

Multivariate Optimization 1

Solution 1: Gradient Descent

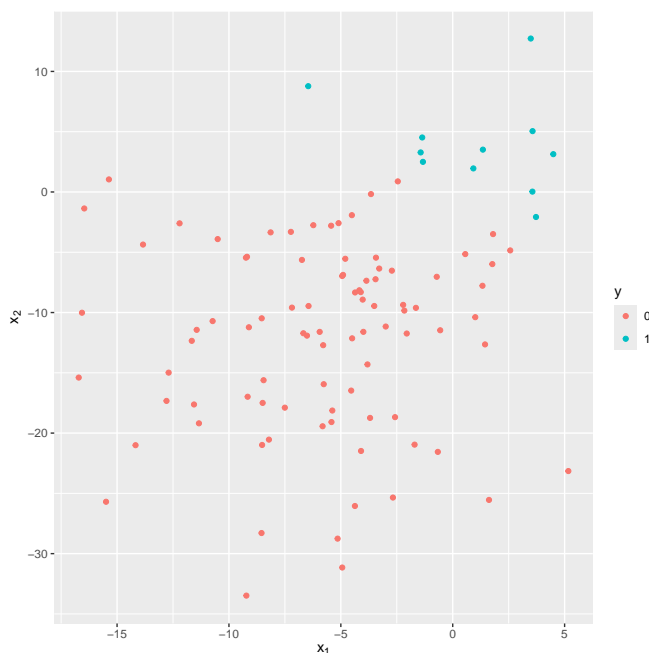
You are given the following data situation:

```
library(ggplot2)

set.seed(314)
n <- 100
X = cbind(rnorm(n, -5, 5),
          rnorm(n, -10, 10))
X_design = cbind(1, X)

z <- 2*X[,1] + 3*X[,2]
pr <- 1/(1+exp(-z))
y <- as.integer(pr > 0.5)
df <- data.frame(X = X, y = y)

ggplot(df) +
  geom_point(aes(x = X[,1], y = X[,2], color=as.factor(y))) +
  xlab(expression(x[1])) +
  ylab(expression(x[2])) +
  labs(colour = "y")
```



(a) We start with

$$\begin{aligned} \mathcal{R}_{\text{emp}}(\tilde{\boldsymbol{\theta}}) &= \sum_{i=1}^n \log(1 + \exp(\tilde{\boldsymbol{\theta}}^\top \mathbf{x}^{(i)})) - y^{(i)} \tilde{\boldsymbol{\theta}}^\top \mathbf{x}^{(i)} \\ &= \sum_{i=1}^n \begin{cases} \log(1 + \exp(\tilde{\boldsymbol{\theta}}^\top \mathbf{x}^{(i)})) & \text{if } y^{(i)} = 0, \\ \log(1 + \exp(\tilde{\boldsymbol{\theta}}^\top \mathbf{x}^{(i)})) - \tilde{\boldsymbol{\theta}}^\top \mathbf{x}^{(i)} & \text{if } y^{(i)} = 1. \end{cases} \end{aligned}$$

Since $\tilde{\theta}$ perfectly classifies the data, we know that

$$\begin{cases} \tilde{\theta}^\top \mathbf{x}^{(i)} < 0 & \text{if } y^{(i)} = 0, \\ \tilde{\theta}^\top \mathbf{x}^{(i)} \geq 0 & \text{if } y^{(i)} = 1. \end{cases}$$

Hence, we can focus on the functions $g(z) = \log(1 + \exp(-z))$ and $h(z) = \log(1 + \exp(z)) - z$ for $z > 0$ and study their monotonicity.

We compute

$$g'(z) = \frac{1}{1 + \exp(-z)} \cdot \exp(-z) \cdot (-1) = - \underbrace{\frac{\exp(-z)}{1 + \exp(-z)}}_{>0} < 0$$

and

$$h'(z) = \underbrace{\frac{\exp(z)}{1 + \exp(z)}}_{<1} - 1 < 0.$$

Therefore, both g and h are strictly monotonically decreasing.

It follows that $\mathcal{R}_{\text{emp}}(\tilde{\theta}) > \mathcal{R}_{\text{emp}}(\alpha \tilde{\theta})$ for $\alpha > 1$.

