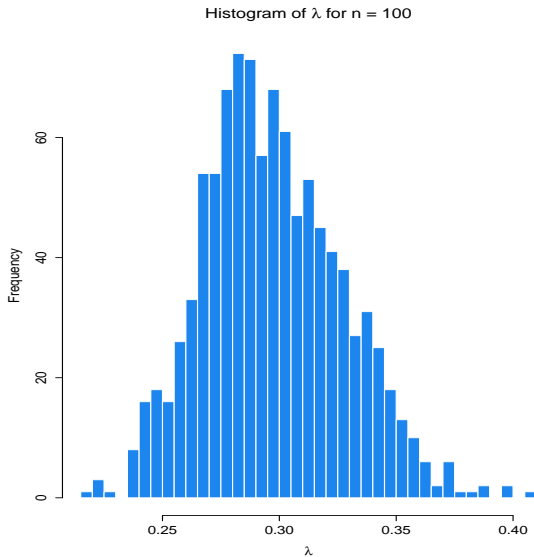
For each of these datasets, we're going to find the maximum likelihood estimate for λ .

```
llexp <- function(param, y){sum(dexp(y, rate=param, log=T))}  
out <- NULL  
for(i in 1:1000){  
  out[i] <- optim(c(1), fn=llexp, y=data[,i],  
    method="BFGS", control=list(fnscale=-1))$par  
}
```

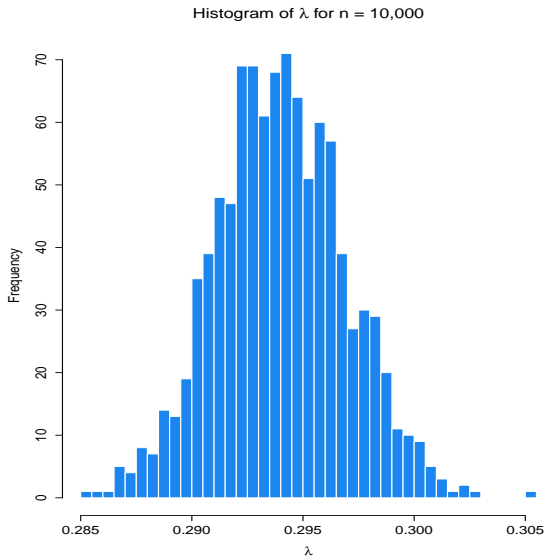
ASYMPTOTIC DISTRIBUTION OF THE MLE



ASYMPTOTIC DISTRIBUTION OF THE MLE



ASYMPTOTIC DISTRIBUTION OF THE MLE



ASYMPTOTIC DISTRIBUTION OF THE MLE

How do we think about this intuitively?

- ▶ The Central Limit Theorem states that the mean of independent random variables will become approximately Normal as n goes to infinity.
- ▶ But we're not talking about the mean!
- ▶ Yes, but the maximum of the log-likelihood is essentially the mean of a lot of likelihoods. And we use this maximum to estimate our parameter.
- ▶ Therefore, as n gets larger, and the more likelihoods we have to conglomerate, the more normal the distribution of the parameter becomes.

ASYMPTOTIC DISTRIBUTION OF THE MLE

So for large n , our parameter θ is distributed Normally with the mean as the true value of θ and variance $[\mathcal{I}(\theta)]^{-1}$:

$$\hat{\theta}_{MLE} \sim \mathcal{N}\left(\theta, \left(\mathcal{I}(\theta)\right)^{-1}\right)$$

VARIANCE OF MLE

Measure of curvature: Hessian Matrix

$$H(\hat{\theta}) = \left. \frac{\partial^2 \ell(\theta)}{\partial^2 \theta} \right|_{\theta=\hat{\theta}}$$

Measure of Information: Fisher Information Matrix

$$I(\hat{\theta}) = -H(\hat{\theta}) = - \left. \frac{\partial^2 \ell(\theta)}{\partial^2 \theta} \right|_{\theta=\hat{\theta}}$$

