

Linear Regression and Regularization

Applied Machine Learning for Educational Data Science

true

10/12/2021

Contents

Linear Regression	2
Model Description	2
Model Estimation	3
Matrix Solution	11
Hat Matrix	16
lm() function	17
Model Interpretation	17
Model Evaluation	17
Linear Regression with Regularization	17
Ridge Penalty	17
Lasso Penalty	17
Elastic Net	17
Wrapping up	17
Building a Prediction Model for Readability Scores	17
Using the Prediction Model for a New Text	17

[Updated: Thu, Oct 14, 2021 - 17:22:29]

In the machine learning literature, the prediction algorithms are classified into two main categories: *supervised* and *unsupervised*. Supervised algorithms are being used when the dataset has an actual outcome of interest to predict (labels), and the goal is to build the “best” model predicting the outcome of interest that exists in the data. On the other side, unsupervised algorithms are being used when the dataset doesn’t have an outcome of interest, and the goal is typically to identify similar groups of observations (rows of data) or similar groups of variables (columns of data) in data. In this course, we plan to cover a number of *supervised* algorithms. Linear regression is one of the simplest approach among supervised algorithms, and also one of the easiest to interpret.

Linear Regression

Model Description

In most general terms, the linear regression model with P predictors ($X_1, X_2, X_3, \dots, X_p$) to predict an outcome (Y) can be written as the following:

$$Y = \beta_0 + \sum_{p=1}^P \beta_p X_p + \epsilon.$$

In this model, Y represents the observed value for the outcome for an observation, X_p represents the observed value of the p^{th} variable for the same observation, and β_p is the associated model parameter for the p^{th} variable. ϵ is the model error (residual) for the observation.

This model includes only the main effects of each predictor and can be easily extended by including a quadratic or higher-order polynomial terms for all (or a specific subset of) predictors. For instance, the model below includes all first-order, second-order, and third-order polynomial terms for all predictors.

$$Y = \beta_0 + \sum_{p=1}^P \beta_p X_p + \sum_{k=1}^P \beta_{k+P} X_k^2 + \sum_{m=1}^P \beta_{m+2P} X_m^3 + \epsilon.$$

The simple first-order, second-order, and third-order polynomial terms can also be replaced by corresponding terms obtained from B-splines or natural splines.

Sometimes, the effect of predictor variables on the outcome variable are not additive, and the effect of one predictor on the response variable can depend on the levels of another predictor. These non-additive effects are also called interaction effects. The interaction effects can also be a first-order interaction (interaction between two variables, e.g., $X_1 * X_2$), second-order interaction ($X_1 * X_2 * X_3$), or higher orders. It is also possible to add the interaction effects to the model. For instance, the model below also adds the first-order interactions.

$$Y = \beta_0 + \sum_{p=1}^P \beta_p X_p + \sum_{k=1}^P \beta_{k+P} X_k^2 + \sum_{m=1}^P \beta_{m+2P} X_m^3 + \sum_{i=1}^P \sum_{j=i+1}^P \beta_{i,j} X_i X_j + \epsilon.$$

If you are not comfortable or confused with notational representation, below is an example for different models you can write with 5 predictors (X_1, X_2, X_3).

A model with only main-effects:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \epsilon.$$

A model with polynomial terms up to the 3rd degree added:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_1^2 + \beta_5 X_2^2 + \beta_6 X_3^2 + \beta_7 X_1^3 + \beta_8 X_2^3 + \beta_9 X_3^3$$