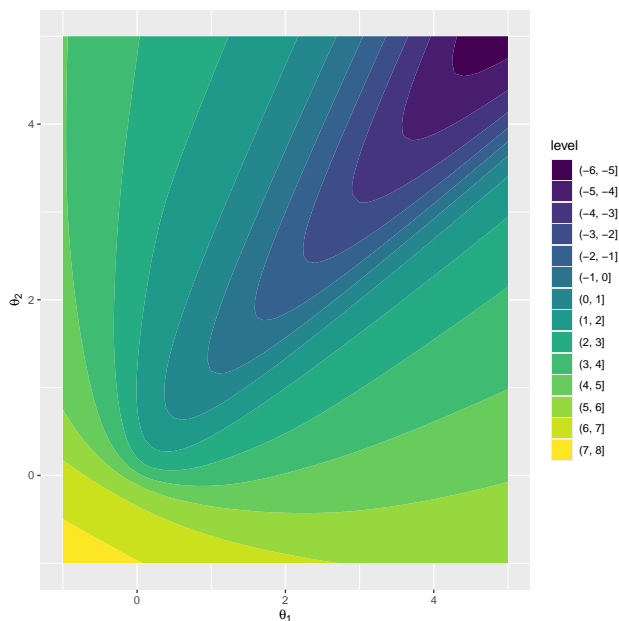
(b) `lambda = 0`

```
f <- function(theta, lambda) lambda * theta %*% theta +
  sum(-y * X %*% theta + log(1 + exp(X %*% theta)))

x = seq(-1, 5, by=0.1)
xx = expand.grid(X1 = x, X2 = x)

fxx = log(apply(xx, 1, function(t) f(t, lambda)))
df = data.frame(xx = xx, fxx = fxx)

ggplot() +
  geom_contour_filled(data = df, aes(x = xx.X1, y = xx.X2, z = fxx)) +
  xlab(expression(theta[1])) +
  ylab(expression(theta[2]))
```



(c) $\frac{\partial}{\partial \theta} \mathcal{R}_{\text{emp}} = \sum_{i=1}^n \frac{\exp(\theta^\top \mathbf{x}^{(i)})}{1 + \exp(\theta^\top \mathbf{x}^{(i)})} \mathbf{x}^{(i)\top} - y^{(i)} \mathbf{x}^{(i)\top}$

- (d) Note that we visualize from the first iteration on ($t = 1$) and not from the initial starting point $\theta^{[0]} = (0, 0)^\top$ but after having already made one GD step.

```
library(gridExtra)

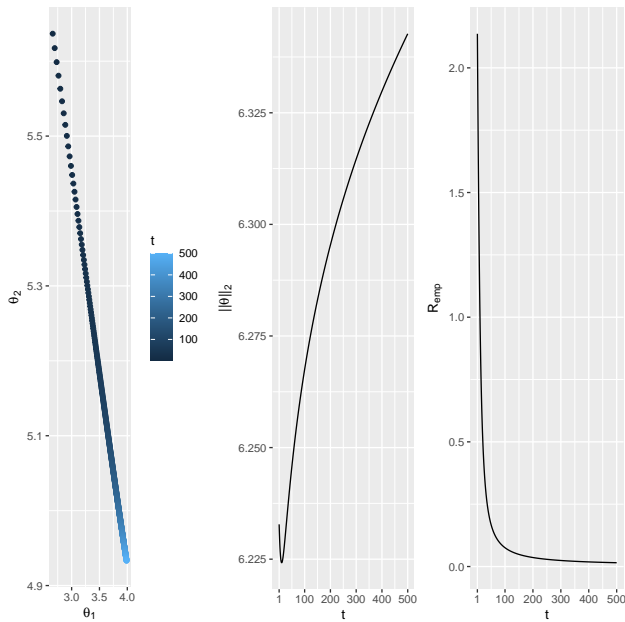
plot_fun <- function(gd_fun, lambda){
  theta = c(0,0)
  thetas = NULL
  thetas_norm = NULL
  fs = NULL
  for(i in 1:500){
    theta = gd_fun(theta)
    thetas_norm = rbind(thetas_norm, sqrt(theta %*% theta))
    thetas = rbind(thetas, t(theta))
    fs = rbind(fs, f(theta, lambda))
  }

  df_trace = as.data.frame(thetas)
  df_trace$t = 1:nrow(df_trace)
  trace_plot = ggplot() +
    geom_point(data = df_trace, aes(x=V1, y=V2, colour=t)) +
    xlab(expression(theta[1])) +
    ylab(expression(theta[2]))
  norm_plot = ggplot(data.frame(norms = thetas_norm, t = 1:nrow(thetas_norm))) +
    geom_line(aes(x = t, y = norms)) +
    scale_x_continuous(breaks = c(1, 100, 200, 300, 400, 500), limits = c(1, 500)) +
    ylab(expression(paste("||", theta, "||"[2])))
  remp_plot = ggplot(data.frame(f = fs, t = 1:nrow(thetas_norm))) +
    geom_line(aes(x = t, y = f)) +
    scale_x_continuous(breaks = c(1, 100, 200, 300, 400, 500), limits = c(1, 500)) +
    ylab(expression(R[emp]))
  grid.arrange(trace_plot, norm_plot, remp_plot, ncol=3)
}

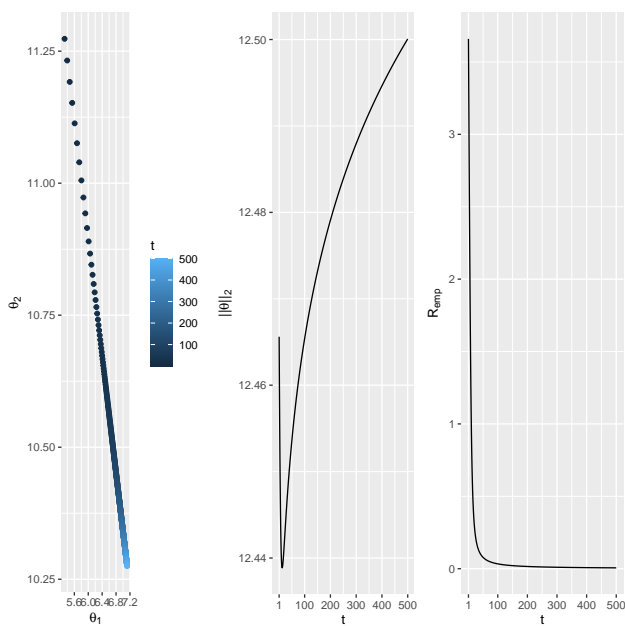
df_t <- function(theta, lambda) lambda * t(theta) -(t(y) %*% X) +
  t(1/(1 + exp(-X %*% theta))) %*% X

gd_step <- function(theta, alpha, lambda) return(theta - alpha * df_t(theta, lambda)[1,])

## Alpha = 0.01
gd_fun <- function(theta) return(gd_step(theta, 0.01, lambda))
plot_fun(gd_fun, 0)
```