Estimate of Variance: Inverse of the Fisher Information Matrix

$$\hat{Var}(\hat{\theta}) \approx [I(\hat{\theta})]^{-1} \approx \left[- \left. \frac{\partial^2 \ell(\theta)}{\partial^2 \theta} \right|_{\theta=\hat{\theta}} \right]^{-1}$$

Estimate of Standard Error: Square root of $\hat{Var}(\hat{\theta})$

$$\hat{SE}(\hat{\theta}) = \sqrt{\hat{Var}(\hat{\theta})} = \sqrt{\left[- \left. \frac{\partial^2 \ell(\theta)}{\partial^2 \theta} \right|_{\theta=\hat{\theta}} \right]^{-1}}$$

VARIANCE OF MLE

- Asymptotically:

$$\text{Var}(\hat{\theta}_{MLE}) = [\mathcal{I}_n(\theta)]^{-1} = \left(-E \left[\frac{\partial^2 \ell(\theta | \mathbf{x})}{\partial \theta^2} \right] \right)^{-1}$$

- That is, it's the inverse of the **Expected Fisher Information** evaluated at the true parameter θ
- Conceptually, this is the *expected* curvature of the log-likelihood curve across repeated samples at the point θ
- We estimate this with the inverse of **Observed Fisher Information**:

$$\hat{\text{Var}}(\hat{\theta}_{MLE}) = [I(\hat{\theta})]^{-1} = \left(- \frac{\partial^2 \ell(\theta)}{\partial^2 \theta} \bigg|_{\theta=\hat{\theta}} \right)^{-1}$$

- As $n \rightarrow \infty$, the observed Fisher information converges to the expected Fisher information and the $\hat{\theta}_{MLE}$ converges to θ_0

7. FIND THE VARIANCE OF $\hat{\lambda}$

Remember that we already calculated the Hessian when we checked that $\hat{\lambda}$ is a maximum:

$$H(\hat{\lambda}) = -\frac{289}{5}$$

We can use this to calculate the standard error of $\hat{\lambda}$.

$$\begin{aligned} I(\hat{\lambda}) &= \frac{289}{5} \\ \hat{Var}(\hat{\lambda}) &= \frac{5}{289} \\ \hat{SE}(\hat{\lambda}) &= \sqrt{\frac{5}{289}} \end{aligned}$$

OUR UNIVARIATE EXAMPLE

Given our model and data

$$Y_i \sim f_{\text{expo}}(y_i | \lambda_i)$$

$$\lambda_i = \lambda$$

$$Y : \{1, 5, 8, 2, 1\}$$

- ▶ $L(\lambda|y) \propto \prod_{i=1}^n \lambda_i e^{-\lambda_i y_i}$
- ▶ $\ell(\lambda|y) = 5 \ln \lambda - 17\lambda + C$
- ▶ $S(\lambda) = \frac{5}{\lambda} - 17$
- ▶ $\hat{\lambda} = \frac{5}{17}$
- ▶ $H(\lambda) = -\frac{5}{\lambda^2}$
- ▶ $\hat{SE}(\hat{\lambda}) = \sqrt{\frac{5}{289}}$

OUTLINE

MULTIVARIATE

Ex. Waiting for the Redline – How long will it take for the next T to get here?



But this time we want to add covariates. What do you think affects the wait for the Redline?

1. WRITE THE MODEL

How would we model this?

- We know the stochastic component:

$$Y_i \sim f_{\text{expo}}(y_i | \lambda_i)$$

$$Y_i \sim \lambda_i e^{-\lambda_i y_i}$$

- Remember, for an Exponential

$$\mu_i = \frac{1}{\lambda_i}$$

- So we're going to set the systematic component

$$\mu_i = \exp(X_i \beta)$$

$$\lambda_i = \frac{1}{\exp(X_i \beta)}$$

What are the parameters?

