# STAT40780 Data Programming with C (online)

## Lab Sheet 10

Dr Marie Galligan

Summer 2015

This week's lab involves learning about the gradient descent algorithm, and getting more practice with Rcpp, and Rcpp sugar.

## 1 Gradient descent algorithm

This task involves implementing a gradient descent algorithm in C++ with the help of Rcpp. But first, some background...

Working with data often involves minimization or maximization of functions with respect to one or more parameters. Consider the following convex function with a single parameter $\theta$:

$$f(\theta) = 1 + 3(\theta + 3)^2$$

A gradient descent algorithm can be used to find the value of $\theta$ that minimizes $f(\theta)$. To implement gradient descent, you need to know the derivative of the function you are trying to minimize, with respect to each parameter. In this case, there is a single parameter $\theta$ and hence only a single derivative is required:
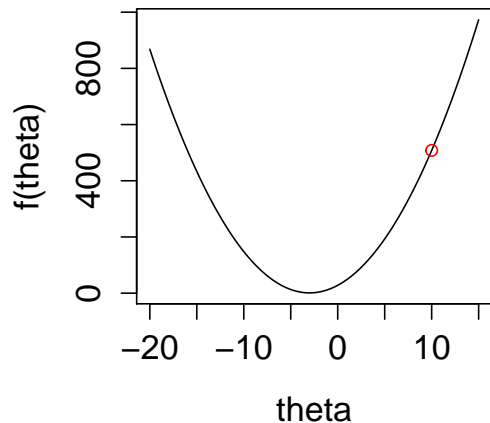
$$\frac{d}{d\theta}f(\theta) = 6(\theta + 3)$$

Remember that the derivative of this function is equal to its slope at a value $\theta$. If the derivative is positive, it means that the function is increasing at theta, while a negative derivative indicates that the function is decreasing at theta.

The gradient descent algorithm for a single parameter works as follows:

1. Choose an initial value for $\theta$ (initial guess)

2. Calculate the derivative of the function at the current 'guess' for $\theta$

3. Update the current 'guess' for theta by taking a step in the direction that decreases $f(\theta)$

Figure 1: Plot of $f(\theta)$. Red dot shows a current 'guess' of $\theta = 10$, where the objective is to find $\theta$ to minimize $f(\theta)$.



4. Repeat steps 2 and 3 until convergence to the minimum

At step 3, the current guess for $\theta$ is updated by taking a step either to the right or left of the current value of theta. Whether the value of $\theta$ is moved to the right or to the left depends on the derivative of $\theta$ at that point.

- if the derivative is positive at theta, $f(\theta)$ is increasing with $\theta$, hence decreasing the value of $\theta$ will decrease $f(\theta)$

- if the derivative is negative at theta, $f(\theta)$ is decreasing with $\theta$, hence increasing the value of $\theta$ will decrease $f(\theta)$

Suppose the current guess for $\theta = 10$ (see red dot in Figure 1). At this point, the function $f(\theta)$ is increasing with theta, and this will be reflected in the derivative (derivative will be positive at $\theta = 10$. Hence, to reach the minimum, the update in gradient descent will reduce the value of $\theta$ by a little bit.

At iteration $j$ of the gradient descent algorithm, denote the previous guess for $\theta$ as $\theta^{j-1}$. The rule for updating $\theta^{j-1}$ to a new guess $\theta^j$ is:

$$\theta^j = \theta^{j-1} - \alpha \frac{d}{d\theta} f(\theta^{j-1})$$

This update moves $\theta$ in the direction that decreases $f(\theta)$. Here, $\alpha$ is called the "learning rate" and determines the size of the step taken in gradient descent. A larger value for $\alpha$ means that the algorithm takes larger steps.

The following R code implements gradient descent to find the minimum of the function $f(\theta)$ defined above. The algorithm plots the current guess for $\theta$ at each iteration, and operates in slow motion to allow you to see individual updates being added to the plot.

Convergence of the algorithm is assessed at each iteration $j$ by monitoring the absolute relative change in $f(\theta)$:

$$rel\_ch = \left| \frac{f(\theta^j) - f(\theta^{j-1})}{f(\theta^{j-1})} \right|$$

If this relative change becomes really small (less than some specified level of tolerance - defined here as $tol$), the algorithm is assumed to have reached convergence.

Copy the code into an R script file and run. Be careful, as copy and paste moves the code around a little - you will probably need to move some things around! Experiment with different 'initial guess' values for theta, and different values for $\alpha$.