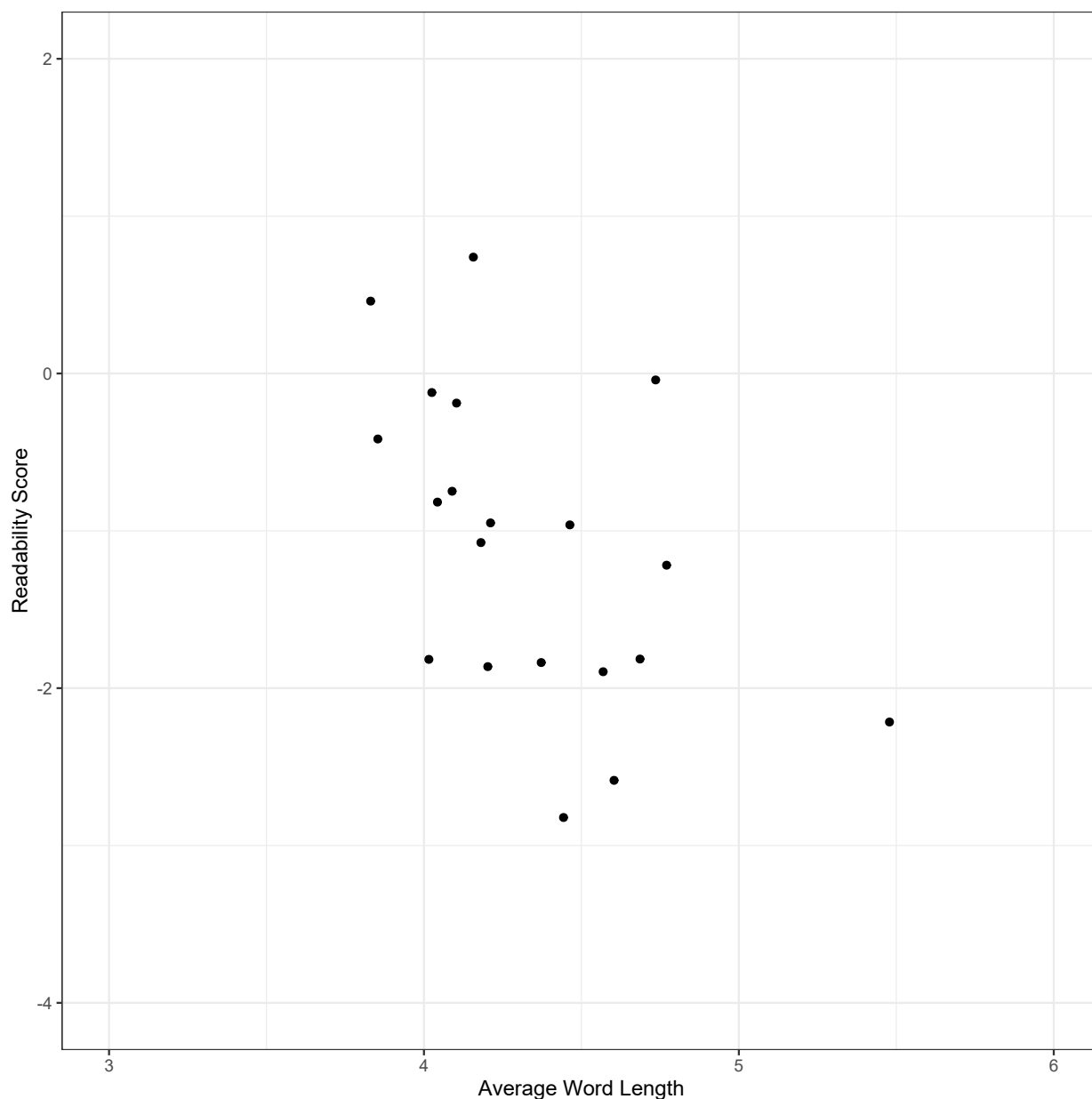
A model with both interaction terms and polynomial terms up to the 3rd degree added:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_1^2 + \beta_5 X_2^2 + \beta_6 X_3^2 + \beta_7 X_1^3 + \beta_8 X_2^3 + \beta_9 X_3^3 + \beta_{1,2} X_1 X_2 + \beta_{1,3} X_1 X_3 + \beta_{2,3} X_2 X_3 + \epsilon$$

Model Estimation

Suppose that we would like to predict the target readability score for a given text from average word length in the text. Below is a scatterplot to show the relationship between these two variables for a random sample of 20 observations. There seems to be a moderate negative correlation. So, we can tell that the higher the average word length is in a given text, the lower the readability score (more difficult to read).

```
readability_sub <- read.csv('https://raw.githubusercontent.com/uo-datasci-specialization/c4-ml-fall-202
```

