

# Multiple Linear Regression

Nate Wells

Math 243: Stat Learning

September 14th, 2020

# Outline

In today's class, we will...

- Generalize the simple regression model to include more than 1 predictor
- Quantify model accuracy for linear regression models (both simple and multiple)
- Implement multiple regression in R

## Section 1

# Multiple Regression

## Many Simple Linear Regression Models?

We are often presented situations where several explanatory variables could be used to predict values of a single response variable.

## Many Simple Linear Regression Models?

We are often presented situations where several explanatory variables could be used to predict values of a single response variable.

- **Response:** Home price
- **Predictors:** square feet, number of bedrooms, number of bathrooms

## Many Simple Linear Regression Models?

We are often presented situations where several explanatory variables could be used to predict values of a single response variable.

- **Response:** Home price
- **Predictors:** square feet, number of bedrooms, number of bathrooms
- **Response:** Professor age in photo
- **Predictors:** number of static lines, proportion gray hair, skin laxity

# Many Simple Linear Regression Models?

We are often presented situations where several explanatory variables could be used to predict values of a single response variable.

- **Response:** Home price
- **Predictors:** square feet, number of bedrooms, number of bathrooms
  
- **Response:** Professor age in photo
- **Predictors:** number of static lines, proportion gray hair, skin laxity

In each case, we could create simple linear regression models for each predictor variable.

## Many Simple Linear Regression Models?

We are often presented situations where several explanatory variables could be used to predict values of a single response variable.

- **Response:** Home price
- **Predictors:** square feet, number of bedrooms, number of bathrooms
- **Response:** Professor age in photo
- **Predictors:** number of static lines, proportion gray hair, skin laxity

In each case, we could create simple linear regression models for each predictor variable.

- But its not clear how to combine estimates from multiple models.



# Many Simple Linear Regression Models?

We are often presented situations where several explanatory variables could be used to predict values of a single response variable.

- **Response:** Home price
- **Predictors:** square feet, number of bedrooms, number of bathrooms
  
- **Response:** Professor age in photo
- **Predictors:** number of static lines, proportion gray hair, skin laxity

In each case, we could create simple linear regression models for each predictor variable.

- But its not clear how to combine estimates from multiple models.
- The results may be misleading. Several explanatory variables may be highly correlated.

## Many Simple Linear Regression Models?

We are often presented situations where several explanatory variables could be used to predict values of a single response variable.

- **Response:** Home price
- **Predictors:** square feet, number of bedrooms, number of bathrooms
- **Response:** Professor age in photo
- **Predictors:** number of static lines, proportion gray hair, skin laxity

In each case, we could create simple linear regression models for each predictor variable.

- But its not clear how to combine estimates from multiple models.
- The results may be misleading. Several explanatory variables may be highly correlated.
- And even if none of the predictors have strong association with the response, it is likely we will observe a significant predictor just due to chance.

## Many Simple Linear Regression Models?

We are often presented situations where several explanatory variables could be used to predict values of a single response variable.

- **Response:** Home price
- **Predictors:** square feet, number of bedrooms, number of bathrooms
- **Response:** Professor age in photo
- **Predictors:** number of static lines, proportion gray hair, skin laxity

In each case, we could create simple linear regression models for each predictor variable.

- But its not clear how to combine estimates from multiple models.
- The results may be misleading. Several explanatory variables may be highly correlated.
- And even if none of the predictors have strong association with the response, it is likely we will observe a significant predictor just due to chance.

Could we get better predictive power by including all explanatory variables in the *same* model?

## Multiple Regression Model

In a **simple linear regression model** (SLR), we express the response variable  $Y$  as a linear function  $f$  of one predictor variable  $X$ :

$$Y = f(X) + \epsilon$$

and estimate  $f$  using

$$\hat{Y} = \hat{f}(X) = \hat{\beta}_0 + \hat{\beta}_1 X$$

## Multiple Regression Model

In a **simple linear regression model** (SLR), we express the response variable  $Y$  as a linear function  $f$  of one predictor variable  $X$ :

$$Y = f(X) + \epsilon$$

and estimate  $f$  using

$$\hat{Y} = \hat{f}(X) = \hat{\beta}_0 + \hat{\beta}_1 X$$

In a **multiple linear regression model** (MLR), we express the response variable  $Y$  as a linear combination  $f$  of  $p$  predictors  $X_1, X_2, \dots, X_p$ :

$$Y = f(X_1, \dots, X_p) + \epsilon$$

and estimate  $f$  using

$$\hat{Y} = \hat{f}(X_1, \dots, X_p) = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \dots + \hat{\beta}_p X_p$$

