

SDS 385: Exercises 2: Online Learning

September 13, 2016

Professor Scott

Spencer Woody

Problem 1

(A) From the previous exercise, we have the gradient of β

$$\nabla \ell(\beta) = \sum_{i=1}^N (m_i w_i - y_i) x_i. \quad (1)$$

We can think of $m_i w_i$ as the fitted value of y_i , or \hat{y}_i , when given β , so the gradient becomes

$$\nabla \ell(\beta) = \sum_{i=1}^N (\hat{y}_i - y_i) x_i \quad (2)$$

$$= \sum_{i=1}^N g_i(\beta) \quad (3)$$

$$g_i(\beta) = (\hat{y}_i - y_i) x_i. \quad (4)$$

(B)

$$E(n g_i(\beta)) = n E(g_i(\beta)) \quad (5)$$

In this expectation, the only random variable is i because the data X and y and the coefficients β are all fixed. The variable i is a random draw so it follows a discrete uniform distribution such that

$$P(i = j) = \begin{cases} \frac{1}{n} & j \in \{1, 2, \dots, n\} \\ 0 & \text{otherwise.} \end{cases} \quad (6)$$

Then we can compute the expectation.

$$E(g_i(\beta)) = \sum_{j=1}^n g_j(\beta) P(i = j) \quad (7)$$

$$= \sum_{j=1}^n g_j(\beta) \frac{1}{n} \quad (8)$$

$$= \frac{1}{n} \sum_{j=1}^n g_j(\beta) \quad (9)$$

$$= \frac{1}{n} \nabla \ell(\beta), \quad (10)$$

$$\Rightarrow E(n g_i(\beta)) = n E(g_i(\beta)) = n \frac{1}{n} \nabla \ell(\beta) = \nabla \ell(\beta) \quad (11)$$

(C) In this case, we perform SGD with 400,000 iterations of samples with replacement and a step size α of 0.0003. Our initial guess for β , $\hat{\beta}_0$ is some distance from the result of β from Newton's method. That distance is 5 plus some noise from the $Exp(1)$ distribution in either the positive or negative direction. Trace plots of all elements of $\hat{\beta}_k$ for each iteration are shown below in Figure 1. The red horizontal lines represent the values of $\hat{\beta}$ obtained from Newton's method. Also shown is trace plot of the exponential moving average (EMA) over the previous 500 iterations of the log-likelihood in Figure 2. We use the EMA to reduce the level of noise in our log-likelihood plot.

$$\hat{\beta}_{k+1} = \hat{\beta}_k - \alpha \nabla \ell(\hat{\beta}_k) \quad (12)$$

