

Questions 1-2 images

Quality of data imp of EDA	How would you eat yourself	Challenge either 1) resolved after assignment 2) hope to work on to improve your analytical skills (1 example)
-------------------------------	----------------------------------	---

— X —

Whether it's a convex -

a) Quadratic loss  $L = \sum_{i=1}^N e_i^2$

For proof of convexity: ST that 2nd derivative  
(or Hessian) of the loss function is non-negative  
for all values in the domain.

$$\rightarrow \frac{d}{de_i} [L] = \frac{d(\sum_{i=1}^N e_i^2)}{de_i} = 2e_i$$

$$\frac{d^2 L}{de_i^2} = \frac{d}{de_i} (e_1^2) + \frac{d}{de_i} (e_2^2) + \dots + \frac{d}{de_i} e_N^2$$

$$\therefore \frac{dL}{de_i} = 2e_i \quad \text{where } L = \sum_{i=1}^N e_i^2$$

$$\frac{d^2 L}{de_i^2} = 2 \quad \text{non-negative value for all domains}$$

$\therefore$  Quadratic loss is convex

$\rightarrow$  Since the second derivative is constant & positive, this confirms it is a convex function. This is good because they have a single global minima. This means that for linear regression, the optimisation will converge to global minima.

Questions 1-2 images

b) Mean absolute error

$$L = \sum_{i=1}^N |e_i|$$

Checking for convex function,

$$\frac{d(L)}{de_i} = \frac{d}{de_i} |e_i| = 1 \text{ for } e_i > 0$$

$$\frac{d^2 L}{de_i^2} = \cancel{\frac{d^2 |e_i|}{de_i^2}} \frac{d^2(1)}{de_i^2} = 0$$

$$\Rightarrow \frac{d(L)}{de_i} = \frac{d}{de_i} |e_i| = -1 \text{ for } e_i < 0$$

$$\frac{d^2 L}{de_i^2} = 0$$

function is non-negative for all values in domain. Hence convex.

→ What does this mean in context of log regression?  
This means that there is a global optima present & convergence is possible

We can also prove convexity of mean abs error by showing that any 2 points  $e_1, e_2$  in the domain of the function f any  $\lambda$  in the interval  $[0,1]$ , the following inequality holds.

$$f(\lambda e_1 + (1-\lambda)e_2) \leq \lambda f(e_1) + (1-\lambda) f(e_2) \text{ where } f(e) = |e|$$

Solving LHS

$$f(\lambda e_1 + (1-\lambda)e_2) = |\lambda e_1 + (1-\lambda)e_2|$$

Solving RHS

$$\lambda f(e_1) + (1-\lambda) f(e_2) = \lambda |e_1| + (1-\lambda) |e_2|$$

Applying 1<sup>st</sup> inequality rule  $|a+b| \leq |a| + |b|$

Questions 1-2 images

$$|\lambda e_1 + (1-\lambda) e_2| \leq \lambda |e_1| + (1-\lambda) |e_2|$$

$\therefore$  The inequality holds. Mean absolute error is hence convex.

c) Huber loss (smooth mean absolute error) with parameter  $s$

$$L = \sum_{i=1}^N l(e_i) \quad \text{where } l(e) = \begin{cases} \frac{1}{2} e^2 & \text{if } |e| \leq s \\ s|e| - \frac{1}{2} s^2 & \text{if } |e| > s \end{cases}$$

$\Rightarrow \frac{d(L)}{de^2}$  for case when  $|e| \leq s$

$$\Rightarrow \frac{d}{de} \left( \frac{1}{2} e^2 \right) = e \Rightarrow \frac{d^2}{de^2} \left( \frac{1}{2} e^2 \right) = 1$$

We can see that for  $|e| \leq s$   $L$  is a quadratic function. Hence it is convex.

$\Rightarrow \frac{d(L)}{de^2}$  for case when  $|e| > s$ , the function

$s|e| - \frac{1}{2} s^2$  is linear. We have already proved convexity for a linear function above.