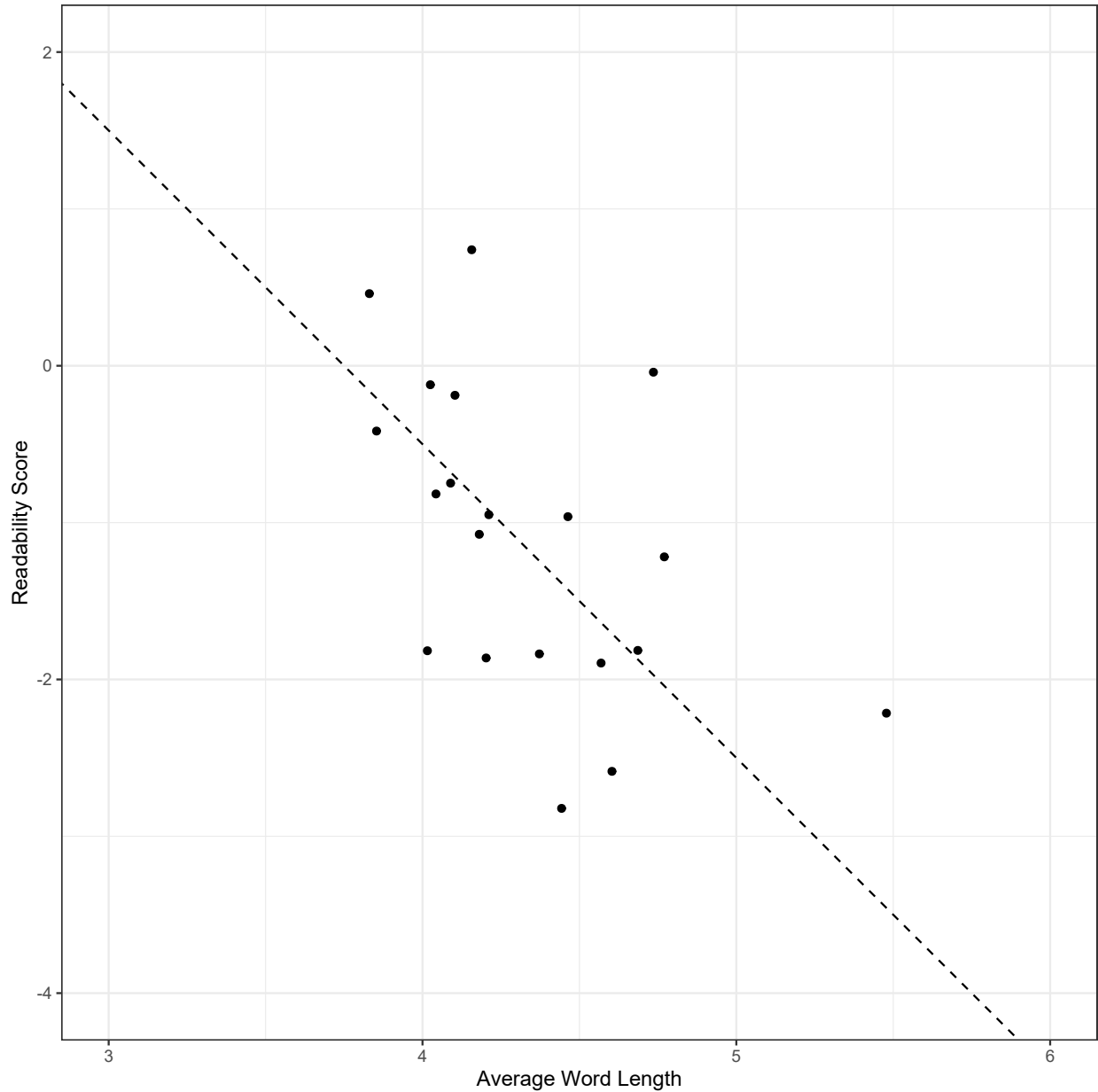


Let's consider a simple linear regression model such that the readability score is the outcome (Y) and average word length is the predictor(X_1). Our regression model would be

$$Y_i = \beta_0 + \beta_1 X_i + \epsilon_i.$$

In this case, the set of coefficients, $\{\beta_0, \beta_1\}$, represents a linear line. We can come up with any set of $\{\beta_0, \beta_1\}$ coefficients and use it as our model. For instance, suppose I guesstimate that these coefficients are $\{\beta_0, \beta_1\} = \{7.5, -2\}$. Then, my model would look like the following.

$$Y_i = 7.5 - 2X_i + \epsilon_i.$$



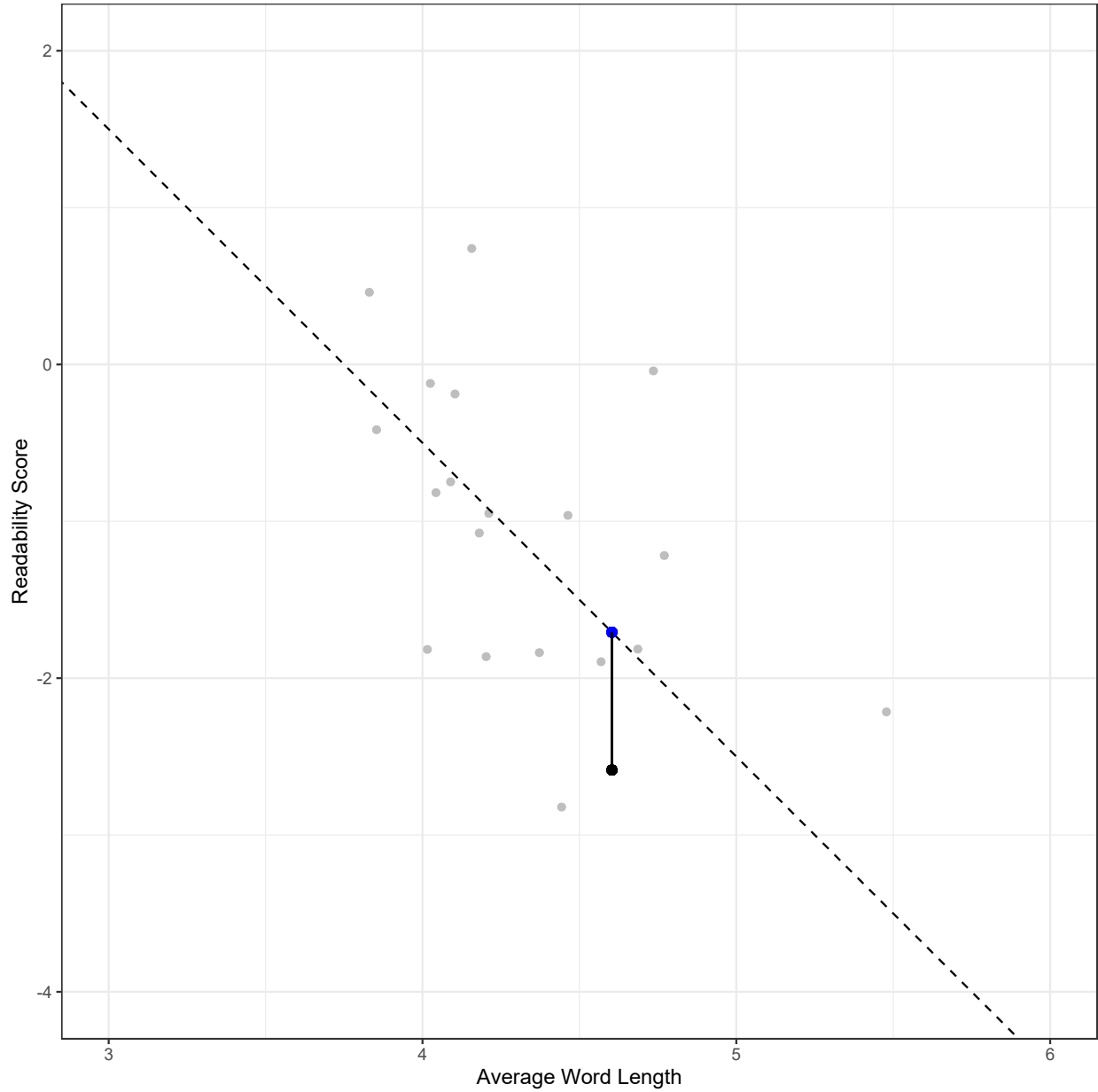
Using this model, I can predict the target readability score for all the observation in my dataset. For instance, the average word length is 4.604 for the first reading passage. Then, my prediction of readability score based on this model would be -1.708. On the other side, the observed value of the target score for this observation is -2.586. This discrepancy between the observed value and my model predicts is the model error (residual) for the first observation and captured in the ϵ term in the model.