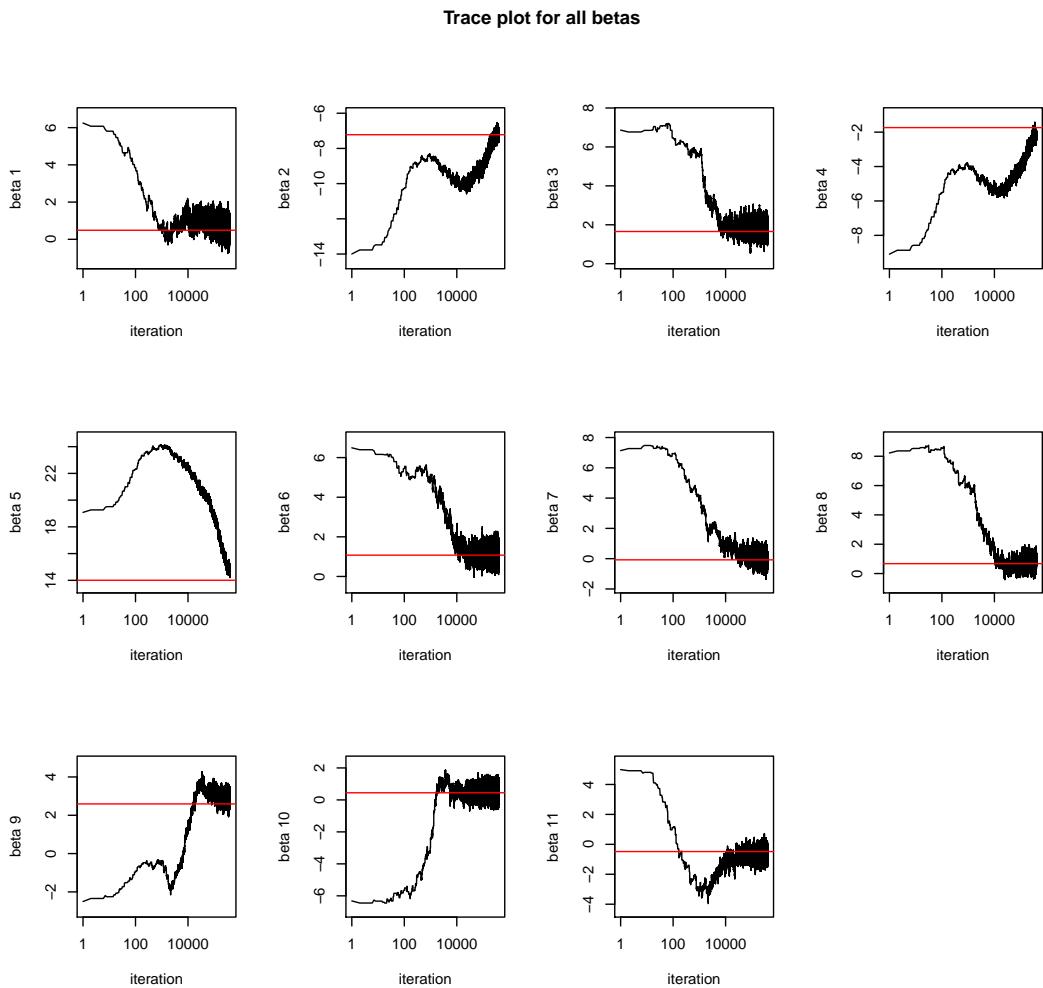


Figure 1: Trace plot of all elements of $\hat{\beta}$ from SGD

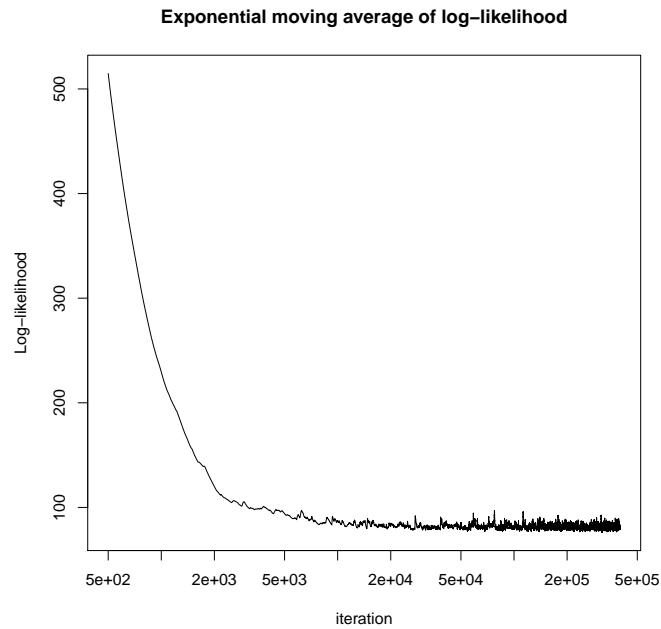


Figure 2: Trace plot of EMA of log-likelihood

(D)