Both the functions are convex individually

$\Rightarrow$  it remains under linear combination.

$\therefore$  Huber loss is convex, which means that its suitable for optimization algorithms that require a global minima.

Questions 1-2 images

2) For linear regression  $y_i = \theta_0 + \theta_1 x_i + e_i$ ,  $i=1,..N$   
 minimizing squared loss ( $e_i$  is random noise variable)

a) Calculate the gradient wrt parameter vector  $\theta$ :  
 Given the linear regression model:

$$y_i = \theta_0 + \theta_1 x_i + e_i$$

The squared loss func for linear regression  
 is typically defined as sum of squared errors.

$$L = \sum_{i=1}^N (y_i - (\theta_0 + \theta_1 x_i))^2$$

The parameter vector  $\theta = [\theta_0 \ \theta_1]^\top$

Calculating the gradient wrt  $\theta$  vector by  
 computing partial derivative wrt each param

$$\frac{\partial L}{\partial \theta_0} = -2 \sum_{i=1}^N (y_i - (\theta_0 + \theta_1 x_i))$$

$$\frac{\partial L}{\partial \theta_1} = -2 \sum_{i=1}^N x_i (y_i - (\theta_0 + \theta_1 x_i))$$

b) Steps of batch gradient descent rule.

→ Initialize parameter vector  $\theta$  to some initial values

→ Calculate the gradient [as shown in (a)] of the  
 loss function wrt  $\theta$

→ Calculating gradient tells us if the curve/slope  
 is decreasing/increasing. This helps us to  
 know if we are reaching global minima

→ Update parameters using gradient update rule

$$\theta_{\text{new}} = \theta_{\text{old}} - \alpha \frac{\partial L}{\partial \theta_{\text{old}}} \quad \alpha = \text{learning rate}$$

Questions 1-2 images

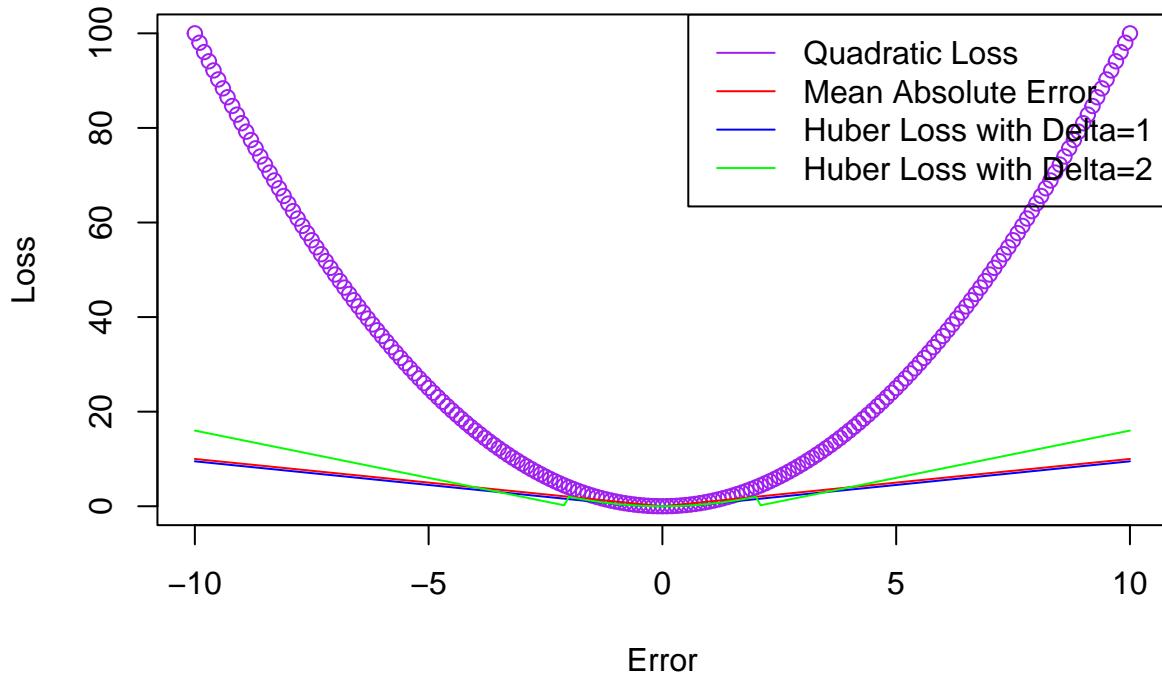
- Repeat this until global minima is reached [when change in  $\theta$  is very small or  $\theta \approx 0$ ]
- c) Steps of stochastic gradient descent rule
- Unlike backprop, this updates after each training example rather than entire dataset.
- Initialize parameter vector  $\theta$ .
- $\frac{\partial L}{\partial \theta_0} = -2(Y_i - (\theta_0 + \theta_1 x_i))$        $\frac{\partial L}{\partial \theta_1} = -2x_i(Y_i - (\theta_0 + \theta_1 x_i))$
- + Update parameter vector using rule  
 $\theta_{\text{new}} = \theta_{\text{old}} - \lambda(\theta_{\text{old}})$        $\lambda$  = learning rate
- Repeat the same until global minima is reached