$$Y_1 = 7.5 - 2X_1 + \epsilon_1.$$

$$\hat{Y}_1 = 7.5 - 2 * 4.604 = -1.708$$

$$\hat{\epsilon}_1 = -2.586 - (-1.708) = -0.878$$

We can visualize this in the plot. The black dot represents the observed data point, and the blue dot on the line represents the model prediction for a given X value. The vertical distance between these two data points is the model error for this particular observation.

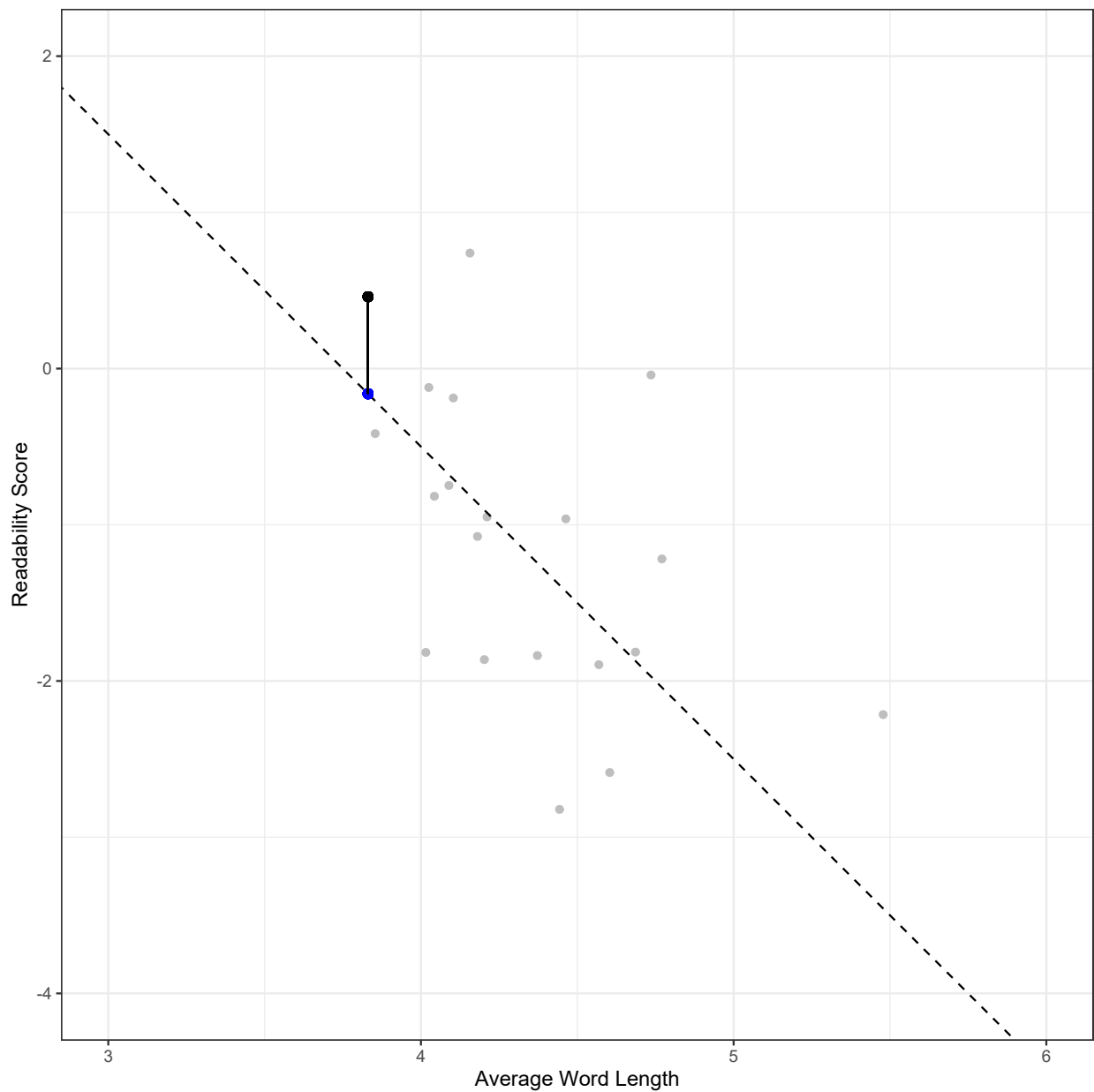


We can do the same experiment for the second observation. The average word length is 3.830 for the second reading passage. The model predicts a readability score of be -0.161. Observed value of the target score for this observation is 0.459. Therefore the model error for the second observation would be 0.62.

$$Y_2 = 7.5 - 2X_2 + \epsilon_2.$$

$$\hat{Y}_2 = 7.5 - 2 * 3.830 = -0.161$$

$$\hat{\epsilon}_2 = 0.459 - (-0.161) = 0.62$$



Using a similar approach, we can calculate the model error for every single observation.

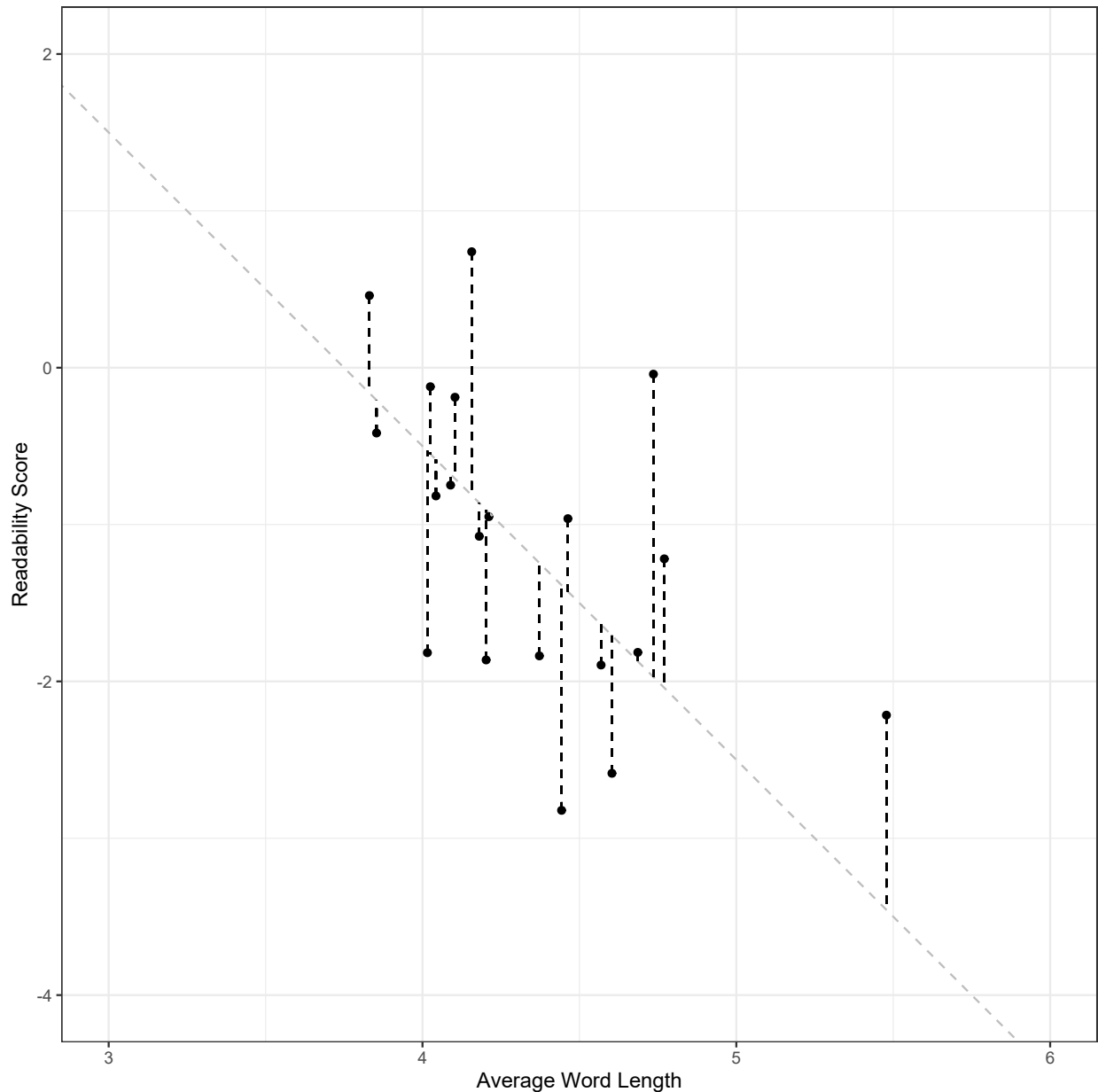
```
d <- readability_sub[,c('mean.wl', 'target')]

d$predicted <- d$mean.wl*-2 + 7.5
d$error <- d$target - d$predicted

d
```

	mean.wl	target	predicted	error
1	4.603659	-2.58590836	-1.7073171	-0.87859129

2	3.830688	0.45993224	-0.1613757	0.62130790
3	4.180851	-1.07470758	-0.8617021	-0.21300545
4	4.015544	-1.81700402	-0.5310881	-1.28591594
5	4.686047	-1.81491744	-1.8720930	0.05717559
6	4.211340	-0.94968236	-0.9226804	-0.02700194
7	4.025000	-0.12103065	-0.5500000	0.42896935
8	4.443182	-2.82200582	-1.3863636	-1.43564218
9	4.089385	-0.74845172	-0.6787709	-0.06968077
10	4.156757	0.73948755	-0.8135135	1.55300107
11	4.463277	-0.96218937	-1.4265537	0.46436430
12	5.478261	-2.21514888	-3.4565217	1.24137286
13	4.770492	-1.21845136	-2.0409836	0.82253224
14	4.568966	-1.89544351	-1.6379310	-0.25751247
15	4.735751	-0.04101056	-1.9715026	1.93049203
16	4.372340	-1.83716516	-1.2446809	-0.59248431
17	4.103448	-0.18818586	-0.7068966	0.51871069
18	4.042857	-0.81739314	-0.5857143	-0.23167886
19	4.202703	-1.86307557	-0.9054054	-0.95767016
20	3.853535	-0.41630158	-0.2070707	-0.20923088



While it is helpful to see the model error for every single observation, we will need to aggregate them in some way to form an overall measure of the total amount of error for this model. Some alternatives for aggregating these individual errors could be using

- a. the sum of the residuals (SR),
- b. the sum of absolute value of residuals (SAR), or
- c. the sum of squared residuals (SSR)

Among these alternatives, (a) is not a useful aggregation as the positive residuals and negative residuals will cancel each other and (a) may misrepresent the total amount of error for all observations. Both (b) and (c) are plausible alternatives and can be used. On the other hand, (b) is less desirable because the absolute values are mathematically more difficult to deal with (ask a calculus professor!). So, (c) seems to be a good way of aggregating the total amount of error, it is mathematically easy to work with. We can show (c) in a mathematical notation as the following.

$$SSR = \sum_{i=1}^N (Y_i - (\beta_0 + \beta_1 X_i))^2$$

$$SSR = \sum_{i=1}^N (Y_i - \hat{Y}_i)^2$$

$$SSR = \sum_{i=1}^N \epsilon_i^2$$

For our model, the sum of squared residuals would be 15.384.

```
sum(d$error^2)
```

```
[1] 15.38364
```

Now, how do we know that the set of coefficients we guesstimate, $\{\beta_0, \beta_1\} = \{7.5, -2\}$, is a good model? Is there any other set of coefficients that would provide less error than this model? The only way of knowing this is to try a bunch of different models and see if we can find a better one that gives us better predictions (smaller residuals). But, there is literally infinite pairs of $\{\beta_0, \beta_1\}$ coefficients, so which ones we should try?

Below, I will do a quick exploration. For instance, suppose the potential range for my intercept (β_0) is from -10 to 10 and I will consider every single possible value from -10 to 10 with increments of .1. Also, suppose the potential range for my slope (β_1) is from -2 to 2 and I will consider every single possible value from -2 to 2 with increments of .01. Given that every single combination of β_0 and β_1 indicates a different model, these settings suggest a total of 80,601 models to explore. If you are crazy enough, you can try every single model and compute the SSR. Then, we can plot them in a 3D by putting β_0 on the X-axis, β_1 on the Y-axis, and SSR on the Z-axis. Check the plot below and tell me if you can explore and find the minimum of this surface.

WebGL is not supported by
your browser - visit
<https://get.webgl.org> for
more info

The finding the best set of $\{\beta_0, \beta_1\}$ coefficients that minimizes the sum of squared residuals is indeed an optimization problem. For any optimization problem, there is a **loss function** we either try to minimize or maximize. In this case, our loss function is the sum of squared residuals.

$$Loss = \sum_{i=1}^N (Y_i - (\beta_0 + \beta_1 X_i))^2$$

In this loss function, X and Y values are observed data, and $\{\beta_0, \beta_1\}$ are unknown parameters. The goal of optimization is to find the set $\{\beta_0, \beta_1\}$ coefficients that provides the minimum value of this function. Once this minima of this function is found, we can argue that the corresponding coefficients are our best solution for the regression model.

In this case, this is a good-looking surface with a single global minima, and it is not difficult to find the minimum of this loss function. We also have an analytical solution to find its minima because of its simplicity. Most of the time, the optimization problems are more difficult, and we solve them using numerical techniques such as steepest ascent (or descent), newton-raphson, quasi-newton, genetic algorithm and many more.

Matrix Solution

For most regression problems, we can find the best set of coefficients with a simple matrix operations. Let's first see how we can represent the regression problem in matrix form. Suppose that I wrote the regression model presented in the earlier section for every single observation in a dataset with a sample size of N.

$$Y_1 = \beta_0 + \beta_1 X_1 + \epsilon_1.$$

$$Y_2 = \beta_0 + \beta_1 X_2 + \epsilon_2.$$

$$Y_3 = \beta_0 + \beta_1 X_3 + \epsilon_3.$$

...

...

...

$$Y_{20} = \beta_0 + \beta_1 X_{20} + \epsilon_{20}.$$

We can write all of these equations in a much simpler format as

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

such that \mathbf{Y} is an N x 1 column vector of observed values for the outcome variable, \mathbf{X} is an N x (P+1) design matrix of observed values for predictor variables, and $\boldsymbol{\beta}$ is an (P+1) x 1 column vector of regression coefficients, and $\boldsymbol{\epsilon}$ is an N x 1 column vector of residuals. For the problem above with our small dataset, these matrix elements would look like the following.

$$\begin{bmatrix} Y_1 \\ Y_2 \\ Y_3 \\ Y_4 \\ Y_5 \\ Y_6 \\ Y_7 \\ Y_8 \\ Y_9 \\ Y_{10} \\ Y_{11} \\ Y_{12} \\ Y_{13} \\ Y_{14} \\ Y_{15} \\ Y_{16} \\ Y_{17} \\ Y_{18} \\ Y_{19} \\ Y_{20} \end{bmatrix} = \begin{bmatrix} 1 & X_1 \\ 1 & X_2 \\ 1 & X_3 \\ 1 & X_4 \\ 1 & X_5 \\ 1 & X_6 \\ 1 & X_7 \\ 1 & X_8 \\ 1 & X_9 \\ 1 & X_{10} \\ 1 & X_{11} \\ 1 & X_{12} \\ 1 & X_{13} \\ 1 & X_{14} \\ 1 & X_{15} \\ 1 & X_{16} \\ 1 & X_{17} \\ 1 & X_{18} \\ 1 & X_{19} \\ 1 & X_{20} \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \epsilon_4 \\ \epsilon_5 \\ \epsilon_6 \\ \epsilon_7 \\ \epsilon_8 \\ \epsilon_9 \\ \epsilon_{10} \\ \epsilon_{11} \\ \epsilon_{12} \\ \epsilon_{13} \\ \epsilon_{14} \\ \epsilon_{15} \\ \epsilon_{16} \\ \epsilon_{17} \\ \epsilon_{18} \\ \epsilon_{19} \\ \epsilon_{20} \end{bmatrix}$$

Or, more specifically, we can replace the observed values of \mathbf{X} and \mathbf{Y} with the corresponding elements.

$$\begin{bmatrix} -2.59 \\ 0.46 \\ -1.07 \\ -1.82 \\ -1.81 \\ -0.95 \\ -0.12 \\ -2.82 \\ -0.75 \\ 0.74 \\ -0.96 \\ -2.22 \\ -1.22 \\ -1.90 \\ -0.04 \\ -1.84 \\ -0.19 \\ -0.82 \\ -1.86 \\ -0.42 \end{bmatrix} = \begin{bmatrix} 1 & 4.60 \\ 1 & 3.83 \\ 1 & 4.18 \\ 1 & 4.02 \\ 1 & 4.69 \\ 1 & 4.21 \\ 1 & 4.03 \\ 1 & 4.44 \\ 1 & 4.09 \\ 1 & 4.16 \\ 1 & 4.46 \\ 1 & 5.48 \\ 1 & 4.77 \\ 1 & 4.57 \\ 1 & 4.74 \\ 1 & 4.37 \\ 1 & 4.10 \\ 1 & 4.04 \\ 1 & 4.20 \\ 1 & 3.85 \end{bmatrix} \begin{bmatrix} \beta_0 \\ \beta_1 \end{bmatrix} + \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \epsilon_3 \\ \epsilon_4 \\ \epsilon_5 \\ \epsilon_6 \\ \epsilon_7 \\ \epsilon_8 \\ \epsilon_9 \\ \epsilon_{10} \\ \epsilon_{11} \\ \epsilon_{12} \\ \epsilon_{13} \\ \epsilon_{14} \\ \epsilon_{15} \\ \epsilon_{16} \\ \epsilon_{17} \\ \epsilon_{18} \\ \epsilon_{19} \\ \epsilon_{20} \end{bmatrix}$$

It can be shown that the set of $\{\beta_0, \beta_1\}$ coefficients that yields the minimum sum of squared residuals for this model can be analytically found using the following matrix operation.

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

If we apply this matrix operation to our small datasets, we will find that the best set of $\{\beta_0, \beta_1\}$ coefficients to predict the readability score with the least amount of error using the average word length as a predictor is $\{\beta_0, \beta_1\} = \{4.494, -1.290\}$. These estimates are also known as the **least square estimates**, and the best linear unbiased estimators (BLUE) for the given regression model.

```
Y <- as.matrix(readability_sub$target)
X <- as.matrix(cbind(1, readability_sub$mean.wl))
```

Y

```
      [,1]
[1,] -2.58590836
[2,]  0.45993224
[3,] -1.07470758
[4,] -1.81700402
[5,] -1.81491744
[6,] -0.94968236
```

```

[7,] -0.12103065
[8,] -2.82200582
[9,] -0.74845172
[10,] 0.73948755
[11,] -0.96218937
[12,] -2.21514888
[13,] -1.21845136
[14,] -1.89544351
[15,] -0.04101056
[16,] -1.83716516
[17,] -0.18818586
[18,] -0.81739314
[19,] -1.86307557
[20,] -0.41630158

```

```
X
```

```

      [,1]      [,2]
[1,]      1 4.603659
[2,]      1 3.830688
[3,]      1 4.180851
[4,]      1 4.015544
[5,]      1 4.686047
[6,]      1 4.211340
[7,]      1 4.025000
[8,]      1 4.443182
[9,]      1 4.089385
[10,]     1 4.156757
[11,]     1 4.463277
[12,]     1 5.478261
[13,]     1 4.770492
[14,]     1 4.568966
[15,]     1 4.735751
[16,]     1 4.372340
[17,]     1 4.103448
[18,]     1 4.042857
[19,]     1 4.202703
[20,]     1 3.853535

```

```
beta <- solve(t(X)%*%X)%*%t(X)%*%Y
```

```
beta
```

```

      [,1]
[1,] 4.493847
[2,] -1.290571

```

Once we find the best estimates for the model coefficients, we can also calculate the model predicted values and residual sum of squares for the given model and dataset.

$$\hat{Y} = X\hat{\beta}$$

$$\hat{\epsilon} = Y - \hat{Y}$$

$$RSS = \hat{\epsilon}^T \hat{\epsilon}$$

```
Y_hat <- X%*%beta
```

```
Y_hat
```

```
      [,1]
[1,] -1.4475035
[2,] -0.4499296
[3,] -0.9018403
[4,] -0.6884998
[5,] -1.5538311
[6,] -0.9411887
[7,] -0.7007034
[8,] -1.2403969
[9,] -0.7837974
[10,] -0.8707449
[11,] -1.2663309
[12,] -2.5762403
[13,] -1.6628138
[14,] -1.4027297
[15,] -1.6179787
[16,] -1.1489710
[17,] -0.8019465
[18,] -0.7237493
[19,] -0.9300414
[20,] -0.4794160
```

```
E <- Y - Y_hat
```

```
E
```

```
      [,1]
[1,] -1.138404820
[2,]  0.909861867
[3,] -0.172867283
[4,] -1.128504242
[5,] -0.261086332
[6,] -0.008493645
[7,]  0.579672713
[8,] -1.581608945
[9,]  0.035345700
[10,]  1.610232426
[11,]  0.304141555
[12,]  0.361091438
[13,]  0.444362421
[14,] -0.492713788
[15,]  1.576968115
[16,] -0.688194163
[17,]  0.613760605
[18,] -0.093643860
[19,] -0.933034170
[20,]  0.063114409
```

```
RSS <- t(E)%*%E
```

```
RSS
```

```
      [,1]  
[1,] 13.81062
```

Note that the matrix formulation is generalized to a regression model for more than one predictor. When there are more predictors in the model, the dimensions of the design matrix (\mathbf{X}) and regression coefficient matrix (β) will be different, but the matrix calculations will be identical. It is difficult to visualize the surface we are trying to minimize beyond two coefficients, but we know that the matrix solution will always provide us the set of coefficients that yields the least amount of error in our predictions.

Let's assume that we would like to expand our model by adding the number of sentences as the second predictor. Our new model will be

$$Y_i = \beta_0 + \beta_1 X_{1i} + \beta_2 X_{2i} + \epsilon_i.$$

Note that I added a subscript for X to differentiate different predictors. Let's say X_1 represents the mean word length and X_2 represents the total number of sentence length. Now, we are looking for the best set of three coefficients, $\{\beta_0, \beta_1, \beta_2\}$ that would yield the least amount of error in predicting the readability. Now, our matrix elements will look like the following:

```
Y <- as.matrix(readability_sub$target)  
X <- as.matrix(cbind(1,readability_sub[,c('mean.wl', 'sents')]))
```

```
Y
```

```
      [,1]  
[1,] -2.58590836  
[2,]  0.45993224  
[3,] -1.07470758  
[4,] -1.81700402  
[5,] -1.81491744  
[6,] -0.94968236  
[7,] -0.12103065  
[8,] -2.82200582  
[9,] -0.74845172  
[10,]  0.73948755  
[11,] -0.96218937  
[12,] -2.21514888  
[13,] -1.21845136  
[14,] -1.89544351  
[15,] -0.04101056  
[16,] -1.83716516  
[17,] -0.18818586  
[18,] -0.81739314  
[19,] -1.86307557  
[20,] -0.41630158
```

```
X
```

```

      1 mean.wl sents
[1,] 1 4.603659    7
[2,] 1 3.830688   23
[3,] 1 4.180851   17
[4,] 1 4.015544    7
[5,] 1 4.686047    6
[6,] 1 4.211340   18
[7,] 1 4.025000   10
[8,] 1 4.443182    4
[9,] 1 4.089385    9
[10,] 1 4.156757   28
[11,] 1 4.463277   15
[12,] 1 5.478261   10
[13,] 1 4.770492   10
[14,] 1 4.568966    8
[15,] 1 4.735751   19
[16,] 1 4.372340   15
[17,] 1 4.103448    6
[18,] 1 4.042857    6
[19,] 1 4.202703    7
[20,] 1 3.853535   19

```

We will get the following estimates for $\{\beta_0, \beta_1, \beta_2\} = \{1.821, -.929, .090\}$ yielding a value of 7.365 for the residual sum of squares.

```
beta <- solve(t(X)%*%X)%*%t(X)%*%Y
```

```
beta
```

```

      [,1]
1      1.82055156
mean.wl -0.92858249
sents    0.09029887

```

```
Y_hat <- X%*%beta
```

```
E <- Y - Y_hat
```

```
RSS <- t(E)%*%E
```

```
RSS
```

```

      [,1]
[1,] 7.365244

```

Hat Matrix

Hat matrix plays an important role in diagnosing observations in the dataset with unusually high influence on predictions, and can be calculated in matrix form from the design matrix (\mathbf{X}). In practice, it can be used to detect observations with potentially extreme values in terms of the predictors.

Remember that we did use the following matrix operation to calculate the predicted values,

$$\hat{\mathbf{Y}} = \mathbf{X}\hat{\boldsymbol{\beta}}.$$

We also know $\hat{\boldsymbol{\beta}}$ is equal,

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y}.$$

If we replace the matrix form of $\hat{\boldsymbol{\beta}}$ in the first equation, we obtain the following equation,

$$\hat{\mathbf{Y}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{Y},$$

or equivalently we can write

$$\hat{\mathbf{Y}} = \mathbf{H}\mathbf{Y},$$

where \mathbf{H} is defined as

$$\mathbf{H} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T.$$

`lm()` function

Model Interpretation

Model Evaluation

Linear Regression with Regularization

Ridge Penalty

Lasso Penalty

Elastic Net

Wrapping up

Building a Prediction Model for Readability Scores

Using the Prediction Model for a New Text