## Multiple Regression Model

In a **simple linear regression model** (SLR), we express the response variable  $Y$  as a linear function  $f$  of one predictor variable  $X$ :

$$Y = f(X) + \epsilon$$

and estimate  $f$  using

$$\hat{Y} = \hat{f}(X) = \hat{\beta}_0 + \hat{\beta}_1 X$$

In a **multiple linear regression model** (MLR), we express the response variable  $Y$  as a linear combination  $f$  of  $p$  predictors  $X_1, X_2, \dots, X_p$ :

$$Y = f(X_1, \dots, X_p) + \epsilon$$

and estimate  $f$  using

$$\hat{Y} = \hat{f}(X_1, \dots, X_p) = \hat{\beta}_0 + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \dots + \hat{\beta}_p X_p$$

- In the MLR model, we allow predictors to either be quantitative or binary categorical (i.e taking values 0 or 1 corresponding to failure or success)

## Finding Parameters

To create an SLR model, we found the equation of a line that minimizes RSS, where

$$\text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i),$$

which has the solution

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

## Finding Parameters

To create an SLR model, we found the equation of a line that minimizes RSS, where

$$\text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i),$$

which has the solution

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

And in R, we computed the coefficients using

```
my_mod <- lm(Y ~ X, data = my_data)
summary(my_mod)
```



## Finding Parameters

To create an SLR model, we found the equation of a line that minimizes RSS, where

$$\text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i),$$

which has the solution

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

And in R, we computed the coefficients using

```
my_mod <- lm(Y ~ X, data = my_data)
summary(my_mod)
```

To create an MLR model...

## Finding Parameters

To create an SLR model, we found the equation of a line that minimizes RSS, where

$$\text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i),$$

which has the solution

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

And in R, we computed the coefficients using

```
my_mod <- lm(Y ~ X, data = my_data)
summary(my_mod)
```

To create an MLR model...

we do the exact same thing!

## Finding Parameters MLR

To create a MLR model, we find the equation of a **hyperplane** in  $\mathbb{R}^{p+1}$  that minimizes RSS, where

$$\text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \cdots - \hat{\beta}_p x_{ip})^2,$$

which has the solution

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

## Finding Parameters MLR

To create a MLR model, we find the equation of a **hyperplane** in  $\mathbb{R}^{p+1}$  that minimizes RSS, where

$$\text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \cdots - \hat{\beta}_p x_{ip})^2,$$

which has the solution

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

- If we have 2 predictors, the equation describes a plane in 3D space.

## Finding Parameters MLR

To create a MLR model, we find the equation of a **hyperplane** in  $\mathbb{R}^{p+1}$  that minimizes RSS, where

$$\text{RSS} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_{i1} - \dots - \hat{\beta}_p x_{ip})^2,$$

which has the solution

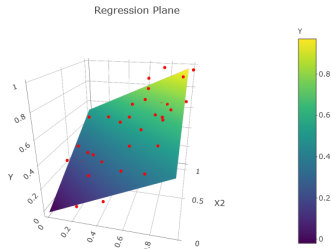
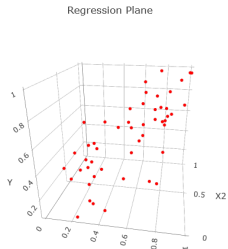
$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- If we have 2 predictors, the equation describes a plane in 3D space.

We even use the exact same R code to fit the linear model:

```
my_mod <- lm(Y ~ X1 + X2 + ... + Xp, data = my_data)
```

# The Plane of Best Fit



An interactive graphic available under Monday 9-14 on schedule page on course website

## Example: Credit Card Debt

The `Credit` dataset in the `ISLR` package contains (fabricated) credit card debt and other financial and demographic information for 400 individuals.

## Example: Credit Card Debt

The `Credit` dataset in the `ISLR` package contains (fabricated) credit card debt and other financial and demographic information for 400 individuals.

**Goal:** Build a model that allows us to predict credit debt given financial and demographic information



## Example: Credit Card Debt

The `Credit` dataset in the `ISLR` package contains (fabricated) credit card debt and other financial and demographic information for 400 individuals.

**Goal:** Build a model that allows us to predict credit debt given financial and demographic information

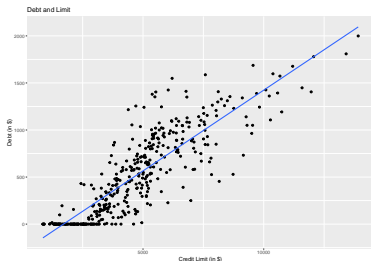
We first consider `balance` as a function of `credit_limit` and `income`

## Example: Credit Card Debt

The `Credit` dataset in the `ISLR` package contains (fabricated) credit card debt and other financial and demographic information for 400 individuals.

**Goal:** Build a model that allows us to predict credit debt given financial and demographic information

We first consider `balance` as a function of `credit_limit` and `income`



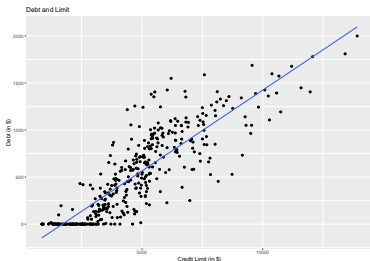
$$R = 0.86 \quad \widehat{Debt} = -292.8 + 0.17 \cdot \text{Limit}$$

## Example: Credit Card Debt

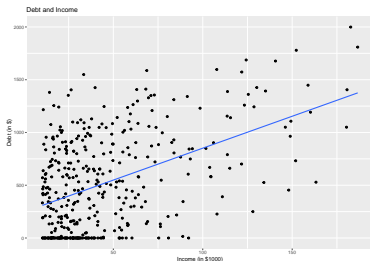
The `Credit` dataset in the `ISLR` package contains (fabricated) credit card debt and other financial and demographic information for 400 individuals.

**Goal:** Build a model that allows us to predict credit debt given financial and demographic information

We first consider `balance` as a function of `credit_limit` and `income`



$$R = 0.86 \quad \widehat{Debt} = -292.8 + 0.17 \cdot \text{Limit}$$



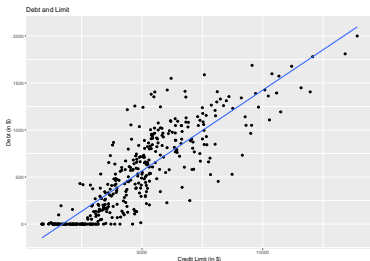
$$R = 0.46 \quad \widehat{Debt} = 246.51 + 6.048 \cdot \text{Income}$$

## Example: Credit Card Debt

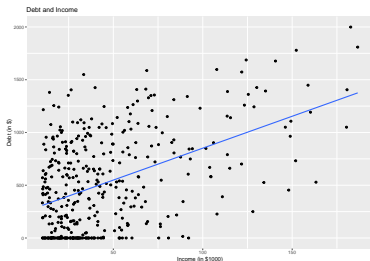
The `Credit` dataset in the `ISLR` package contains (fabricated) credit card debt and other financial and demographic information for 400 individuals.

**Goal:** Build a model that allows us to predict credit debt given financial and demographic information

We first consider `balance` as a function of `credit_limit` and `income`



$$R = 0.86 \quad \hat{Debt} = -292.8 + 0.17 \cdot \text{Limit}$$



$$R = 0.46 \quad \hat{Debt} = 246.51 + 6.048 \cdot \text{Income}$$

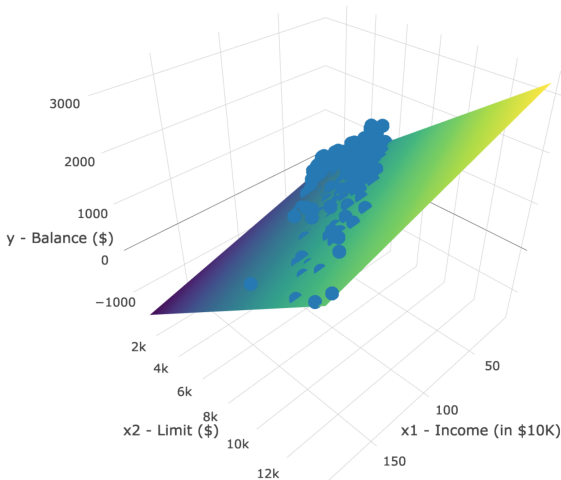
Both variables have some explanatory power for `Debt`.

## The Regression Plane

How do Limit and Income *together* explain Debt?

# The Regression Plane

How do Limit and Income *together* explain Debt?



## Multiple Regression for Debt

Let's find the MLR model

```
mod<-lm(Balance ~ Limit + Income, data = Credit)
```

## Multiple Regression for Debt

Let's find the MLR model

```
mod<-lm(Balance ~ Limit + Income, data = Credit)
```

And investigate the regression table

```
summary(mod)$coefficients
```

##	Estimate	Std. Error	t value	Pr(> t )
## (Intercept)	-385.1792604	19.464801525	-19.78850	3.878764e-61
## Limit	0.2643216	0.005879729	44.95471	7.717386e-158
## Income	-7.6633230	0.385072058	-19.90101	1.260933e-61



## Multiple Regression for Debt

Let's find the MLR model

```
mod<-lm(Balance ~ Limit + Income, data = Credit)
```

And investigate the regression table

```
summary(mod)$coefficients
```

##	Estimate	Std. Error	t value	Pr(> t )
## (Intercept)	-385.1792604	19.464801525	-19.78850	3.878764e-61
## Limit	0.2643216	0.005879729	44.95471	7.717386e-158
## Income	-7.6633230	0.385072058	-19.90101	1.260933e-61

Which gives us the regression equation:

$$\hat{Debt} = -385.179 + 0.264 \cdot \text{Limit} - 0.7663 \cdot \text{Income}$$

## Multiple Regression for Debt

Let's find the MLR model

```
mod<-lm(Balance ~ Limit + Income, data = Credit)
```

And investigate the regression table

```
summary(mod)$coefficients
```

##	Estimate	Std. Error	t value	Pr(> t )
## (Intercept)	-385.1792604	19.464801525	-19.78850	3.878764e-61
## Limit	0.2643216	0.005879729	44.95471	7.717386e-158
## Income	-7.6633230	0.385072058	-19.90101	1.260933e-61

Which gives us the regression equation:

$$\hat{Debt} = -385.179 + 0.264 \cdot \text{Limit} - 0.7663 \cdot \text{Income}$$

- For **fixed** value of Income, increasing Credit Limit by \$1 increases debt by an average of \$0.264.

## Multiple Regression for Debt

Let's find the MLR model

```
mod<-lm(Balance ~ Limit + Income, data = Credit)
```

And investigate the regression table

```
summary(mod)$coefficients
```

##	Estimate	Std. Error	t value	Pr(> t )
## (Intercept)	-385.1792604	19.464801525	-19.78850	3.878764e-61
## Limit	0.2643216	0.005879729	44.95471	7.717386e-158
## Income	-7.6633230	0.385072058	-19.90101	1.260933e-61

Which gives us the regression equation:

$$\hat{Debt} = -385.179 + 0.264 \cdot \text{Limit} - 0.7663 \cdot \text{Income}$$

- For **fixed** value of Income, increasing Credit Limit by \$1 increases debt by an average of \$0.264.
- While for **fixed** value of Limit, increasing Income by \$1000 decreases debt by an average of \$7.66.

# Comparing MLR and SLR

Wait...

## Comparing MLR and SLR

Wait...

- The SLR for Debt and Income was

$$\hat{Debt} = 246.51 + 6.048 \cdot \text{Income}$$

## Comparing MLR and SLR

Wait...

- The SLR for Debt and Income was

$$\hat{Debt} = 246.51 + 6.048 \cdot \text{Income}$$

- That is, increasing Income by \$1000 **INCREASED** debt by \$6.05.

## Comparing MLR and SLR

Wait...

- The SLR for Debt and Income was

$$\hat{Debt} = 246.51 + 6.048 \cdot \text{Income}$$

- That is, increasing Income by \$1000 **INCREASED** debt by \$6.05.
- But the MLR is

$$\hat{Debt} = -385.179 + 0.264 \cdot \text{Limit} - 0.7663 \cdot \text{Income}$$

## Comparing MLR and SLR

Wait...

- The SLR for Debt and Income was

$$\hat{Debt} = 246.51 + 6.048 \cdot \text{Income}$$

- That is, increasing Income by \$1000 **INCREASED** debt by \$6.05.
- But the MLR is

$$\hat{Debt} = -385.179 + 0.264 \cdot \text{Limit} - 0.7663 \cdot \text{Income}$$

- Not only has MLR given us a new rate of change, but it's completely switched the direction!



## Comparing MLR and SLR

Wait...

- The SLR for Debt and Income was

$$\hat{Debt} = 246.51 + 6.048 \cdot \text{Income}$$

- That is, increasing Income by \$1000 **INCREASED** debt by \$6.05.
- But the MLR is

$$\hat{Debt} = -385.179 + 0.264 \cdot \text{Limit} - 0.7663 \cdot \text{Income}$$

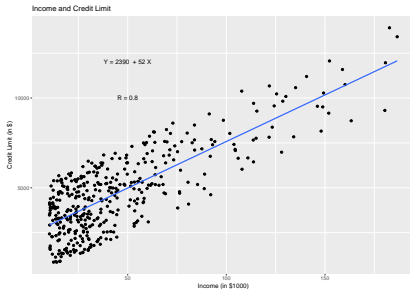
- Not only has MLR given us a new rate of change, but it's completely switched the direction!
- How is this possible?

## Income and Credit Limit

Let's consider the relationship between income and credit limit

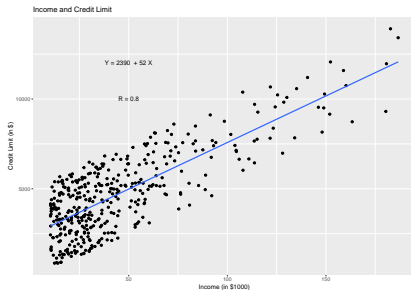
# Income and Credit Limit

Let's consider the relationship between income and credit limit



# Income and Credit Limit

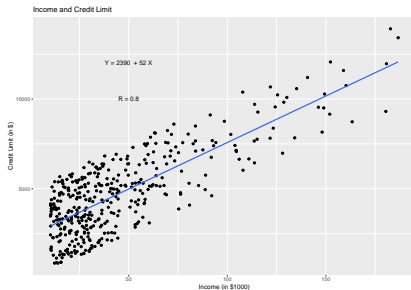
Let's consider the relationship between income and credit limit



In a vacuum, as income increases, so too does credit limit.

# Income and Credit Limit

Let's consider the relationship between income and credit limit

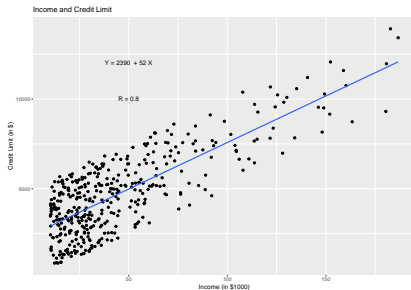


In a vacuum, as income increases, so too does credit limit.

- So in the SLR model, when we assess the change in Debt due to increase in Income, we are implicitly also increasing Credit Limit

## Income and Credit Limit

Let's consider the relationship between income and credit limit

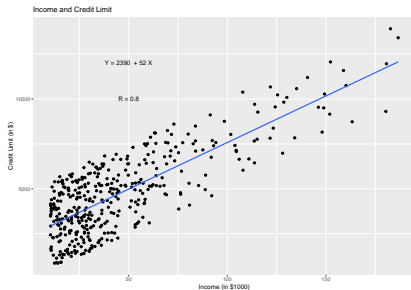


In a vacuum, as income increases, so too does credit limit.

- So in the SLR model, when we assess the change in Debt due to increase in Income, we are implicitly also increasing Credit Limit
  - We could say Credit Limit is a confounding variable in the SLR model.

# Income and Credit Limit

Let's consider the relationship between income and credit limit



In a vacuum, as income increases, so too does credit limit.

- So in the SLR model, when we assess the change in Debt due to increase in Income, we are implicitly also increasing Credit Limit
  - We could say Credit Limit is a confounding variable in the SLR model.

## The Regression Plane Revisited

In the MLR model, we may freely change both Income and Credit Limit



## The Regression Plane Revisited

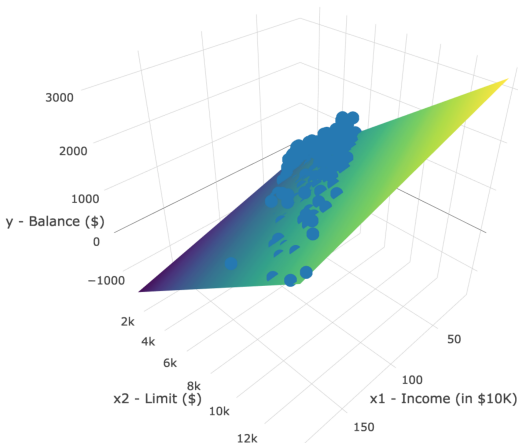
In the MLR model, we may freely change both Income and Credit Limit

- This corresponds to the fact that there is a unique Debt point on the regression plane for each pair of Income / Credit Limit values.

# The Regression Plane Revisited

In the MLR model, we may freely change both Income and Credit Limit

- This corresponds to the fact that there is a unique Debt point on the regression plane for each pair of Income / Credit Limit values.

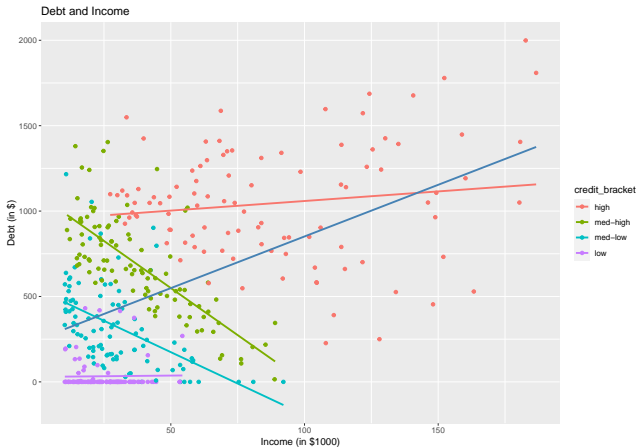


## Debt vs. Income Revisited

We can lump Credit Limits into 4 brackets (low, med-low, med-high, high) to create a categorical variable and analyze the SLR for Debt and Income for each level of Credit Limit

## Debt vs. Income Revisited

We can lump Credit Limits into 4 brackets (low, med-low, med-high, high) to create a categorical variable and analyze the SLR for Debt and Income for each level of Credit Limit



## Section 2

# Assessing Model Accuracy

## How Strong is a Linear Model?

In an linear model model,

$$Y = f(X) + \epsilon$$

So even if we could perfectly predict  $f$  using  $\hat{f}$ , our model would still have non-zero MSE.

## How Strong is a Linear Model?

In an linear model model,

$$Y = f(X) + \epsilon$$

So even if we could perfectly predict  $f$  using  $\hat{f}$ , our model would still have non-zero MSE.

The **Residual Standard Error** (RSE) measures the average size of deviations of the response from the linear regression line, is given by

$$\text{RSE} = \sqrt{\frac{1}{n-1-p} \text{RSS}} = \sqrt{\frac{1}{n-1-p} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

## How Strong is a Linear Model?

In an linear model model,

$$Y = f(X) + \epsilon$$

So even if we could perfectly predict  $f$  using  $\hat{f}$ , our model would still have non-zero MSE.

The **Residual Standard Error** (RSE) measures the average size of deviations of the response from the linear regression line, is given by

$$\text{RSE} = \sqrt{\frac{1}{n-1-p} \text{RSS}} = \sqrt{\frac{1}{n-1-p} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

It has the property that

$$E(\text{RSE}) \approx \text{Var}(\epsilon)$$



# Poll

Which of the following is most likely to decrease as more and more predictors are added to a linear model?

- a MSE
- b RSS
- c RSE
- d  $\text{Var}(\epsilon)$

## The $R^2$ statistic

Large RSE indicates poor model fit, while small RSE indicates good fit. But how do we determine how small is **small**?

## The $R^2$ statistic

Large RSE indicates poor model fit, while small RSE indicates good fit. But how do we determine how small is **small**?

- The answer depends on the units of  $Y$

## The $R^2$ statistic

Large RSE indicates poor model fit, while small RSE indicates good fit. But how do we determine how small is **small**?

- The answer depends on the units of  $Y$

An alternative, standardized measure of goodness of fit is the  $R^2$  statistic:

$$R^2 = 1 - \frac{\text{RSS}}{\text{TSS}} \quad \text{where TSS} = \sum_{i=1}^n (y_i - \bar{y})^2$$

## The $R^2$ statistic

Large RSE indicates poor model fit, while small RSE indicates good fit. But how do we determine how small is **small**?

- The answer depends on the units of  $Y$

An alternative, standardized measure of goodness of fit is the  $R^2$  statistic:

$$R^2 = 1 - \frac{\text{RSS}}{\text{TSS}} \quad \text{where TSS} = \text{sum}_{i=1}^n (y_i - \bar{y})^2$$

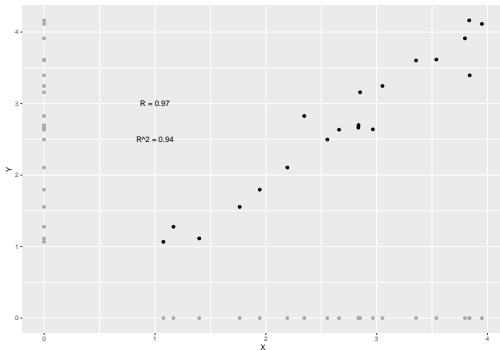
- The value of  $R^2$  is always between 0 and 1, and represents the percentage of variability in values of the response just due to variability in the predictors.

## Values of $R^2$

If  $R^2 \approx 1$ : nearly all the variability in response is due to variability in the predictor variable.

# Values of $R^2$

If  $R^2 \approx 1$ : nearly all the variability in response is due to variability in the predictor variable.



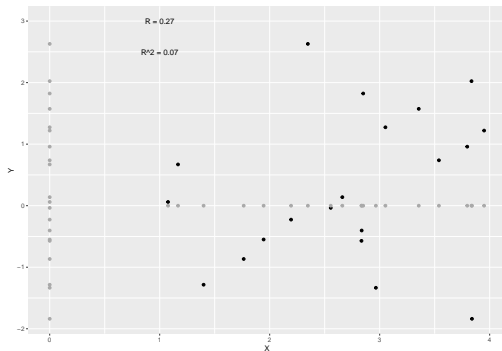
## Values of $R^2$

If  $R^2 \approx 0$ : almost none of the variability in response is due to variability in the predictor variable.



## Values of $R^2$

If  $R^2 \approx 0$ : almost none of the variability in response is due to variability in the predictor variable.



Formulas for  $R^2$  in terms of correlation

For SLR,

$$R^2 = [\text{Cor}(X, Y)]^2 = \left[ \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} \right]^2 = \left[ \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \right]^2$$

## Formulas for $R^2$ in terms of correlation

For SLR,

$$R^2 = [\text{Cor}(X, Y)]^2 = \left[ \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} \right]^2 = \left[ \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \right]^2$$

For MLR,

$$R^2 = [\text{Cor}(Y, \hat{Y})]^2$$

## Formulas for $R^2$ in terms of correlation

For SLR,

$$R^2 = [\text{Cor}(X, Y)]^2 = \left[ \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)\text{Var}(Y)}} \right]^2 = \left[ \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}} \right]^2$$

For MLR,

$$R^2 = [\text{Cor}(Y, \hat{Y})]^2$$

We will usually use software to compute  $R^2$ .

## Model Accuracy in R

```
mod_credit<-lm(Balance ~ Income + Limit , data = Credit)

summary(mod_credit)

##
## Call:
## lm(formula = Balance ~ Income + Limit, data = Credit)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -232.79 -115.45  -48.20   53.36  549.77
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -385.17926   19.46480  -19.79  <2e-16 ***
## Income       -7.66332    0.38507  -19.90  <2e-16 ***
## Limit         0.26432    0.00588   44.95  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 165.5 on 397 degrees of freedom
## Multiple R-squared:  0.8711, Adjusted R-squared:  0.8705
## F-statistic: 1342 on 2 and 397 DF, p-value: < 2.2e-16
```

## Model Accuracy in R

```
mod_credit<-lm(Balance ~ Income + Limit , data = Credit)
```

```
summary(mod_credit)
```

```
##  
## Call:  
## lm(formula = Balance ~ Income + Limit, data = Credit)  
##  
## Residuals:  
##      Min       1Q   Median       3Q      Max   
## -232.79 -115.45  -48.20   53.36  549.77   
##  
## Coefficients:  
##              Estimate Std. Error t value Pr(>|t|)      
## (Intercept) -385.17926   19.46480  -19.79  <2e-16 ***   
## Income       -7.66332    0.38507  -19.90  <2e-16 ***   
## Limit         0.26432    0.00588   44.95  <2e-16 ***   
## ---  
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1  
##  
## Residual standard error: 165.5 on 397 degrees of freedom  
## Multiple R-squared:  0.8711, Adjusted R-squared:  0.8705   
## F-statistic: 1342 on 2 and 397 DF,  p-value: < 2.2e-16
```

We can use `summary(mod)$r.sq` or `summary(mod)$sigma` to access  $R^2$  and RSE directly.

## Adjusted $R^2$

- It turns out that the sample's  $R^2$  gives a **biased** estimate of the variability in the *population* explained by the model.

## Adjusted $R^2$

- It turns out that the sample's  $R^2$  gives a **biased** estimate of the variability in the *population* explained by the model.
- Instead, we use the adjusted R:

$$R_{\text{adjusted}}^2 = 1 - \frac{\text{RSS}}{\text{TSS}} \frac{n-1}{n-p-1}$$



## Adjusted $R^2$

- It turns out that the sample's  $R^2$  gives a **biased** estimate of the variability in the *population* explained by the model.
- Instead, we use the adjusted  $R$ :

$$R_{\text{adjusted}}^2 = 1 - \frac{\text{RSS}}{\text{TSS}} \frac{n-1}{n-p-1}$$

- This adjusted  $R^2$  is usually a bit smaller than  $R^2$ , and the difference decreases as  $n$  gets large.

## Testing Significance

Suppose we wish to test whether at least one predictor has a significant linear relationship with the response.

## Testing Significance

Suppose we wish to test whether at least one predictor has a significant linear relationship with the response.

Why would it be incorrect to conduct  $p$  many significant tests comparing each predictor to the response?

## The Hypothesis Test

Goal: test whether any predictors are significant.

# The Hypothesis Test

Goal: test whether any predictors are significant.

$$H_0 : \beta_1 = \cdots = \beta_p = 0 \quad H_a : \text{at least one of } \beta_i \neq 0$$

## The Hypothesis Test

Goal: test whether any predictors are significant.

$$H_0 : \beta_1 = \cdots = \beta_p = 0 \quad H_a : \text{at least one of } \beta_i \neq 0$$

Test statistic:

$$F = \frac{(TSS - RSS)/p}{RSS/(n - p - 1)}$$

## The Hypothesis Test

Goal: test whether any predictors are significant.

$$H_0 : \beta_1 = \dots = \beta_p = 0 \quad H_a : \text{at least one of } \beta_i \neq 0$$

Test statistic:

$$F = \frac{(TSS - RSS)/p}{RSS/(n - p - 1)}$$

Under the null hypothesis,  $F$  is approximately  $F$ -distributed with  $p, n - p - 1$  parameters.

# The Hypothesis Test

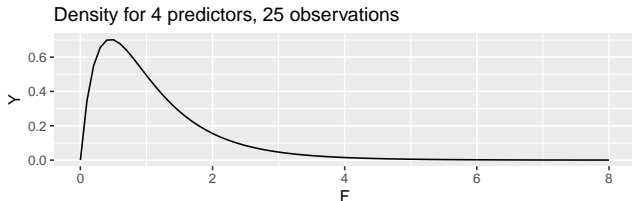
Goal: test whether any predictors are significant.

$$H_0 : \beta_1 = \dots = \beta_p = 0 \quad H_a : \text{at least one of } \beta_i \neq 0$$

Test statistic:

$$F = \frac{(TSS - RSS)/p}{RSS/(n - p - 1)}$$

Under the null hypothesis,  $F$  is approximately  $F$ -distributed with  $p, n - p - 1$  parameters.





# The Hypothesis Test

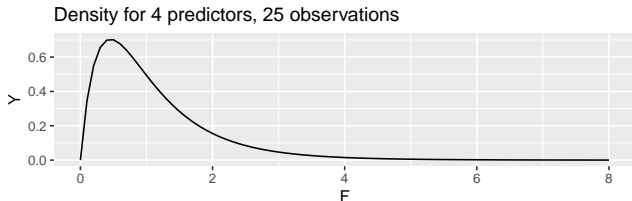
Goal: test whether any predictors are significant.

$$H_0 : \beta_1 = \dots = \beta_p = 0 \quad H_a : \text{at least one of } \beta_i \neq 0$$

Test statistic:

$$F = \frac{(TSS - RSS)/p}{RSS/(n - p - 1)}$$

Under the null hypothesis,  $F$  is approximately  $F$ -distributed with  $p, n - p - 1$  parameters.



If  $H_0$  is true, then  $F$  is typically close to 1, and rarely much greater.