2. CALCULATE $L(\lambda|y)$ AND 3. $\ell(\lambda|y)$

$$L(\lambda|y) \propto \prod_{i=1}^n \lambda_i e^{-\lambda_i y_i}$$

$$\ell(\lambda|y) \propto \sum_{i=1}^n (\ln \lambda_i - \lambda_i y_i)$$

4. PLUG IN SYSTEMATIC COMPONENT

$$\begin{aligned}\ell(\beta|y) &\propto \sum_{i=1}^n (\ln \lambda_i - \lambda_i y_i) \\ &\propto \sum_{i=1}^n \left(\ln \left(\frac{1}{\exp(X_i \beta)} \right) - \frac{1}{\exp(X_i \beta)} y_i \right) \\ &\propto \sum_{i=1}^n \left(\ln(1) - \ln(\exp(X_i \beta)) - \frac{1}{\exp(X_i \beta)} y_i \right) \\ &\propto \sum_{i=1}^n \left(-(X_i \beta) - \frac{1}{\exp(X_i \beta)} y_i \right)\end{aligned}$$

5. BRING IN OBSERVED DATA

I'm going to say whether or not it is Friday and the minutes behind schedule are important covariates.

Let's create some fake data:

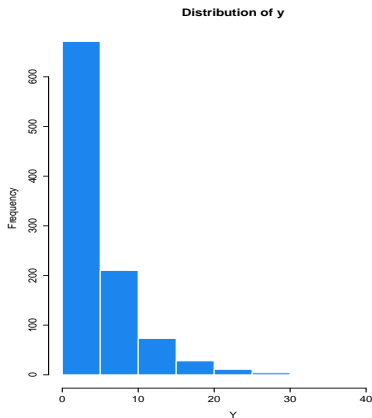
```
set.seed(02139)
n <- 1000
Friday <- sample(c(0,1), n, replace=T)
minsSch <- rnorm(n, 3, .5)

Y <- rexp(n, rate = 1/exp(1.25 - .5*Friday +.2*minsSch))
data <- as.data.frame(cbind(Y, Friday, minsSch))
```

5. BRING IN OBSERVED DATA

Let's look at Y:

```
hist(Y, col = "dodgerblue2", border = "white" ,  
main = "Distribution of y")
```



6. MAXIMIZE $\ell(\beta|y)$

Remember the log-likelihood we solved for before:

$$\ell(\beta|y) = \sum_{i=1}^n \left(-(X_i\beta) - \frac{1}{\exp(X_i\beta)} y_i \right)$$

We can find the MLE by setting the score to 0. But with multiple parameters, the score is now:

► Score:

$$S(\beta) = \nabla \ell(\beta) = \begin{pmatrix} \frac{\partial \ell(\beta)}{\partial \beta_1} \\ \frac{\partial \ell(\beta)}{\partial \beta_2} \\ \vdots \\ \frac{\partial \ell(\beta)}{\partial \beta_k} \end{pmatrix}$$

Set **each** element to 0 and solve the system of equations to get $\hat{\beta}$.

6. MAXIMIZE $\ell(\beta|y)$

OR, we can solve for the MLE in R by first programming the log-likelihood.

```
llexp <- function(param, y, x){  
  rate <- 1/exp(x%*%param)  
  sum(dexp(y, rate=rate, log=T))  
}  
  
llexp2 <- function(param, y,x){  
  cov <- x%*%param  
  sum(-cov - 1/exp(cov)*y)  
}
```

6. MAXIMIZE $\ell(\beta|y)$

We can maximize our function using `optim`:

```
#Create X with an intercept
X <- cbind(1, Friday, minsSch)

#Specify starting values for all three parameters
param <- c(1,1,1)

#Solve using optim
out <- optim(param, fn=llexp, y=Y, x=X, method="BFGS",
             hessian=T, control=list(fnscale=-1))

out$par
[1] 1.0885871 -0.4634621 0.2120591
```


7. FIND THE VARIANCE OF $\hat{\beta}$

We can find the variance of the MLE with:

$$\hat{Var}(\hat{\beta}) = [-H(\hat{\beta})]^{-1}$$

But with multiple parameters, the Hessian is now:

► Hessian:

$$H(\hat{\beta}) = \left(\begin{array}{cccc} \frac{\partial^2}{\partial \beta_1^2} & \frac{\partial^2}{\partial \beta_1 \partial \beta_2} & \cdots & \frac{\partial^2}{\partial \beta_1 \partial \beta_k} \\ \frac{\partial^2}{\partial \beta_2 \partial \beta_1} & \frac{\partial^2}{\partial \beta_2^2} & \cdots & \frac{\partial^2}{\partial \beta_2 \partial \beta_k} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2}{\partial \beta_k \partial \beta_1} & \frac{\partial^2}{\partial \beta_k \partial \beta_2} & \cdots & \frac{\partial^2}{\partial \beta_k^2} \end{array} \right) \bigg|_{\beta=\hat{\beta}}$$

7. FIND THE VARIANCE OF $\hat{\beta}$

We can calculate this in R using the output from `optim`.

```
# Get the Hessian from optim
H <- out$hessian

# Calculate the observed fisher information
I <- -H

# Calculate the variance-covariance matrix
V <- solve(I)

# Get the standard errors
ses <- sqrt(diag(V))
```

OUTLINE

LIKELIHOOD RATIO TEST

- ▶ Useful for when you are comparing two models.
- ▶ We'll call these restricted and unrestricted:

$$\textit{Unrestricted} \quad : \quad \beta_0 + \beta_1 X_1 + \beta_2 X_2$$

$$\textit{Restricted} \quad : \quad \beta_0 + \beta_2 X_2$$

- ▶ We want to test the usefulness of the parameters in the unrestricted model but omitted in the restricted model (eg. β_1).

LIKELIHOOD RATIO TESTS

Here's how to operationalize this:

- ▶ Let L_u^* be the maximum of the unrestricted likelihood, and let L_r^* the maximum of the restricted likelihood.
- ▶ But adding more variables can only increase the likelihood.
- ▶ Thus, $L_u^* \geq L_r^*$, or $\frac{L_r^*}{L_u^*} \leq 1$ always.
- ▶ If the likelihood ratio is exactly 1, then there's no effect of the extra parameters at all ($L_u^* = L_r^*$).

LIKELIHOOD RATIO TEST

Now, let's define a test statistic:

$$\begin{aligned}\text{define : } \mathfrak{R} &= -2 \ln \frac{L_r^*}{L_u^*} \\ &= 2(\ln L_u^* - \ln L_r^*)\end{aligned}$$

$$\mathfrak{R} \sim \chi_m^2$$

$$\text{Reject if : } \mathfrak{R} > \chi_{m,1-\alpha}^2$$

- ▶ m is the number of restrictions.
- ▶ \mathfrak{R} will always be greater than zero.
- ▶ Key question: how much greater than zero does \mathfrak{R} have to be in order to convince us that the difference is due to systematic differences between the two models?

BACK TO OUR EXAMPLE

What if we wanted to test whether the minutes behind schedule should be in our model at all?

```
unrestricted <- optim(param, fn=llexp, y=Y, x=X,  
  method="BFGS", hessian=T, control=list(fnscale=-1))  
unrestricted$value  
[1] -2503.445
```

V.

```
restricted <- optim(c(1,1), fn=llexp, y=Y,  
  x=cbind(1, Friday), method="BFGS",  
  hessian=T, control=list(fnscale=-1))  
restricted$value  
[1] -2509.471
```

BACK TO OUR EXAMPLE

Under the null that the restrictions are valid, the test statistic would be distributed χ^2 with one degree of freedom:

```
# Calculate our test statistic
r <- 2*(unrestricted$value - restricted$value)

# Calculate the p-value for this test statistic
1-pchisq(r,df=1)
[1] 0.0005176814
```

So the probability of getting this test statistic under the null is extremely small. We reject.

QUESTIONS

Questions?

