

Sam Leonard

In this paper I'll explore a four different algorithms for obtaining estimates for the parameters of a probit regression: Newton-Raphson, Fisher scoring, coordinate descent, and stochastic gradient.

Recall that probit regression estimates a categorical outcome:

$$y_i | \theta_i \sim \text{Ber}(\theta_i)$$

$$\theta_i = \Phi(x_i^T \beta)$$

Section one is the tough part: the math.

In the second section I will show the code I wrote for executing each of these algorithms.

And finally, section three will be a brief comparison of the results against a simulated dataset.

Section 1 - the Math

Newton method:

$$\log(p(y|\beta)) = \sum_{i=1}^n (y_i \log(\Phi(x_i^T \beta)) + (1 - y_i) \log(1 - \Phi(x_i^T \beta)))$$

For some $r \in \{1, \dots, p\}$:

$$\frac{\partial l}{\partial \beta_r} = \sum_{i=1}^n \left(\frac{y_i \phi(x_i^T \beta)}{\Phi(x_i^T \beta)} - \frac{(1 - y_i) \phi(x_i^T \beta)}{1 - \Phi(x_i^T \beta)} \right) x_{i,r}$$

This expression appears for each β_r in the matrix:

$$\nabla l(\beta) = \begin{bmatrix} \frac{\partial l}{\partial \beta_1} \\ \frac{\partial l}{\partial \beta_2} \\ \vdots \\ \frac{\partial l}{\partial \beta_p} \end{bmatrix}$$

For some s where $s \in \{1, \dots, p\}$:

$$\frac{\partial^2 l}{\partial \beta_r \partial \beta_s} = \sum_{i=1}^n \left(\frac{y_i \phi(x_i^T \beta) (-2 * \frac{1}{2} * x_i^T \beta)}{\Phi(x_i^T \beta)} - \frac{y_i \phi(x_i^T \beta) \phi(x_i^T \beta)}{\Phi(x_i^T \beta)^2} - \frac{(1 - y_i) \phi(x_i^T \beta) (-x_i^T \beta)}{1 - \Phi(x_i^T \beta)} - \frac{(1 - y_i) \phi(x_i^T \beta) (-\phi(x_i^T \beta))}{(1 - \Phi(x_i^T \beta))^2} \right) x_{i,r} x_{i,s}$$

So

$$\frac{\partial^2 l}{\partial \beta_r \partial \beta_s} = \sum_{i=1}^n \left(- \frac{y_i \phi(x_i^T \beta) x_i^T \beta}{\Phi(x_i^T \beta)} - \frac{y_i \phi(x_i^T \beta)^2}{\Phi(x_i^T \beta)^2} + \frac{(1 - y_i) \phi(x_i^T \beta) x_i^T \beta}{1 - \Phi(x_i^T \beta)} - \frac{(1 - y_i) \phi(x_i^T \beta)^2}{(1 - \Phi(x_i^T \beta))^2} \right) x_{i,r} x_{i,s}$$

This expression for each r and s completes the Hessian matrix:

$$Hes(\beta) = \begin{bmatrix} \frac{\partial^2 l}{\partial \beta_1 \partial \beta_1} & \frac{\partial^2 l}{\partial \beta_1 \partial \beta_2} & \cdots & \frac{\partial^2 l}{\partial \beta_1 \partial \beta_p} \\ \frac{\partial^2 l}{\partial \beta_2 \partial \beta_1} & \frac{\partial^2 l}{\partial \beta_2 \partial \beta_2} & \cdots & \frac{\partial^2 l}{\partial \beta_2 \partial \beta_p} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 l}{\partial \beta_p \partial \beta_1} & \frac{\partial^2 l}{\partial \beta_p \partial \beta_2} & \cdots & \frac{\partial^2 l}{\partial \beta_p \partial \beta_p} \end{bmatrix}$$

To implement Newton-Raphson, we start with some guess $\beta^{(0)}$. Then

$$\beta^{(1)} = \beta^{(0)} - [Hes(\beta^{(0)})]^{-1} * (\nabla l(\beta^{(0)}))$$

And generally

$$\beta^{(t+1)} = \beta^{(t)} - [Hes(\beta^{(t)})]^{-1} * (\nabla l(\beta^{(t)}))$$

This sequence converges to the minimum of $\nabla l(\beta)$. The log of the probit likelihood is unimodal, therefore this minimum will be the MLE for β .

A reasonable stopping condition is when the relative change in log likelihood is very small: $\left| \frac{l(\beta^{(t+1)}) - l(\beta^{(t)})}{l(\beta^{(t+1)})} \right| < \epsilon$, where ϵ is chosen in advance.

Fisher Scoring

Each y_i is Bernoulli with $p = \Phi(x_i^T \beta)$. Thus $E[y_i] = \Phi(x_i^T \beta)$. Entries of the expected information matrix take the form

$$E\left[\frac{\partial^2 l}{\partial \beta_r \partial \beta_s}\right] = \sum_{i=1}^n \left(-\frac{\Phi(x_i^T \beta) \phi(x_i^T \beta) x_i^T \beta}{\Phi(x_i^T \beta)} - \frac{\Phi(x_i^T \beta) \phi(x_i^T \beta)^2}{\Phi(x_i^T \beta)^2} + \frac{(1 - \Phi(x_i^T \beta)) \phi(x_i^T \beta) x_i^T \beta}{1 - \Phi(x_i^T \beta)} - \frac{(1 - \Phi(x_i^T \beta)) \phi(x_i^T \beta)^2}{(1 - \Phi(x_i^T \beta))^2} \right) x_{i,r} x_{i,s}$$

So

$$E\left[\frac{\partial^2 l}{\partial \beta_r \partial \beta_s}\right] = \sum_{i=1}^n \left(\frac{-\Phi(x_i^T \beta)(1 - \Phi(x_i^T \beta)) - \phi(x_i^T \beta)^2 \Phi(x_i^T \beta)}{\Phi(x_i^T \beta)(1 - \Phi(x_i^T \beta))} \right) x_{i,r} x_{i,s} = - \sum_{i=1}^n \left(\frac{\phi(x_i^T \beta)^2}{\Phi(x_i^T \beta)(1 - \Phi(x_i^T \beta))} \right) x_{i,r} x_{i,s}$$

The negative of this expression completes the expected information matrix:

$$I_E(\beta) = - \begin{bmatrix} E\left[\frac{\partial^2 l}{\partial \beta_1 \partial \beta_1}\right] & E\left[\frac{\partial^2 l}{\partial \beta_1 \partial \beta_2}\right] & \cdots & E\left[\frac{\partial^2 l}{\partial \beta_1 \partial \beta_p}\right] \\ E\left[\frac{\partial^2 l}{\partial \beta_2 \partial \beta_1}\right] & E\left[\frac{\partial^2 l}{\partial \beta_2 \partial \beta_2}\right] & \cdots & E\left[\frac{\partial^2 l}{\partial \beta_2 \partial \beta_p}\right] \\ \vdots & \vdots & \ddots & \vdots \\ E\left[\frac{\partial^2 l}{\partial \beta_p \partial \beta_1}\right] & E\left[\frac{\partial^2 l}{\partial \beta_p \partial \beta_2}\right] & \cdots & E\left[\frac{\partial^2 l}{\partial \beta_p \partial \beta_p}\right] \end{bmatrix}$$

To execute the Fisher Scoring algorithm start with some guess $\beta^{(0)}$. Then, because we defined expected information as the negative of the expectation of the Hessian, we have:

$$\beta^{(1)} = \beta^{(0)} + [I_E(\beta^{(0)})]^{-1} * \nabla l(\beta^{(0)})$$

And

$$\beta^{(t+1)} = \beta^{(t)} + [I_E(\beta^{(t)})]^{-1} * \nabla l(\beta^{(t)})$$

The algorithm stops again when the relative change in the log likelihood from $l(\beta^t)$ to $l(\beta^{t+1})$ is smaller than some pre-selected ϵ .

Coordinate Descent

Start with some guess for the vector β , call it $\beta^{(0)}$. Assume that $l(\beta)$ is a function of only β_r . All other elements of the vector β are fixed. Let β_{-r} be β excluding β_r . To execute a coordinate descent algorithm we start by finding β_r that minimizes $\frac{dl}{d\beta_r}$ given β_{-r} .

The first derivative of $l(\beta)$ with respect to β_r is

$$\frac{dl}{d\beta_r} = \sum_{i=1}^n \left(\frac{y_i \phi(x_i^T \beta)}{\Phi(x_i^T \beta)} - \frac{(1 - y_i) \phi(x_i^T \beta)}{1 - \Phi(x_i^T \beta)} \right) x_{i,r}$$

No closed form solution for $\frac{dl}{d\beta_r} = 0$ is apparent, so we can use Newton-Raphson in one dimension to find the β_r that minimizes $\frac{dl}{d\beta_r}$. The second derivative is

$$\frac{d^2 l}{d^2 \beta_r} = \sum_{i=1}^n \left(-\frac{y_i \phi(x_i^T \beta) x_i^T \beta}{\Phi(x_i^T \beta)^2} - \frac{y_i \phi(x_i^T \beta)^2}{\Phi(x_i^T \beta)^2} + \frac{(1 - y_i) \phi(x_i^T \beta) x_i^T \beta}{1 - \Phi(x_i^T \beta)} - \frac{(1 - y_i) \phi(x_i^T \beta)^2}{(1 - \Phi(x_i^T \beta))^2} \right) x_{i,r} x_{i,r}$$

The one-dimensional Newton-Raphson algorithm for finding β_r that minimizes $\frac{dl}{d\beta_r}$ given β_{-r} is

$$\beta_r^{(1)} = \beta_r^{(0)} - \frac{\frac{dl}{d\beta_r}(\beta^{(0)})}{\frac{d^2 l}{d^2 \beta_r}(\beta^{(0)})}$$

and generally

$$\beta_r^{(t+1)} = \beta_r^{(t)} - \frac{\frac{dl}{d\beta_r}(\beta^{(t)})}{\frac{d^2 l}{d^2 \beta_r}(\beta^{(t)})}$$

Continue this process until the relative change in $\beta_r^{(t)}$ is less than some pre-selected ϵ .

Once we have minimized $\frac{dl}{d\beta_r}$, the r^{th} element of the vector β is now this minimizing value. Move to another element of the vector β , β_s . Use the same process to minimize $\frac{dl}{d\beta_s}$ given β_{-s} .

Continue this process until each of the p Newton-Raphson sequences have converged. A reasonable stopping condition is when each of the p sequences stop after one iteration sequentially.

Stochastic Gradient

Again start with a reasonable guess for β , $\beta^{(0)}$.

To implement a stochastic gradient algorithm select m data points at random from your dataset. Then update $\beta^{(0)}$ by finding

$$\beta^{(1)} = \beta^{(0)} + \lambda_0 * [\nabla l_m(\beta^{(0)})]$$

and generally

$$\beta^{(t+1)} = \beta^{(t)} + \lambda_t * [\nabla l_m(\beta^{(t)})]$$

Here $\nabla l_m(\beta)$ the gradient of $l(\beta)$ evaluated over the data points m . In our case this is:

$$\nabla l_m(\beta) = \sum_{i \in m} \left(\frac{y_i \phi(x_i^T \beta)}{\Phi(x_i^T \beta)} - \frac{(1 - y_i) \phi(x_i^T \beta)}{1 - \Phi(x_i^T \beta)} \right) x_i$$

And $\lambda_t = \frac{\alpha}{1+\eta t}$ is a learning rate which influences the scale of movement from $\beta^{(t)}$ to $\beta^{(t+1)}$.

The m data points selected at each iteration are discarded after the completion of that iteration, so that no data point is randomly selected twice. The algorithm stops when a certain number of iterations has passed, or we see a stabilization in the log likelihood.

Section 2 - code

Here is the code for executing each algorithm.

First I'll simulate a dataset with

$$n = 100,000$$

$$x_{i,j} \sim N(0, 2)$$

$$\theta_i = \Phi(.2 + 2.4x_{i,1} - .5x_{i,3} + 1.2x_{i,7})$$

```
n = 100000
p = 7
x = matrix(rnorm(700000, 0, 2), ncol = 7)
x = cbind(1, x)
totals = .2*x[, 1]+2.4*x[, 2] - .5*x[, 4] + 1.2*x[, 8]
theta = pnorm(totals)
y = sapply(theta, function(t)rbinom(1, 1, prob = t))
true.beta = c(.2, 2.4, 0, -.5, 0, 0, 0, 1.2)

likelihood.beta <- function(beta){
  sums = apply(x, 1, function(t)sum(t*beta))
  l.term1 = pnorm(sums, log = TRUE)
  l.term2 = pnorm(-sums, log = TRUE)
  return(sum(y*l.term1+(1-y)*l.term2))
}
```

Newton-Raphson

```
loop = TRUE
beta = rep(0, 8)
l = 0

start <- Sys.time()
while (loop){
  info_matrix <- matrix(0, nrow = p+1, ncol = p+1)
  sums = apply(x, 1, function(t)sum(t*beta))
  dd.term1 = dnorm(sums)*sums/pnorm(sums)
  dd.term2 = (dnorm(sums)/pnorm(sums))**2
  dd.term3 = dnorm(sums)*sums/pnorm(-sums)
```

```

dd.term4 = (dnorm(sums)/pnorm(-sums))**2
d.term1 = dnorm(sums)/pnorm(sums)
d.term2 = dnorm(sums)/pnorm(-sums)

for (i in 1:(p+1)){
  for (j in i:(p+1)){
    x.1 = x[, i]
    x.2 = x[, j]
    info_matrix[i, j] = info_matrix[j, i] =
      sum((-y*dd.term1 - y*dd.term2 + (1 - y)*dd.term3 - (1 - y)*dd.term4)*x.1*x.2)
  }
}

info.inv = solve(info_matrix)
dl.beta.vec = rep(0, p+1)
for (i in 1:(p+1)){
  x.iter.1 = x[, i]
  dl.beta.vec[i] = sum((y*d.term1 - (1-y)*d.term2)*x.iter.1)
}

beta.t1 = beta - info.inv %*% dl.beta.vec
l = l + 1
if (abs((likelihood.beta(beta.t1)-likelihood.beta(beta))/likelihood.beta(beta.t1)) < 10**(-4)){
  end <- Sys.time()
  cat('beta0: ', beta.t1[1], '\n')
  cat('beta1: ', beta.t1[2], '\n')
  cat('beta2: ', beta.t1[3], '\n')
  cat('beta3: ', beta.t1[4], '\n')
  cat('beta4: ', beta.t1[5], '\n')
  cat('beta5: ', beta.t1[6], '\n')
  cat('beta6: ', beta.t1[7], '\n')
  cat('beta7: ', beta.t1[8], '\n')
  cat('iterations: ', l, '\n')
  print(end - start)
  loop = FALSE
}
beta = beta.t1
}

## beta0: 0.1945283
## beta1: 2.405935
## beta2: -0.001431506
## beta3: -0.5042822
## beta4: -0.0007086811
## beta5: 0.008327334
## beta6: 0.002512334
## beta7: 1.199156
## iterations: 7
## Time difference of 12.65184 secs

beta.newton = beta.t1
time.newton = end - start

```

Coordinate Descent

```

beta.0 = rep(0, 8)
beta.1 = beta.0
error.vector = rep(1, p+1)
errors = c()
location = 1
l = 0
loop = TRUE

start <- Sys.time()
while (loop){
  likelihood.0 = likelihood.beta(beta.0)
  sums = apply(x, 1, function(t)sum(t*beta.0))
  dd.term1 = dnorm(sums)*sums/pnorm(sums)
  dd.term2 = (dnorm(sums)/pnorm(sums))**2
  dd.term3 = dnorm(sums)*sums/pnorm(-sums)
  dd.term4 = (dnorm(sums)/pnorm(-sums))**2
  d.term1 = dnorm(sums)/pnorm(sums)
  d.term2 = dnorm(sums)/pnorm(-sums)

  beta.1[location] = beta.0[location] -
    sum((y*d.term1 - (1-y)*d.term2)*x[, location])/
    sum((-y*dd.term1 - y*dd.term2 + (1 - y)*dd.term3 - (1 - y)*dd.term4)*x[, locat:

  likelihood.1 = likelihood.beta(beta.1)
  err = abs((likelihood.1 - likelihood.0)/likelihood.1)
  errors = append(errors, err)
  if(err < 10**-4){
    error.vector[location] = max(errors)
    errors = c()
    location = (location %% 8) + 1
  }
  if(all(error.vector < 10**-4)){
    l = l + 1
    end <- Sys.time()
    cat('beta0: ', beta.1[1], '\n')
    cat('beta1: ', beta.1[2], '\n')
    cat('beta2: ', beta.1[3], '\n')
    cat('beta3: ', beta.1[4], '\n')
    cat('beta4: ', beta.1[5], '\n')
    cat('beta5: ', beta.1[6], '\n')
    cat('beta6: ', beta.1[7], '\n')
    cat('beta7: ', beta.1[8], '\n')
    cat('iterations: ', l, '\n')
    print(end - start)
    loop = FALSE
  }
  beta.0 = beta.1
  l = l + 1
}

## beta0: 0.1899355
## beta1: 2.357691
## beta2: -0.001422168
## beta3: -0.4920079

```

```
## beta4: -0.000836956
## beta5: 0.008106912
## beta6: 0.002361728
## beta7: 1.170991
## iterations: 176
## Time difference of 2.813666 mins
beta.coord = beta.1
time.coord = end - start
```

Stochastic Gradient

I ran the stochastic gradient algorithms against standardized data.

In a simple case of one intercept and one covariate, against standardized data we have $\theta_i = \Phi(\beta_0 + \beta_1(\frac{x_{i,1} - \bar{x}_1}{s_1})) = \Phi((\beta_0 - \frac{\bar{x}_1}{s_1}) + \frac{\beta_1}{s_1}x_{i,1})$.

Thus to compare scaled MLEs to the MLEs already generated, we simply need to divide by the standard deviation of the data column x , except in the case of the intercept. In the case of the intercept we simply subtract the mean of each column divided by its standard deviation.

This dataset was generated by sampling from a normal distribution with mean 0 and standard deviation 2. Therefore scaling the data will not change the intercept parameter, and the remaining parameters generated from standardized data will be twice the magnitude of those generated from unstandardized data.

For $m = 1000$, $\alpha = .0045$ and $\eta = .05$ stabilized around good estimates of the β used to generate the data after just 100 iterations.

```
m = .01*n
x.scaled <- cbind(rep(1, n), scale(x[, -1]))
location = 1
beta.0 = rep(0, 8)
beta.1 = beta.0
alpha = .0083
eta = .05
niter = 0
indices = seq(1:100000)
likelihood.vec.1 <- rep(0, 100)

start <- Sys.time()
while(niter < 100){
  beta.0 = beta.1
  samp = sample(indices, m)
  indices = indices[!indices %in% samp]
  mat = x.scaled[samp, ]
  sums = apply(mat, 1, function(t)sum(t*beta.0))
  d.term1 = apply(mat, 1, function(t)exp(dnorm(sum(t*beta.0), log = TRUE) - pnorm(sum(t*beta.0), log = TRUE)))
  d.term2 = apply(mat, 1, function(t)exp(dnorm(sum(t*beta.0), log = TRUE) - (pnorm(-sum(t*beta.0), log = TRUE))))
  for (location in 1:(p+1)){
    beta.1[location] = beta.0[location] + alpha/(1+niter*eta) * sum((y[samp]*d.term1 - (1-y[samp])*d.term2))
  }
  location = 1
  beta.scaled = .5*beta.1
  beta.scaled[1] = beta.1[1]
  likelihood.vec.1[niter+1] = likelihood.beta(beta.scaled)
  niter = niter + 1
}
```

```

test = TRUE
if(test){
  end <- Sys.time()
  cat('beta0: ', beta.1[1], '\n')
  cat('beta1: ', beta.1[2], '\n')
  cat('beta2: ', beta.1[3], '\n')
  cat('beta3: ', beta.1[4], '\n')
  cat('beta4: ', beta.1[5], '\n')
  cat('beta5: ', beta.1[6], '\n')
  cat('beta6: ', beta.1[7], '\n')
  cat('beta7: ', beta.1[8], '\n')
  cat('iterations: ', niter, '\n')
  print(end - start)
  beta.m1000 = beta.1
  time.m1000 = end - start
}

```

```

## beta0: 0.1829191
## beta1: 4.823844
## beta2: -0.001875432
## beta3: -0.9771004
## beta4: -0.02555673
## beta5: 0.009351828
## beta6: 0.03160992
## beta7: 2.411478
## iterations: 100
## Time difference of 29.52099 secs

```

For $m = 1$, I found $\eta = .0009$ and $\alpha = .086$ to yield reasonable convergence.

```

m = 1
location = 1
beta.0 = rep(0, 8)
beta.1 = beta.0
alpha = .086
eta = .0009
niter = 0
indices = seq(1:100000)
beta.matrix = matrix(0, ncol = 8, nrow = 30000)

start <- Sys.time()
while(niter < 30000){
  beta.0 = beta.1
  samp = sample(indices, m)
  indices = indices[!indices %in% samp]
  total = sum(x.scaled[samp, ]*beta.0)
  d.term1 = exp(dnorm(total, log = TRUE) - pnorm(total, log = TRUE))
  d.term2 = exp(dnorm(total, log = TRUE) - pnorm(-total, log = TRUE))

  beta.1 = beta.0 + alpha/(1+eta*niter) * sum((y[samp]*d.term1 - (1-y[samp])*d.term2))*x.scaled[s:

  location = 1
  beta.scaled = .5*beta.1
  beta.scaled[1] = beta.1[1]

```



```

        beta.matrix[niter+1, ] = beta.scaled
        niter = niter + 1
    }
    test = TRUE
    end <- Sys.time()
    if(test){
        cat('beta0: ', beta.1[1], '\n')
        cat('beta1: ', beta.1[2], '\n')
        cat('beta2: ', beta.1[3], '\n')
        cat('beta3: ', beta.1[4], '\n')
        cat('beta4: ', beta.1[5], '\n')
        cat('beta5: ', beta.1[6], '\n')
        cat('beta6: ', beta.1[7], '\n')
        cat('beta7: ', beta.1[8], '\n')
        cat('iterations: ', niter, '\n')
        print(end - start)
        beta.m1 = beta.1
        time.m1 = end-start
    }

```

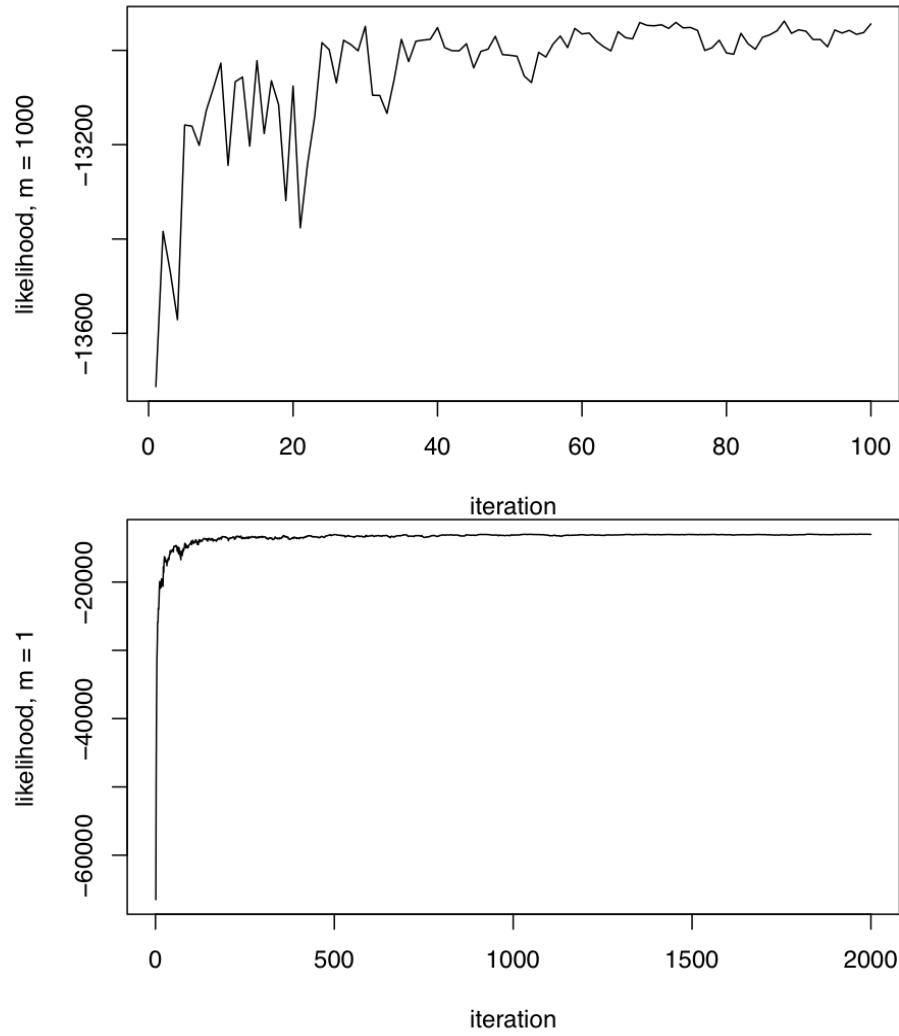
```

## beta0:  0.132605
## beta1:  4.56089
## beta2: -0.05280738
## beta3: -0.9442366
## beta4:  0.01761679
## beta5:  0.02645344
## beta6:  0.02807576
## beta7:  2.221518
## iterations:  30000
## Time difference of 2.056379 mins

```

Here are plots of the log likelihood of $\beta^{(t)}$ at each point in the stochastic gradient algorithms. In the $m = 1$ case I stored the likelihood at each 15th iteration to save computation time.

As expected, we see an initial increase as the algorithm moves towards an MLE, and stabilization after the neighborhood of this MLE is discovered and the number of iterations increases.



Tuning α and η was challenging because the stochastic gradient algorithm was sensitive to these parameters. Too large values for α led to movement away from the correct values of the parameters. I found small values for α and smaller decaying rate η yielded more accurate results.

Section 3 - Comparison of Results

Method	Execute Time	$ \beta - \hat{\beta} $
Newton-Raphson	12.65 secs	0.013
Coordinate Descent	2.81 mins	0.054
Stochastic Gradient $m = 1000$	29.52 secs	0.224
Stochastic Gradient $m = 1$	2.06 mins	0.41

The table above summarizes the results of the four algorithms. I used Euclidean distance between β , the vector used to generate the data, and $\hat{\beta}$, the MLE produced by the algorithm, to measure the accuracy of each method.

In this case Newton-Raphson was both the most accurate and speediest algorithm. This is due to p being small. Because Newton-Raphson requires the computation the inverse of a $p \times p$ matrix at each step, a growing p would significantly slow down this approach.

The stochastic gradient algorithms involve only the computation of the gradient. This computation is easily vectorized and therefore the time required to complete this computation will not increase quickly as p increases.

The coordinate descent algorithm can benefit from a stricter stopping condition. As the algorithm approaches the correct parameters, the magnitude of the steps each parameter takes diminishes rapidly. Thus the algorithm may still be approaching a better solution, but slowly enough that the relative change in the log likelihood is very small.

If the stopping condition is made to be a relative change in the log likelihood that is less than 10^{-7} , the coordinate descent algorithm yields an estimate much closer to the parameters used to generate the data:

```
## beta0: 0.1943842
## beta1: 2.404458
## beta2: -0.001431121
## beta3: -0.5038944
## beta4: -0.0007128089
## beta5: 0.008320743
## beta6: 0.002507493
## beta7: 1.198272
## iterations: 304
## Time difference of 4.812995 mins
```