**gradient descent algorithm**

```
1
2
3   #FIRST, DEFINE THE FUNCTION TO BE MINIMIZED f(theta)
4   objfun <- function(theta){
5     return(1 + 3*(theta + 3)^2 )
6   }
7
8   #DEFINE A FUNCTION THAT CALCULATES THE DERIVATIVE f'(theta)
9   # AT INPUT VALUE theta
10  deriv <- function(theta)
11  {
12    return( 6*(theta + 3 ) )
13  }
14
15  #PLOT THE FUNCTION WE ARE SEEKING TO MINIMIZE OVER THE RANGE [-15, 15]
16  theta_seq <- seq(-15,15, len = 1000)
17  plot(theta_seq, objfun(theta_seq), type = "l",
18       ylab = "f(theta)", xlab = "theta")
19
20   tol <- 0.0000001 #CONVERGENCE THRESHOLD
21   alpha <- 0.01  #LEARNING RATE
22   theta0 <- 10  #SELECT INITIAL GUESS FOR THETA
23   newval <- objfun( theta0 ) #inital value of f(theta)
24   points(theta0, newval, col = "red") #add current theta to plot
25   rel_ch <- 1 #initialize relative change marker to any value larger than tol
26   j <- 1 #iteration counter
27   theta <- c() #vector to store theta at each iteration
28   theta[j] <- theta0 #theta[j] stores current value of theta
29
30   #update theta while relative change in f(theta) is greater than tol
31   while( rel_ch > tol )
32   {
33    j <- j+1; #increment j
34    #update theta
35    #set theta_new =  theta_previous - ( learning rate )* f'(theta_previous)
36    theta[j] <- theta[j-1] - alpha * deriv(theta[ j-1 ])
37
38    #test relative absolute change in target function
39    oldval <- newval #store f(theta_previous)
40    newval <- objfun(theta[j]) #calculate f(theta_new)
41    points(theta[j], newval, col = "red") #add new theta to plot
42    Sys.sleep(0.1) #pause algorithm to give you time to see dot appear on plot
43    #calculate relative change in f(theta)
44    rel_ch <- abs(  ( newval - oldval ) / oldval ) #use to test convergence
45   }
```

Be sure that you understand the gradient descent algorithm. The following resources give a more detailed explanation:

http://bayen.eecs.berkeley.edu/bayen/?q=webfm_send/262 https://www.youtube.com/watch?v=Fn8qXpIcdnI

and there are many others!

**TASK: Implement this gradient descent algorithm in C++ with help from Rcpp.**

# 2   Choosing the learning rate $\alpha$

In the gradient descent algorithm implemented in the previous task, a fixed learning rate was chosen e.g. $\alpha = 0.01$. Choosing $\alpha$ can have a large impact on

the algorithm. If $\alpha$ is too large, the function may not converge, instead bouncing around, missing the minimum. If $\alpha$ is too small, the algorithm can be slow to converge. See the following video for details:

http://www.dailymotion.com/video/x2clpb9_4-4-machine-learning-gradient-descent-in-practice-ii-learning-rate_school

From the video, you can see that if the learning rate is too large, the function you are trying to minimize could actually increase. In addition, ideally, the algorithm should take large steps when far away from the minimum value to move faster towards it, and take smaller steps as it approaches the minimum (to be sure it doesn't miss it).

The learning rate can be adapted as the algorithm progresses to ensure that this happens, using the Bold Driver method. At each iteration, check whether the target function $f(\alpha)$ has decreased since the previous iteration i.e. check whether

$$f(\theta^j) - f(\theta^{j-1}) < 0$$

If the the target function HAS decreased (i.e. got closer to minimum), then increase the learning rate a little on the next iteration by 5% so that larger steps are taken. However, if the target function has actually increased (BAD!), the algorithm has missed the minimum, so the learning rate should be decreased drastically by 50% on the next iteration (to ensure smaller steps are taken).

**Task: implement the Bold Driver method in your C++ gradient descent function to adapt the learning rate at each iteration.**

# 3   Rcpp sugar

1. Write a C++ function that calculates the min, max, mean and standard deviation (using Rcpp sugar functions) of a numeric vector passed from R, and returns the results to R in a named vector (elements named "min", "max", "mean" and "sd"). Pass this function a vector with some missing elements and observe its behaviour.

2. Rewrite the function to handle missing values.