OUTLINE

OPTIMIZATION STRATEGIES

This appendix provides alternative methods of optimization.

NEWTON'S METHOD

Newton's method: a pretty good approach for a continuous and twice-differentiable function. We'll look at a univariate function here. Suppose we know our function $f(\cdot)$ and we have

a starting value of x_0 . Our goal is to find move from x_0 to x_1 such that $f'(x_1) = 0$ (or is at least closer to 0 than at x_0). This will be a sequential process of approximation and

eventually $f'(x_n)$ will be close enough to zero to let us declare that x_n a critical value.

TAYLOR SERIES EXPANSION: A STAPLE OF CALCULUS

We can approximate function $f(\cdot)$ at point a using a Taylor series expansion:

$$f(x) = f(a) + f'(a)(x - a) + \frac{f''(a)}{2!}(x - a)^2 + \frac{f'''(a)}{3!}(x - a)^3 + \dots$$

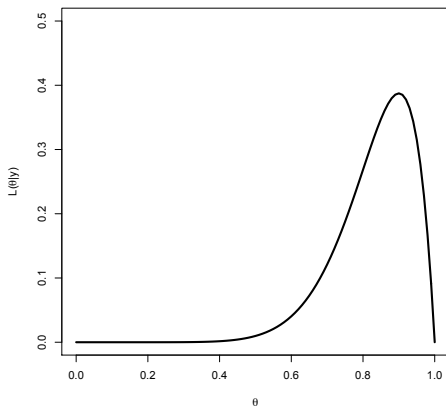
In this course, we'll work with the first- and second-order Taylor polynomials:

$$f(x) \approx f(a) + f'(a)(x - a)$$

$$f(x) \approx f(a) + f'(a)(x - a) + \frac{f''(a)}{2!}(x - a)^2$$

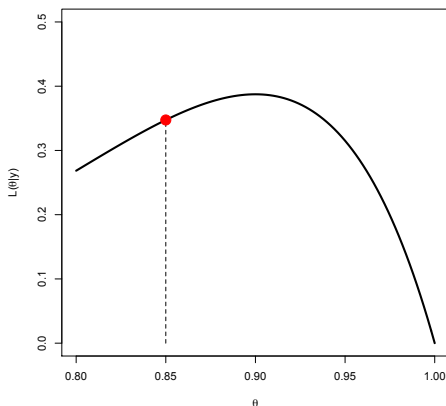
LIKELIHOOD OF BINOMIAL DISTRIBUTION

This is the likelihood of a binomial distribution with 10 trials ($N = 10$) from which we drew one observation: $y = 9$.



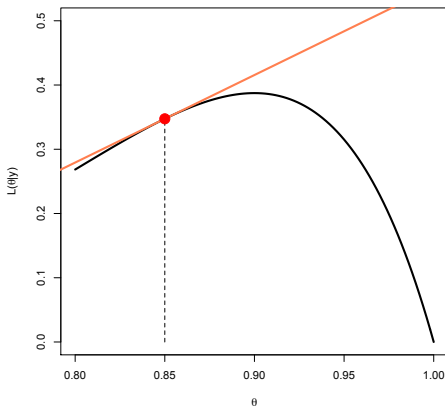
LIKELIHOOD OF BINOMIAL DISTRIBUTION: ZOOMED IN

We want to approximate the likelihood curve around $\theta_0 = 0.85$.



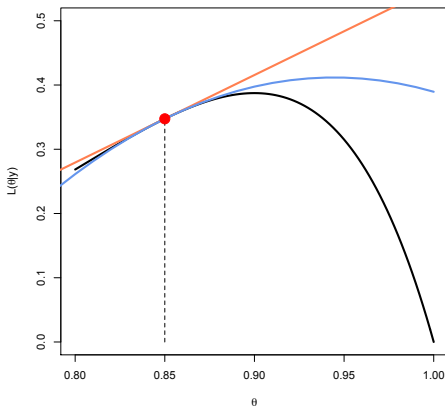
LIKELIHOOD OF BINOMIAL DISTRIBUTION: ZOOMED IN

We can approximate the likelihood around $\theta_0 = 0.85$ using a first-order Taylor polynomial.