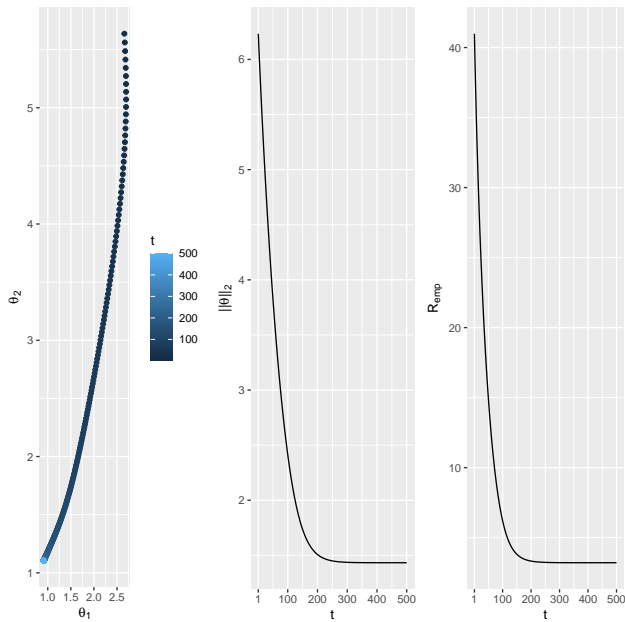


```
## Alpha = 0.02
gd_fun <- function(theta) return(gd_step(theta, 0.02, lambda))
plot_fun(gd_fun, 0)
```

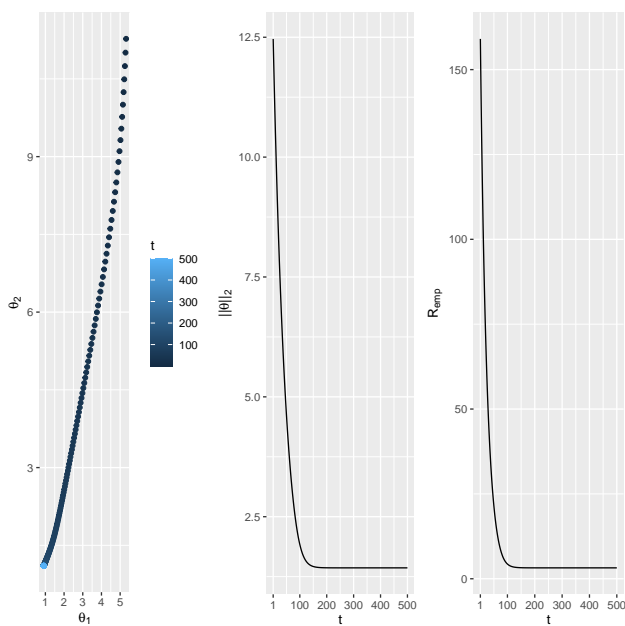


Gradient descent will in theory not converge since \mathcal{R}_{emp} has no minimum (a)

```
(e) ## Lambda = 1, alpha = 0.01
gd_fun <- function(theta) return(gd_step(theta, 0.01, 1))
plot_fun(gd_fun, 1)
```



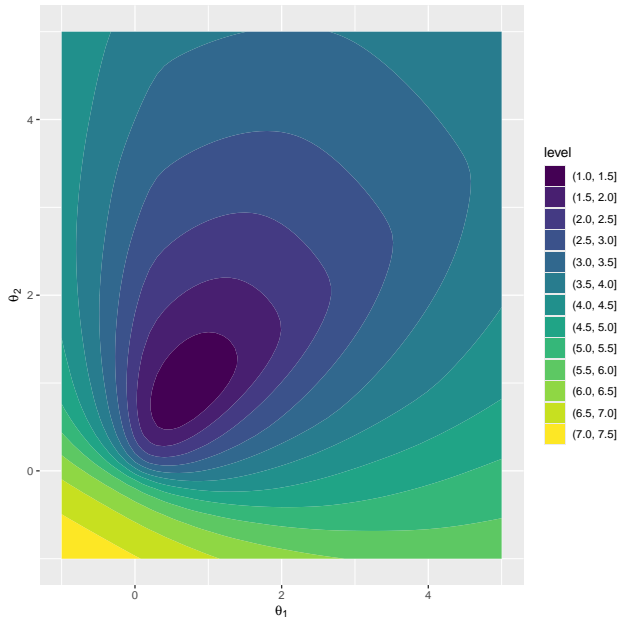
```
## Lambda = 1, alpha = 0.02
gd_fun <- function(theta) return(gd_step(theta, 0.02, 1))
plot_fun(gd_fun, 1)
```



(f) `lambda = 1`

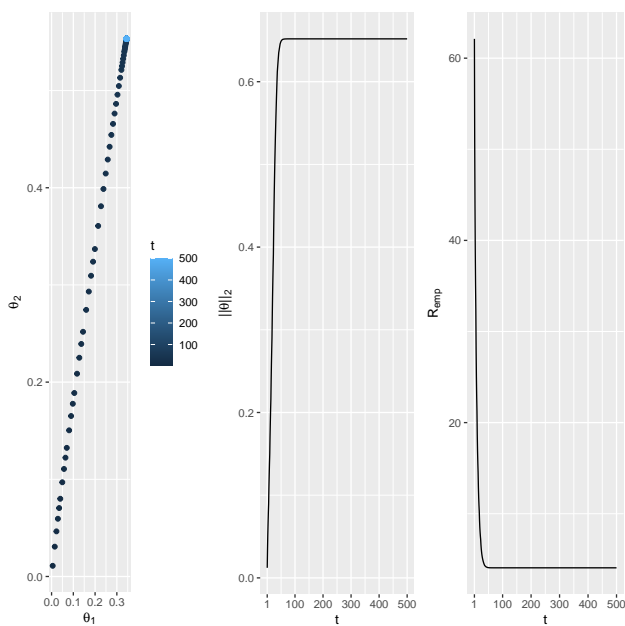
```
fxx_reg = log(apply(xx, 1, function(t) f(t, lambda)))
df_reg = data.frame(xx = xx, fxx = fxx_reg)

ggplot() +
  geom_contour_filled(data = df_reg, aes(x = xx.X1, y = xx.X2, z = fxx)) +
  xlab(expression(theta[1])) +
  ylab(expression(theta[2]))
```



```
(g) gd_backtracking_step <- function(theta, alpha, gamma, tau, lambda){
  ftheta = f(theta, lambda)
  dftheta = df_t(theta, lambda)[1,]
  for(i in 1:1000) {
    theta_prop = theta - alpha * dftheta
    if(f(theta_prop, lambda) <= ftheta - gamma * alpha * t(dftheta) %*% dftheta){
      return(theta_prop)
    }else{
      alpha = tau * alpha
    }
  }
  return(theta)
}

## Lambda = 1, alpha = 0.01
gd_fun <- function(theta) return(gd_backtracking_step(theta, 0.01, 0.9, 0.5, 1))
plot_fun(gd_fun, 1)
```



```
## Lambda = 1, alpha = 0.02
gd_fun <- function(theta) return(gd_backtracking_step(theta, 0.02, 0.9, 0.5, 1))
plot_fun(gd_fun, 1)
```