## Assignment 02 R Notebook

Overlay graphs of the loss functions in question 1 for a range of  $\epsilon$  (consider two different values of  $\delta$  for Huber loss). Use the graph to discuss the relative advantages and disadvantages of these loss functions for linear regression

```
quadratic_loss <- function(err){  
  result <- err^2  
  return (result)  
}  
  
mean_abs_error <- function(err){  
  result <- abs(err)  
  return (result)  
}  
  
huber_loss <- function(err, delta){  
  ifelse(abs(err) <= delta, 0.5 * err^2, ifelse(abs(err) > delta, (delta * (abs(err) - 0.5*delta^2)),NA))  
}  
  
errs <- seq(-10,10,0.1)  
  
plot(errs, sapply(errs,quadratic_loss), col = 'purple',xlab = 'Error',ylab = 'Loss',main = 'Comparison of Loss Functions')  
lines(errs, sapply(errs,mean_abs_error), col = 'red')  
lines(errs,sapply(errs, huber_loss,delta=1), col = 'blue')  
lines(errs,sapply(errs, huber_loss,delta=2), col = 'green')  
legend('topright', legend = c('Quadratic Loss', 'Mean Absolute Error', 'Huber Loss with Delta=1','Huber Loss with Delta=2'))
```

## Comparison of Loss Function



Question 3b)

```

quadratic_GD <- function(x,y, learningRate, epochs) {
  n <- length(y) #no of features
  m <- length(x) #training rows
  theta <- rep(0,m) #should this be n?

  for (err in 1:epochs){
    gradient <- 2/n * t(x) %*% (x %*% theta - y)
    theta <- theta - learningRate * gradient #Updating theta
  }
  return (theta)
}

MAE_GD <- function(x,y, learningRate, epochs){
  n <- length(y) #no of features
  m <- length(x) #training rows
  theta <- rep(0,m) #should this be n?

  for (err in 1:epochs){
gradient <- 1/n * t(x) %*% sign(x %*% theta - y)
    theta <- theta - learningRate * gradient #Updating theta. This remains the same throughout.
  }
  return (theta)
}

huber_GD <- function(x,y, learningRate, epochs, delta){
```

```

n <- length(y) #no of features
m <- length(x) #training rows
theta <- rep(0,m) #should this be n?

for (err in 1:epochs){
  error <- x %*% theta - y
  gradient <- (1/n) * t(x) %*% ifelse(abs(error) <= delta, error, delta * sign(error))
  theta <- theta - learningRate * gradient #Updating theta. This remains the same throughout.
}
return (theta)
}