LIKELIHOOD OF BINOMIAL DISTRIBUTION: ZOOMED IN

We can improve our approximation of the likelihood around $\theta_0 = 0.85$ by using a second-order Taylor polynomial.



WHY DON'T WE JUST MAXIMIZE THE SECOND-ORDER TAYLOR POLYNOMIAL?

We can write a second-order Taylor expansion around θ_0 as:

$$f(\theta) \approx f(\theta_0) + f'(\theta_0)(\theta - \theta_0) + \frac{f''(\theta_0)}{2!}(\theta - \theta_0)^2$$

To maximize, take the derivative with respect to θ and set it equal to 0:

$$f'(\theta) = f'(\theta_0) + f''(\theta_0)(\theta - \theta_0) = 0$$

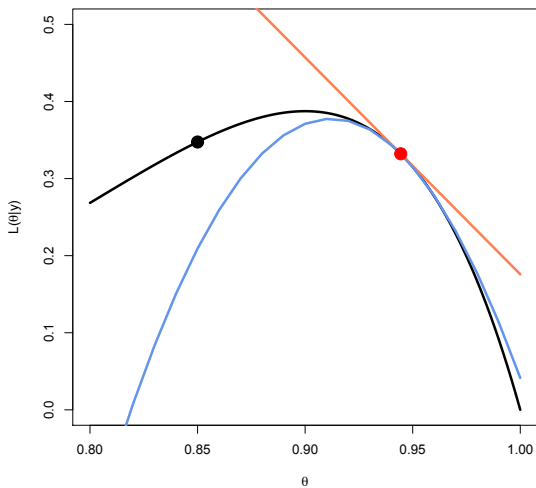
Rearranging:

$$\theta = \theta_0 - \frac{f'(\theta_0)}{f''(\theta_0)}$$

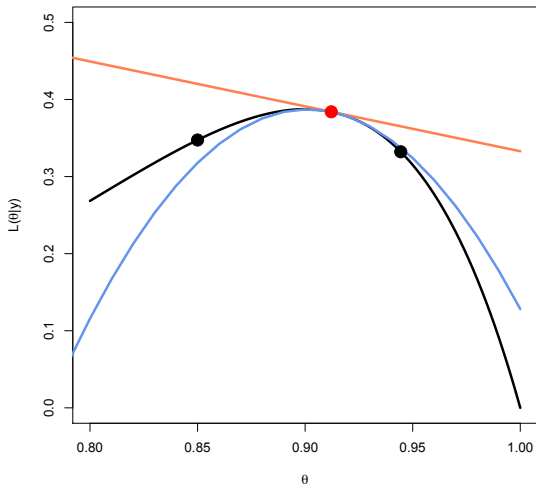
WE'VE JUST DERIVED THE UPDATE STEP WE WANTED!

