# Lab 1

*August 11, 2016*

**Part 1.** We will work through some details on logistic regression. We first set some notation. Let $x_i \in \mathbb{R}^p, y_i \in \{0, 1\}$ for $i = 1, \ldots, n$. Let $X \in \mathbb{R}^{n \times p}$ denote the matrix with $x_i$ as its $i$th row. Recall that the negative log-likelihood $\ell(\beta)$ can be written as follows:

$$\ell(\beta) = \sum_{i=1}^{n} \left[ -y_i x_i^T \beta + \log(1 + \exp(x_i^T \beta)) \right].$$

1. Write the gradient and Hessian of $\ell(\beta)$.

2. What is the computational complexity for a calculating the gradient and Hessian of $\ell(\beta)$?

3. Under what condition is $\ell(\beta)$ strictly convex?

4. Prove that $\ell(\beta)$ is $L$-Lipschitz differentiable with $L = \frac{1}{4}\|X\|_{\text{op}}^2$.

5. Suppose that there is a vector $w \in \mathbb{R}^p$ such that $x_i^T w > 0$ if $y_i = 1$ and $x_i^T w < 0$ if $y_i = 0$. Prove that $\ell(\beta)$ does not have a global minimum. In other words, when the classification problem is completely separable, there is no maximum likelihood estimator.

**Part 2.** Gradient Descent

You will now implementation of gradient descent to your R package. Your function will include using both a fixed step-size as well as one chosen by backtracking.

Please complete the following steps.

**Step 1:** Write a function "gradient_step."

```
#' Gradient Step
#'
#' @param gradf handle to function that returns gradient of objective function
#' @param x current parameter estimate
#' @param t step-size
gradient_step <- function(gradf, t, x) {

}
```

Your function should return $x^+ = x - t\nabla f(x)$.

**Step 2:** Write a function "gradient_descent_fixed." The function should terminate either when a maximum number of iterations has been taken or if the relative change in the objective function has fallen below a tolerance

$$\frac{|f(x_k) - f(x_{k-1})|}{|f(x_{k-1})| + 1)} < \text{tolerance}$$

```
#' Gradient Descent (Fixed Step-Size)
#'
#' @param fx handle to function that returns objective function values
#' @param gradf handle to function that returns gradient of objective function
#' @param x0 initial parameter estimate
#' @param t step-size
#' @param max_iter maximum number of iterations
#' @param tol convergence tolerance
gradient_descent_fixed <- function(fx, gradf, t, x0, max_iter=1e2, tol=1e-3) {

}
```

Your function should return

- The final iterate value
- The objective function values
- The 2-norm of the gradient values
- The relative change in the function values
- The relative change in the iterate values

**Step 3:** Write a function "backtrack."

```
#' Backtracking
#'
#' @param fx handle to function that returns objective function values
#' @param x current parameter estimate
#' @param t current step-size
#' @param df the value of the gradient of objective function evaluated at the current x
#' @param alpha the backtracking parameter
#' @param beta the decrementing multiplier
backtrack <- function(fx, t, x, df, alpha=0.5, beta=0.9) {


}
```

Your function should return the selected step-size.

**Step 4:** Write a function "gradient_descent_backtrack" that performs gradient descent using backtracking.

```
#' Gradient Descent (Backtracking Step-Size)
#'
#' @param fx handle to function that returns objective function values
#' @param gradf handle to function that returns gradient of objective function
#' @param x0 initial parameter estimate
#' @param max_iter maximum number of iterations
#' @param tol convergence tolerance
gradient_descent_backtrack <- function(fx, gradf, x0, max_iter=1e2, tol=1e-3) {


}
```

Your function should return

- The final iterate value
- The objective function values
- The 2-norm of the gradient values
- The relative change in the function values
- The relative change in the iterate values

**Step 5:** Write a function "gradient_descent" that is a wrapper function for "gradient_descent_fixed" and "gradient_descent_backtrack." The default should be to use the backtracking.

```
#' Gradient Descent
#'
#' @param fx handle to function that returns objective function values
#' @param gradf handle to function that returns gradient of objective function
#' @param x0 initial parameter estimate
#' @param t step-size
#' @param max_iter maximum number of iterations
#' @param tol convergence tolerance
gradient_descent <- function(fx, gradf, x0, t=NULL, max_iter=1e2, tol=1e-3) {


}
```

Your function should return

- The final iterate value
- The objective function values
- The 2-norm of the gradient values
- The relative change in the function values
- The relative change in the iterate values

**Step 6:** Write functions 'fx_logistic' and 'gradf_logistic' to perform ridge logistic regression

```
#' Objective Function for Logistic Regression
#'
#' @param y binary response
#' @param X design matrix
#' @param beta regression coefficient vector
#' @param lambda regularization parameter
fx_logistic <- function(y, X, beta, lambda=0) {


}

#' Gradient for Logistic Regession
#'
#' @param y binary response
#' @param X design matrix
#' @param beta regression coefficient vector
#' @param lambda regularization parameter
gradf_logistic <- function(y, X, beta, lambda=0) {


}
```
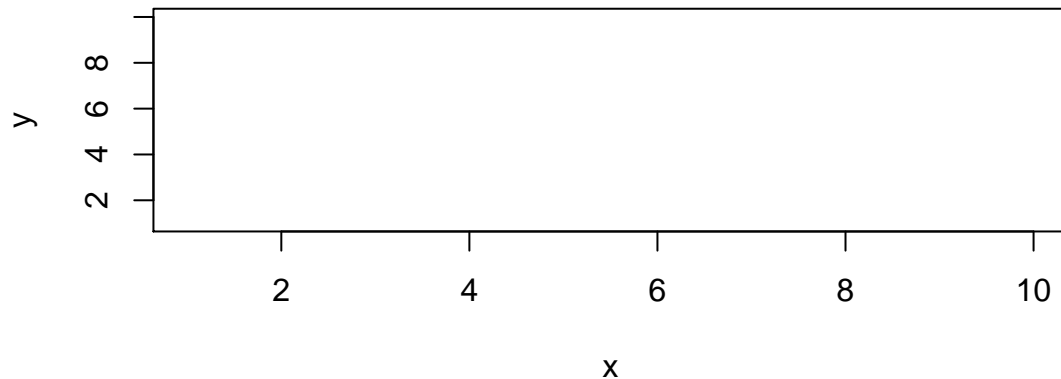
**Step 7:** Perform logistic regression (with $\lambda = 0$) on the following data example $(y, X)$ using the fixed step-size. Use your answers to Part 1 to choose an appropriate fixed step-size. Plot the difference $\ell(\beta_k) - \ell(\beta_{10000})$ versus the iteration $k$. Comment on the shape of the plot given what you know about the iteration complexity of gradient descent with a fixed step size.
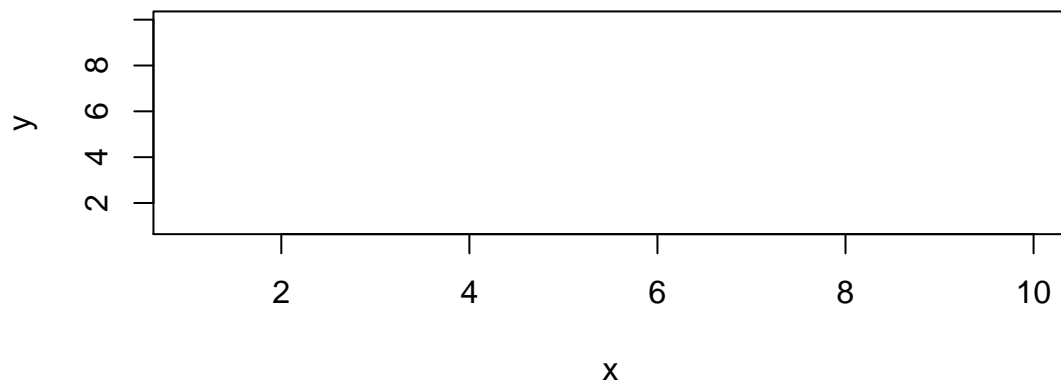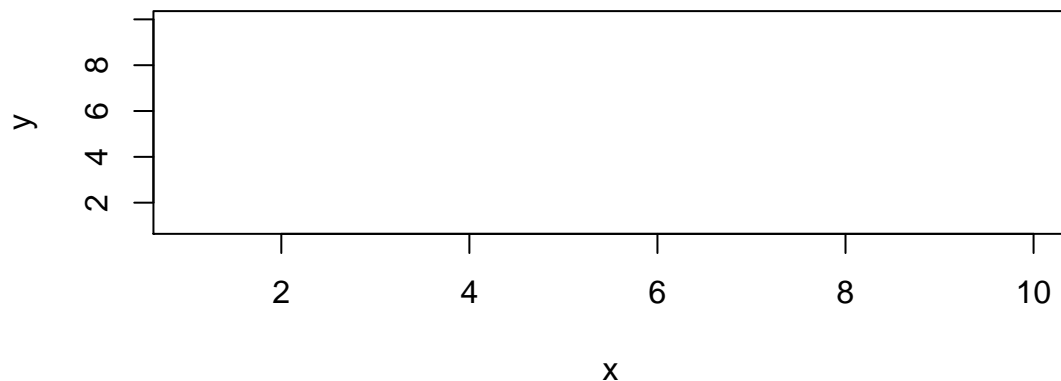
```
set.seed(12345)
n <- 100
p <- 2
```

```
X <- matrix(rnorm(n*p),n,p)
beta0 <- matrix(rnorm(p),p,1)
y <- (runif(n) <= plogis(X%*%beta0)) + 0
```



**Step 8:** Perform logistic regression (with $\lambda = 0$) on the simulated data above using backtacking. Plot the difference $\ell(\beta_k) - \ell(\beta_{10000})$ versus the iteration $k$. Comment on the shape of the plot given what you know about the iteration complexity of gradient descent with backtracking.



**Step 9:** Perform logistic regression (with $\lambda = 10$) on the simulated data above using the fixed step-size. Plot the difference $\ell_\lambda(\beta_k) - \ell_\lambda(\beta_{10000})$ versus the iteration $k$. Comment on the shape of the plot given what you know about the iteration complexity of gradient descent applied to strongly convex functions.

**Part 3.** If you have time, try out higher order methods. Install the R package lbfgs. How does the run time compare with gradient descent? You might find the 'system.time()' command useful. Try implementing Newton's method as discussed in lecture. How does the run time compare with gradient descent? BFGS?