```

Question 3c)

```

quadratic_sto <- function(x,y,learningRate, epochs){
  n <- length(y) #no of features
  m <- length(x) #training rows
  theta <- rep(0,m) #should this be n?

  for (err in 1:epochs){
    for(i in 1:n){
      random <- sample(1:n,1)
      xi <- x[random,]
      yi <- y[random]
      gradient <- 2 * t(xi) %*% (xi %*% theta - yi)
      theta <- theta - learningRate * gradient
    }
  }
  return (theta)
}

MAE_sto <- function(x,y,learningRate, epochs){
  n <- length(y) #no of features
  m <- length(x) #training rows
  theta <- rep(0,m) #should this be n?

  for (err in 1:epochs){
    for(i in 1:n){
      random <- sample(1:n,1)
      xi <- x[random,]
      yi <- y[random]
      gradient <- 2 * t(xi) %*% sign(xi %*% theta - yi)
      theta <- theta - learningRate * gradient
    }
  }
  return (theta)
}

huber_sto <- function(x,y,learningRate, epochs){
  n <- length(y) #no of features
  m <- length(x) #training rows
  theta <- rep(0,m) #should this be n?
}

```

```

for (err in 1:epochs){
  for(i in 1:n){
    random <- sample(1:n,1)
    xi <- x[random,]
    yi <- y[random]
    gradient <- (1/n) * t(x) %*% ifelse(abs(error) <= delta, error, delta * sign(error))
    theta <- theta - learningRate * gradient
  }
}
return (theta)
}

```

Question 4a)

```

#fabricate data
set.seed(0)
N <- 50
x <- runif(N, -2, 2)
epsilon <- rnorm(N, 0, 2)
y <- 3 + 2 * x + epsilon

x <- cbind(1,x)
analytical <- solve(t(x) %*% x) %*% t(x) %*% y #Analytical Solution

batch_GD <- function(x, y, alpha = 0.01, iterations = 1000){
  m <- nrow(x)
  theta <- matrix(0, ncol(x), 1)
  for(i in 1:iterations){
    gradient <- t(x) %*% (x %*% theta - y) / m
    theta <- theta - alpha * gradient
  }
  return(theta)
}

batch <- batch_GD(x, y)

#SGD
sto_GD <- function(X, y, learningRate = 0.01, iterations = 50 * N) {
  m <- nrow(X)
  theta <- matrix(0, ncol(X), 1) # Initial parameters
  for (i in 1:iterations) {
    idx <- sample(1:m, 1)
    X_i <- X[idx, , drop = FALSE] # Ensure X_i is always a matrix
    y_i <- y[idx]
    gradient <- t(X_i) %*% (X_i %*% theta - y_i)
    theta <- theta - learningRate * gradient
  }
  return(theta)
}

sto <- sto_GD(x, y)