$$\theta_{n+1} = \theta_n - \frac{f'(\theta_n)}{f''(\theta_n)}$$

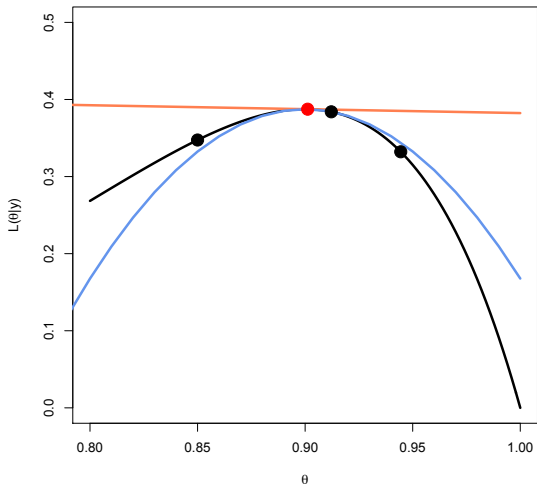
FIRST ITERATION



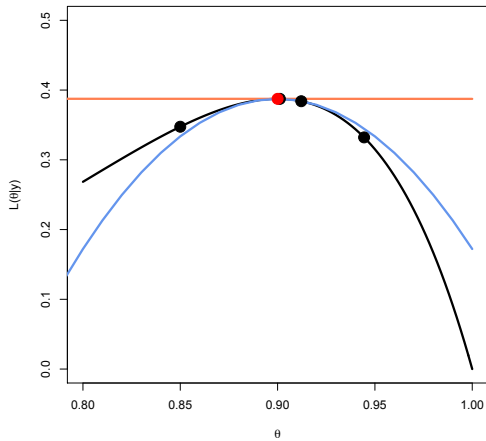
SECOND ITERATION



THIRD ITERATION



FOURTH ITERATION



ANOTHER EXAMPLE OF TAYLOR SERIES APPROXIMATION

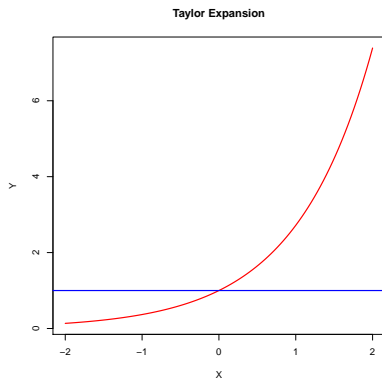


Figure : The exponential function, $g(x) = e^x$, and the Taylor Series approximation: $x_0 = 0, g_0(x_1) = 1$ (from Wikipedia)

ANOTHER EXAMPLE OF TAYLOR SERIES APPROXIMATION

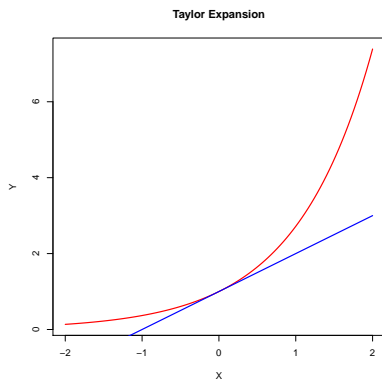


Figure : The exponential function, $g(x) = e^x$, and the Taylor Series approximation: $x_0 = 0, g_1(x_1) = 1 + x_1$ (from Wikipedia)

ANOTHER EXAMPLE OF TAYLOR SERIES APPROXIMATION

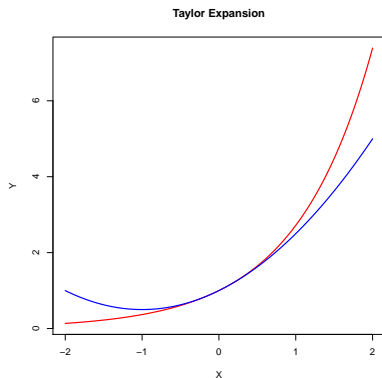


Figure : The exponential function, $g(x) = e^x$, and the Taylor Series approximation: $x_0 = 0, g_2(x_1) = 1 + x_1 + \frac{x_1^2}{2}$ (from Wikipedia)

ANOTHER EXAMPLE OF TAYLOR SERIES APPROXIMATION

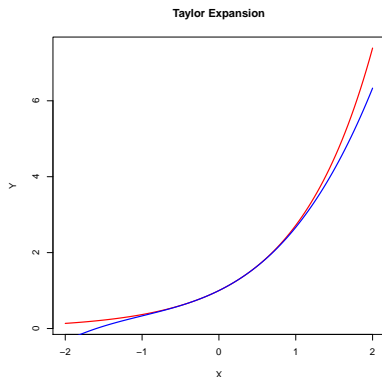


Figure : The exponential function, $g(x) = e^x$, and the Taylor Series approximation: $x_0 = 0, g_3(x_1) = 1 + x_1 + \frac{x_1^2}{2} + \frac{x_1^3}{6}$ (from Wikipedia)