(E)

```
#####
##### Created by Spencer Woody on 13 Sep 2016 #####
#####

5 library(TTR)

# Read in data file, scale X (y = 1 represents a malignant tumor)

wdbc <- read.csv("wdbc.csv", header = FALSE)

10 X <- as.matrix(wdbc[, 3:12])
X <- scale(X)
X <- cbind(rep(1, nrow(X)), X)

15 y <- wdbc[, 2]
y <- y == "M"
m.i <- 1

# # # # # # # # # # # # # # # #
# # # # # # # # # # # # # # # #
# # # # # # # # # # # # # # # #

# Function for computing w.i (logit transform of Xbeta)

25 comp.wi <- function (X, beta) {
  wi <- 1 / (1 + exp(-X %*% beta))
  return(wi)
}

30 # Function for computing likelihood

loglik <- function(beta, y, X, m.i) {
  loglik <- apply((m.i - y) * (X %*% beta)+ m.i*log(1 + exp(-X %*% beta)), 2, sum)
  return(loglik)
}

35 # Function for computing gradient of likelihood

grad.loglik <- function(beta, y, X, mi){
  grad <- array(NA, dim = length(beta))
  wi <- comp.wi(X, beta)
  grad <- apply(X*as.numeric(mi * wi - y), 2, sum)
  return(grad)
}

40 # # # # # # # # # # # # # # #
# # # # # # # # # # # # # # #
# # # # # # # # # # # # # # #

45 #### NEWTON'S METHOD (results are referred to as beta.N)

beta.N <- as.matrix(rep(0, ncol(X)))
```

```

newton.steps <- 10
55
for (step in 1:newton.steps) {
  Newton.wi <- as.numeric(comp.wi(X, beta.N))
  Hessian   <- t(X) %*% diag(Newton.wi*(1-Newton.wi)) %*% X
  beta.N    <- beta.N - solve(Hessian, grad.loglik(beta.N, y, X, m.i))
60
}

# Function for making traceplots of all betas

graph.betatrace <- function(beta, beta.N, graphname) {
  num.vars <- nrow(beta)
  n.graphrows <- floor(sqrt(num.vars))
  n.graphcols <- ceiling(num.vars / n.graphrows)
  pdf(graphname, width = n.graphcols * 2, height = n.graphrows * 2.5)
  par(mfrow = c(n.graphrows, n.graphcols), oma=c(0,0,2,0))
  70
  for (j in 1:num.vars) {
    plot(beta[j, ],
          xlab = "iteration",
          ylab = paste("beta", sprintf("%i", j)),
          type = "l",
          75
          ylim = c(min(beta.N[j], min(beta[j, ])) - 0.5,
          max(beta.N[j], max(beta[j, ])) + 0.5),
          log = "x"
          )
    abline(h = beta.N[j], col = "red")
  }
  title("Trace plot for all betas", outer = TRUE)
  dev.off()
}

85 # Stochastic gradient

SGD <- function(n.iterSGD, beta.init, step.size, X, y, m.i) {
  #' Perform stochastic gradient descent for a binomial logistic regression
  #
  #' @param n.iterSGD Number of iterations to loop through
  #' @param beta.init Initial guess for coefficients (betas)
  #' @param step.size Constant step size
  #' @param X N by P matrix of covariate data
  #' @param y P-vector of responses
  #' @param m.i n-parameter of binomial (1 for case of binary logistic)
  #
  #' @return A matrix, each column is an iteration of computed betas.
  #
  #
  100
  # Create zero matrix to store iterative values of beta; initialize beta
  betaSGD <- matrix(rep(0, ncol(X) * (n.iterSGD+ 1)), nrow = ncol(X))
  betaSGD[, 1] <- beta.init
  #
  for (step in 1:n.iterSGD) {
    #
    105
    # Draw random sample of single row of data (with replacement)
}

```

```

    i <- sample(nrow(X), 1)

    # Compute w.i, fitted value of p-parameter of binomial
    w.i <- 1 / (1 + exp(-crossprod(X[i, ], betaSGD[, step])))

    # Compute gradient, the descent direction
    grad <- nrow(X) * (m.i * w.i - y[i]) * X[i, ]

115   # Next set of betas
    betaSGD[, step + 1] <- betaSGD[, step] - step.size * grad
}
return(betaSGD)
}

120 # Create initial values of beta which are some distance away from "true" values
# of beta computed from Newton's method. At least 5 units away plus some
# exponential noise, in either positive or negative direction
beta1 <- beta.N + (-1) ^ rbinom(ncol(X), 1, 0.5) * (5 + rexp(ncol(X), rate = 1))

125 # Set number of iterations and stepsize
my.iter <- 4e5
my.stepsize <- 0.0003

130 # Perform SGD, graph trace plot of betas
sgd1 <- SGD(my.iter, beta1, my.stepsize, X, y, m.i)
graph.betatrace(sgd1, beta.N, "SGD.pdf")

135 # Compute loglikelihood from results of SGD
trace1 <- loglik(sgd1, y, X, m.i)

# Set time period for moving average
my.time <- 500

140 # Create plot of EMA of log-likelihood
pdf("EMALoglik.pdf")
plot(EMA(trace1, n = my.time),
      type = "l",
      main = "Exponential moving average of log-likelihood",
      xlab = "iteration",
      ylab = "Log-likelihood",
      xlim = c(my.time, length(trace1)),
      log = "x")
dev.off()

145 # Decaying steps

150 # Decaying steps

SGD.decay <- function(n.iterSGD, beta.init, C, t.0, alpha, X, y, m.i) {
  #' Perform stochastic gradient descent for a binomial logistic regression
  #
  #' @param n.iterSGD Number of iterations to loop through
  #' @param beta.init Initial guess for coefficients
  #' @param C Constant C in Robbins-Monro rule
}

```

```

160  #' @param t.0 Constant t.0 in Robbins-Monro rule
161  #' @param alpha Constant alpha in Robbins-Monro rule
162  #' @param X N by P matrix of covariate data
163  #' @param y P-vector of responses
164  #' @param m.i n-parameter of binomial (1 for case of binary logistic)
165  #
166  #' @return A matrix, each column is an iteration of computed betas.
167  #
168  #
169  # Create NULL matrix to store iterative values of beta; initialize beta
170  betaSGD.decay <- matrix(rep(0, ncol(X) * (n.iterSGD+ 1)), nrow = ncol(X))
171  betaSGD.decay[, 1] <- beta.init
172  #
173  for (step.decay in 1:n.iterSGD) {
174
175      # Draw random sample of single row of data (with replacement)
176      i <- sample(nrow(X), 1)
177
178      # Compute w.i, fitted value of p-parameter of binomial
179      w.i <- 1 / (1 + exp(-crossprod(X[i, ], betaSGD.decay[, step.decay])))
180
181      # Compute gradient, the descent direction
182      grad <- nrow(X) * (m.i * w.i - y[i]) * X[i, ]
183
184      # Compute next step size
185      stepsize.decay <- C * (step.decay + t.0) ^ -alpha
186
187      # Next set of betas
188      betaSGD.decay[, step.decay + 1] <- betaSGD.decay[, step] - stepsize.decay * grad
189  }
190
191  return(betaSGD.decay)
192 }

beta1 <- beta.N + (-1) ^ rbinom(ncol(X), 1, 0.5) * (3 + rexp(ncol(X), rate = 1))

195 decay1 <- SGD.decay(1e5, beta1, C = 10, t.0 = 1, alpha = 0.75, X, y, m.i)
graph.betatrace(decay1[, 10000:1e5], beta.N, "decay1.pdf")

# Running average of beta with decaying steps

200 burn.in <- 3e5
beta.burn <- sgd1[, burn.in:4e5]
beta.burnbar <- apply(beta.burn, 1, mean)

while (TRUE) {
205     if (A & A | j > 1000) {
          break
        }
}

```