```

```

print(analytical)

##      [,1]
## 2.931669
## x 1.656737

print(batch)

##      [,1]
## 2.931514
## x 1.656808

print(sto)

##      [,1]
## 2.840583
## x 1.693194

4b)

fit <- function() {
  x <- runif(N, -2, 2)
  epsilon <- rnorm(N, 0, sqrt(4))
  y <- 3 + 2*x + epsilon
  X <- cbind(1, x)

  # Analytical Soln
  analytical <- solve(t(X) %*% X) %*% t(X) %*% y

  # Batch Gradient Descent
  batch_GD <- batch_GD(X, y)

  # Stochastic Gradient Descent
  sto_GD <- sto_GD(X, y)

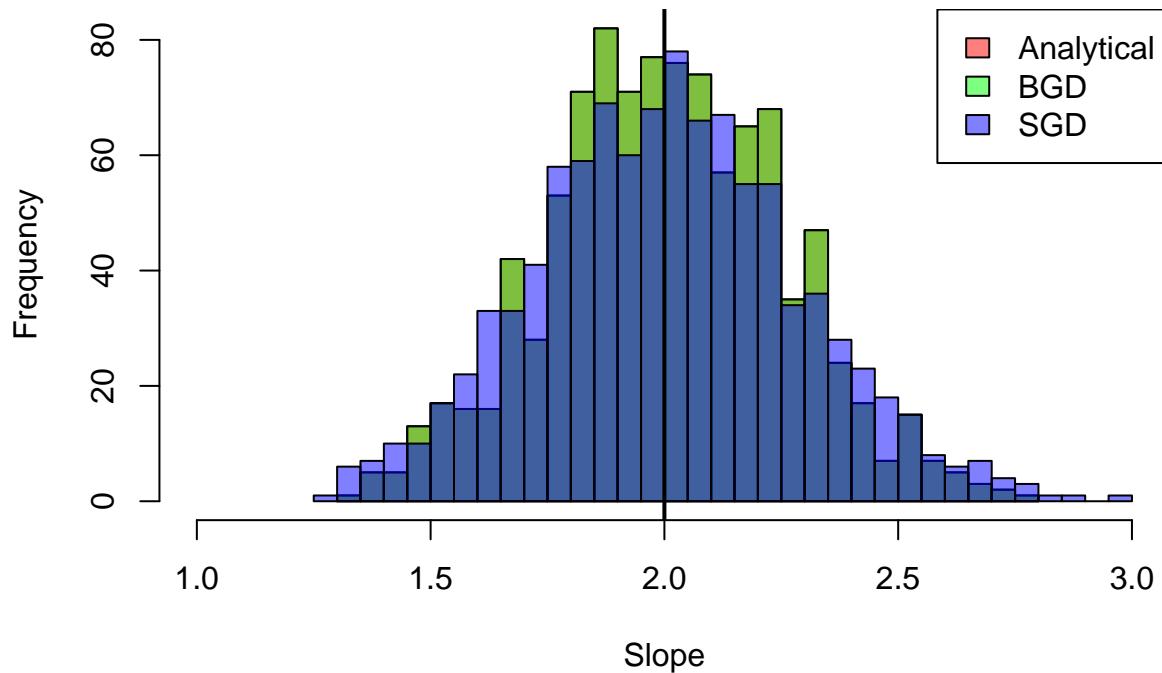
  c(analytical[2], batch_GD[2], sto_GD[2])
}

set.seed(0)
results <- replicate(1000, fit())

# Plot results
hist(results[1,], breaks = 30, col = rgb(1,0,0,0.5), xlim = c(1,3), main = "Slope Estimates", xlab = "Slope Estimate")
hist(results[2,], breaks = 30, col = rgb(0,1,0,0.5), add = TRUE)
hist(results[3,], breaks = 30, col = rgb(0,0,1,0.5), add = TRUE)
abline(v = 2, col = "black", lwd = 2) # True value
legend("topright", legend = c("Analytical", "BGD", "SGD"), fill = c(rgb(1,0,0,0.5), rgb(0,1,0,0.5), rgb(0,0,1,0.5)))

```

## Slope Estimates



Depending on the learning rate and number of epochs used, batch gradient descent and stochastic gradient descent may produce somewhat different estimates, even if the analytical solution offers objective and accurate estimates. SGD may offer improved generalization and faster convergence, however the estimates can vary somewhat. Batch gradient descent could yield more reliable estimates, but for best results, hyperparameters might need to be adjusted.

4c)

```
set.seed(0)
N <- 50
x <- runif(N, -2, 2)
epsilon <- rnorm(N, 0, sqrt(4))
y <- 3 + 2*x + epsilon
X <- cbind(1, x)

analytical <- solve(t(X) %*% X) %*% t(X) %*% y

gradient_l1_loss <- function(Y, X, 0) { # MAE with Batch Gradient Descent
  gradient <- numeric(length(0))
  Y_hat <- 0[1] + 0[2]*X
  gradient[1] <- sum(ifelse(Y > Y_hat, 1, ifelse(Y < Y_hat, -1, 0)))
  gradient[2] <- sum(ifelse(Y > Y_hat, X, ifelse(Y < Y_hat, -X, 0)))
  return(gradient)
}

l1_gd <- function(Y, X, 0, alpha = 0.01, num_iterations = 500) {
  for (i in 1:num_iterations) {
```

```

gradient <- gradient_l1_loss(Y, X, 0)
O <- O - alpha * gradient
}
return(0)
}

b_mae = c()

b_hat_mae <- l1_gd(X, y, O = c(0.5,0.5))

gradient_hubert_loss <- function(Y, X, O, delta){
  # Huber Loss with Batch Gradient Descent
  gradient <- numeric(length(O))
  Y_hat <- O[1] + O[2]*X
  is_small <- abs(Y - Y_hat) <= delta
  gradient[1] <- -X * ifelse(is_small, Y-Y_hat, delta * sign(Y-Y_hat))
  gradient[2] <- -1 * ifelse(is_small, Y-Y_hat, delta * sign(Y-Y_hat))
  return(gradient)
}

h_gd <- function(Y, X, O, alpha = 0.01, num_iterations = 500, delta = 1.5) {
  for (i in 1:num_iterations) {
    gradient <- gradient_hubert_loss(Y, X, O, delta)
    O <- O - alpha * gradient
  }
  return(O)
}

suppressWarnings(b_hat_hubert <- h_gd(X, y, O = c(0.5, 0.5), delta = 1.5))

print(analytical)

```

```

##      [,1]
## 2.931669
## x 1.656737

```

```
print(b_hat_mae)
```

```
## [1] 360.780 1731.747
```

```
print(b_hat_hubert)
```

```
## [1] -1.0451984 0.2848151
```

4d)

```

choice_loss <- function() {
  x <- runif(N, -2, 2)
  epsilon <- rnorm(N, 0, sqrt(4))
  y <- 3 + 2*x + epsilon
  X <- cbind(1, x)
}
```

```

beta_hat_analytical <- solve(t(X) %*% X) %*% t(X) %*% y

beta_hat_mae <- l1_gd(X, y, 0 = c(0.5, 0.5)) #MAE

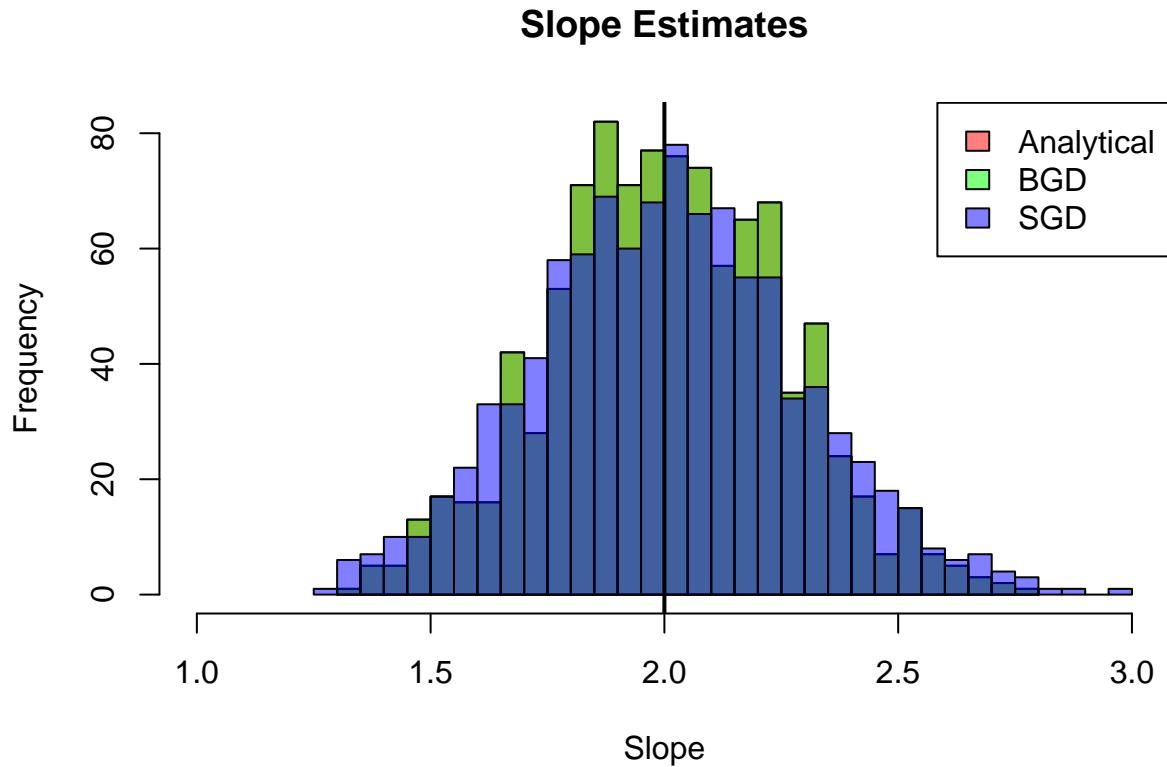
suppressWarnings(beta_hat_hub <- h_gd(X, y, 0 = c(0.5, 0.5), delta = 1.5)) #Huber

c(beta_hat_analytical[2], beta_hat_mae[2], beta_hat_hub[2])
}

set.seed(0)
results_part_c <- replicate(1000, choice_loss())

# Plotting hists
hist(results[1,], breaks = 30, col = rgb(1,0,0,0.5), xlim = c(1,3), main = "Slope Estimates", xlab = "Slope")
hist(results[2,], breaks = 30, col = rgb(0,1,0,0.5), add = TRUE)
hist(results[3,], breaks = 30, col = rgb(0,0,1,0.5), add = TRUE)
abline(v = 2, col = "black", lwd = 2) # True value
legend("topright", legend = c("Analytical", "BGD", "SGD"), fill = c(rgb(1,0,0,0.5), rgb(0,1,0,0.5), rgb(0,0,1,0.5)))

```



The loss function used influences the slope parameter estimates in linear regression by changing the trade-off between resilience to outliers and optimization stability. Squared loss is sensitive to outliers but produces a convex optimization issue, whereas absolute loss is resilient to outliers but may result in non-convex optimization. Huber loss strikes a balance between these two extremes, providing resilience to outliers while retaining convergence features.

4e)

```

outliers <- function(out_prop = 0.1) {
  x <- runif(N, -2, 2)
  epsilon <- rnorm(N, 0, sqrt(4))
  y <- 3 + 2*x + epsilon
  X <- cbind(1, x)

  #Outliers
  num_outliers <- ceiling(N * out_prop)
  outlier_indices <- sample(1:N, num_outliers, replace = FALSE)

  # Apply outlier effect
  for (i in outlier_indices) {
    # Randomly decide to add or subtract the magnitude
    if (runif(1) < 0.5) {
      y[i] <- 2 * y[i]
    } else {
      y[i] <- y[i] / 2
    }
  }

  #Analytical
  analytical <- solve(t(X) %*% X) %*% t(X) %*% y

  mae <- l1_gd(X, y, 0 = c(0.5, 0.5)) # MAE

  suppressWarnings(huber <- h_gd(X, y, 0 = c(0.5, 0.5), delta = 1.5)) #huber

  return(c(analytical[2], mae[2], huber[2]))
}

set.seed(0)
ans <- outliers()
print(ans)

```

```
## [1] 1.7316679 1832.7799157 0.2674075
```

4f)

```

outliers <- function(out_prop = 0.1) {
  x <- runif(N, -2, 2)
  epsilon <- rnorm(N, 0, sqrt(4))
  y <- 3 + 2*x + epsilon
  X <- cbind(1, x)

  #Outliers
  num_outliers <- ceiling(N * out_prop)
  outlier_indices <- sample(1:N, num_outliers, replace = FALSE)

  # Apply outlier effect
  for (i in outlier_indices) {
    # Randomly decide to add or subtract the magnitude
    if (runif(1) < 0.5) {
      y[i] <- 2 * y[i]
    }
  }
}
```

```

    } else {
      y[i] <- y[i] / 2
    }
}

# Squared Loss Analytical
beta_hat_analytical <- solve(t(X) %*% X) %*% t(X) %*% y

# MAE with Batch Gradient Descent
beta_hat_mae <- l1_gd(X, y, 0 = c(0.5, 0.5))

# Huber Loss with Batch Gradient Descent
suppressWarnings(beta_hat_huber <- h_gd(X, y, 0 = c(0.5, 0.5), delta = 1.5))

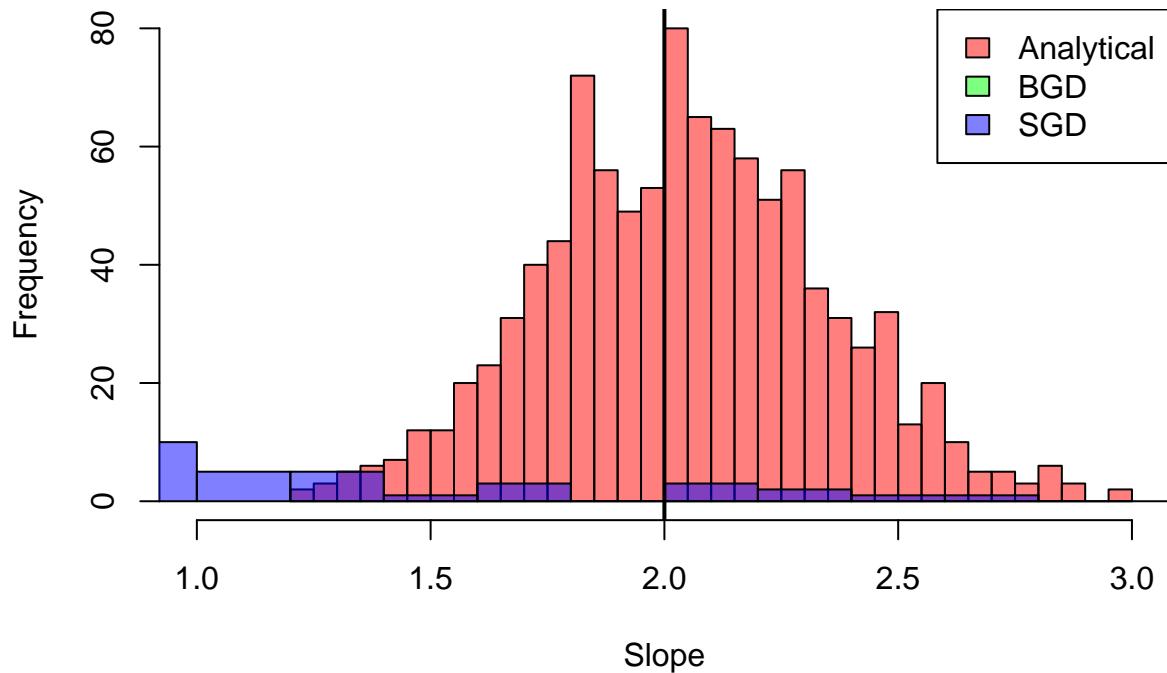
return(c(beta_hat_analytical[2], beta_hat_mae[2], beta_hat_huber[2]))
}

set.seed(0)
results <- replicate(1000, outliers())

# Plot results
hist(results[1,], breaks = 30, col = rgb(1,0,0,0.5), xlim = c(1,3), main = "Slope Estimates", xlab = "Slope Estimate")
hist(results[2,], breaks = 30, col = rgb(0,1,0,0.5), add = TRUE)
hist(results[3,], breaks = 30, col = rgb(0,0,1,0.5), add = TRUE)
abline(v = 2, col = "black", lwd = 2) # True value
legend("topright", legend = c("Analytical", "BGD", "SGD"), fill = c(rgb(1,0,0,0.5), rgb(0,1,0,0.5), rg

```

## Slope Estimates



Loss function selection in the presence of outliers influences the estimations of the slope parameter. Compared to mean absolute error and Huber loss, which handle outliers more leniently, squared loss is more resilient to outliers and may provide estimates that are skewed. When working with data that contains outliers, huber loss is a common option because it offers a flexible balance between the efficiency of squared loss and the robustness of absolute loss.