ANOTHER EXAMPLE OF TAYLOR SERIES APPROXIMATION

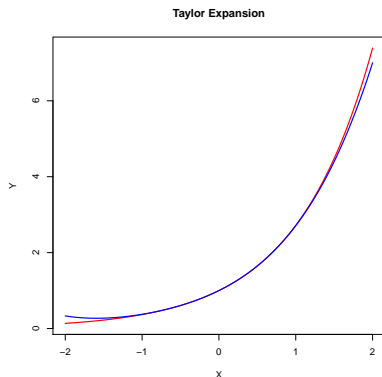


Figure : The exponential function, $g(x) = e^x$, and the Taylor Series approximation: $x_0 = 0, g_4(x) = 1 + x + \frac{x^2}{2} + \frac{x^3}{6} + \frac{x^4}{24}$ (from Wikipedia)

NEWTON IN ACTION: BERNOULLI EXAMPLE REVISITED

Let's maximize our likelihood: $3 \ln \pi + 2 \ln(1 - \pi)$.

Recall that $L'(\pi) = \frac{3}{\pi} - \frac{2}{1-\pi}$ and $L''(\pi) = -\frac{3}{\pi^2} - \frac{2}{(1-\pi)^2}$.

Starting at $\pi_0 = .3$, we use our updating formula:

$$\pi_1 = \pi_0 - \frac{L'(\pi_0)}{L''(\pi_0)} = .3 - \frac{L'(.3)}{L''(.3)} = 0.4909.$$

Now use $\pi_1 = .4909$ as a starting value.

$$\pi_2 = \pi_1 - \frac{L'(\pi_1)}{L''(\pi_1)} = .4909 - \frac{L'(.4909)}{L''(.4909)} = 0.5991.$$

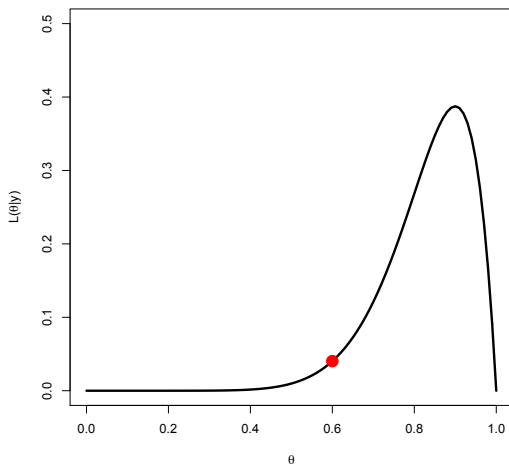
So we're already there!

PROPERTIES OF NEWTON-RAPHSON

- ▶ Converges quickly
- ▶ Can get stuck in local minima/maxima
- ▶ Can have troubles with root jumping
- ▶ Won't walk at all on a flat space

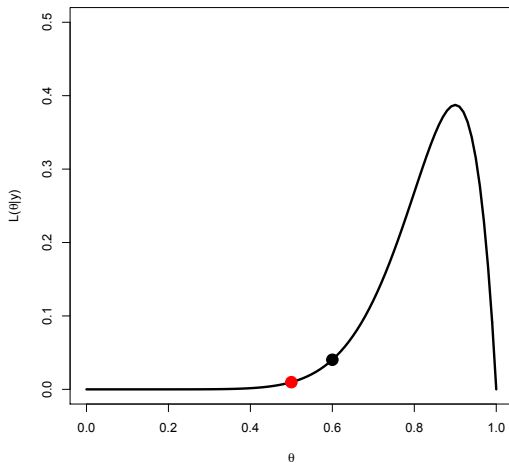
NEWTON-RAPHSON GONE AWRY

What if we had taken $\theta_0 = 0.60$ to be the starting point for the Newton-Raphson maximization for the binomial likelihood?



IT WALKS THE WRONG WAY...

The first iteration:



AND THEN GETS STUCK!

After 10 iterations:

