
Assignment Project Exam Help

Regression
<https://powcoder.com>

Add WeChat powcoder



Agenda

| Start | End | Item |
|-------|-----|---|
| | | Regression in Action |
| | | Partitioning Data |
| | | Understanding Regression (simple univariate) |
| | | Build & Evaluate a Multiple Linear Regression |
| | | Housekeeping |

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Let's Practice

Open A Regression.R



Predict Diamond Prices with linear regression.



Agenda

| Start | End | Item |
|-------|-----|---|
| | | Regression in Action |
| | | Partitioning Data |
| | | Understanding Regression (simple univariate) |
| | | Build & Evaluate a Multiple Linear Regression |
| | | Housekeeping |

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



The Problem of Overfitting

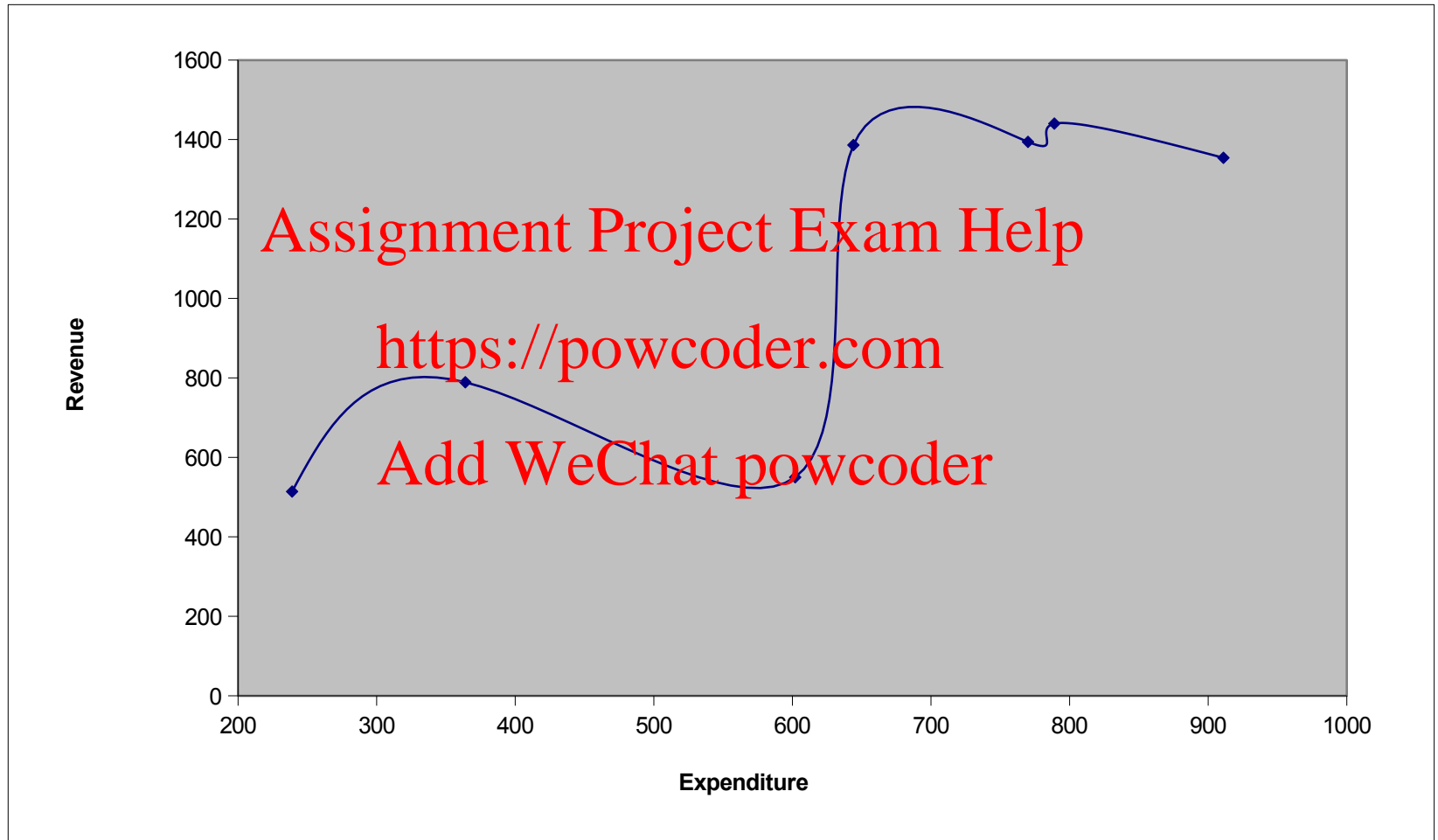
- Statistical models can produce highly complex explanations of relationships between variables
- The “fit” may be excellent
- When used with new data, models of great complexity may do not do so well.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

100% fit – not useful for new data



Another view of overfitting to a problem...



Overfitting, continued.

Causes:

- Too many predictors start to inject noise not signal
- Not adhering to a priori partitioning – *up next*
- Lack of data knowledge & problem understanding

Consequence: Deployed model will not work as well as expected with completely new data.

<https://powcoder.com>

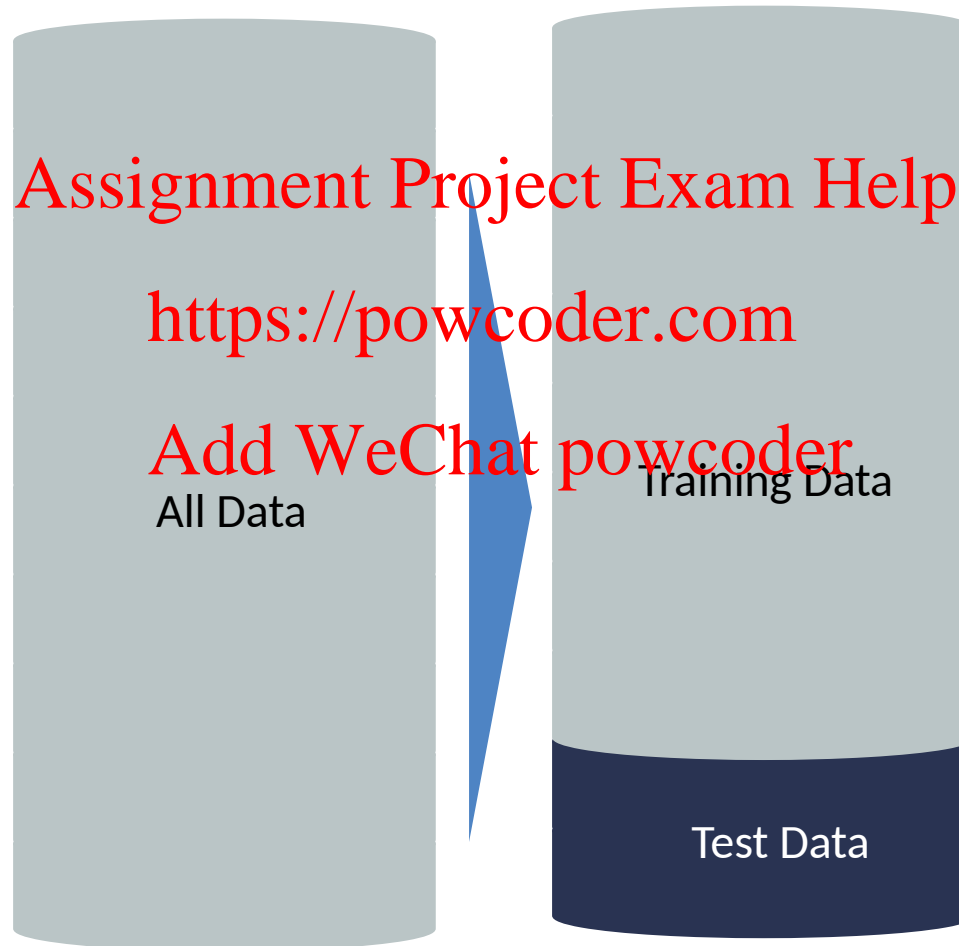
Add WeChat powcoder



Minimize Overfitting - Partitioning

Divide data into training portion and validation portion

Test model on the test portion

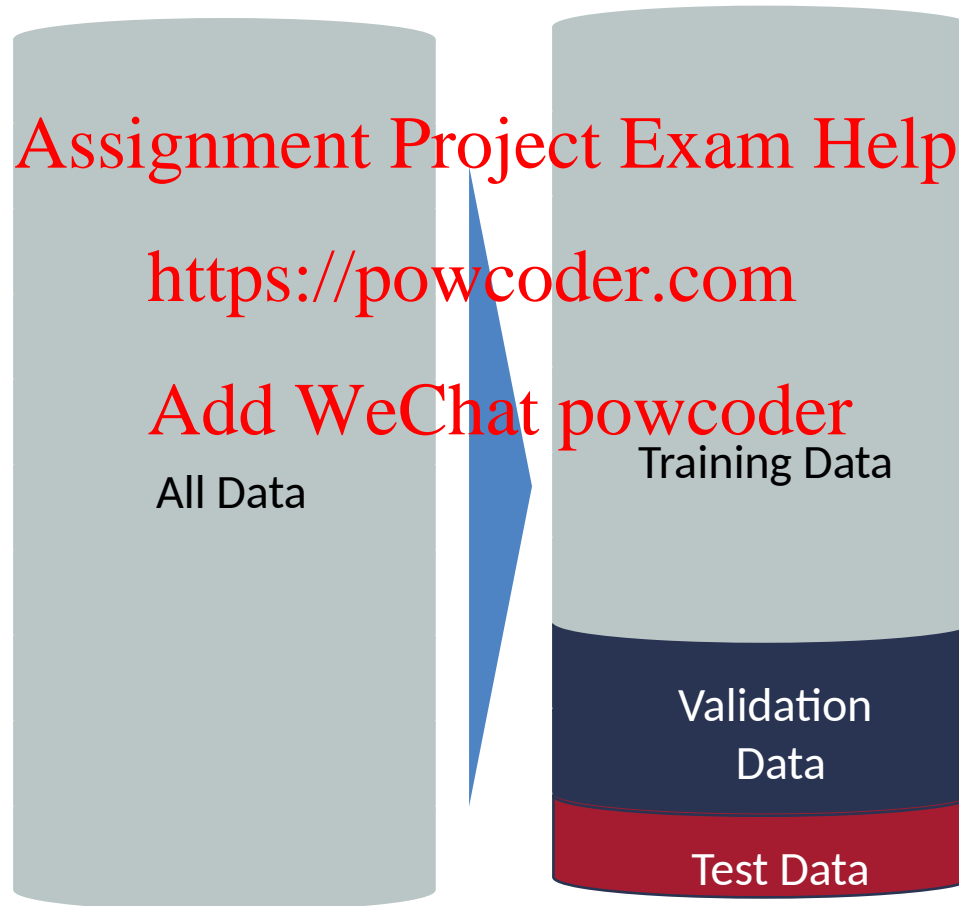


Minimize Overfitting - Partitioning

Divide data into training portion , validation & test portions

Tune a model and/or compare models with the validation portion

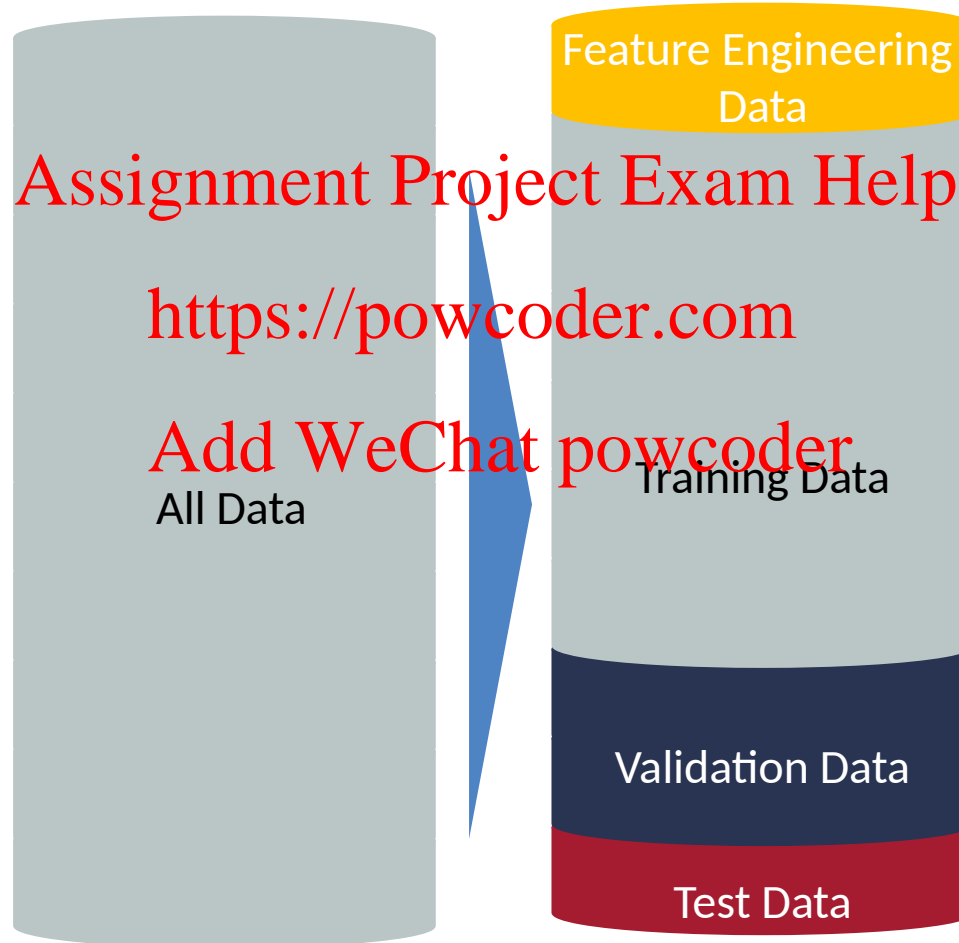
The “true” way a model will behave when launched on new data.



Best Practice

If you have enough data and the model impact is large, this is a good partitioning schema

However, this much effort is seldom undertaken.



Next Glass

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Next Glass is not out of business but leveraged chemistry & modeling for it's business model.

Let's Practice

Open B_anotherGlass.R

Train/Test

Train/Validation/Test

Engineering/Train/
Validation/Test

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Review

- Data Mining:
 - Supervised - Classification & Prediction
 - Unsupervised- Association Rules, Data Reduction, Data Exploration & Visualization
- Before algorithms can be applied, data must be explored then pre-processed (treated)
- To evaluate performance and to avoid overfitting, data partitioning is used
- Models are fit to the training partition and assessed on the validation and test partitions

Today's lesson explores partitioning and simple prediction.



Agenda

| Start | End | Item |
|-------|-----|---|
| | | Regression in Action |
| | | Partitioning Data |
| | | Understanding Regression (simple univariate) |
| | | Build & Evaluate a Multiple Linear Regression |
| | | Housekeeping |

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Supervised Learning

- Goal: Predict a single “target” or “outcome” variable

- Training data, where target value is known

<https://powcoder.com>

- Score to data where value is not known

Add WeChat powcoder

- Methods: Classification and Prediction



Supervised Learning

Inferring a function from labeled data.

“Learn from telling”, “Look at my data and I will tell you what to predict”

Business Context

Marketing- Will a customer buy yes or no? How much will a customer spend?


Operations- Will an applicant default? When will a machine break?

Sports Analytics- How many points will the Bears' QB score? What is the Bears' probability of winning?

*Requires expertise
and stakeholder buy in*

Data

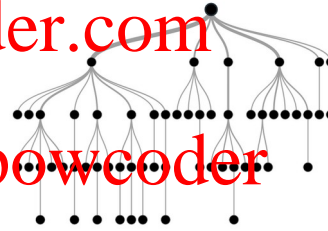
Setup

A bar chart with four bars of varying heights. A magnifying glass is positioned over the third bar from the left.

Flat “Excel” file. Each row is a record or observation. Each column is an attribute of the record.

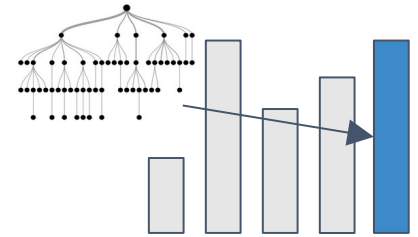
One column is the outcome, y or target attribute.

Algorithm



Modeling e.g. K-NN, linear regression, decision tree, random forest etc.

Application



Use the model to make predictions for the target label on the new data.

Supervised Learning Example

Inferring a function from labeled data.

“Learn from telling”, “Look at my data and I will tell you what to predict”

Assignment Project Exam Help



#

$=f(\dots)$

<https://powcoder.com>

Add WeChat powcoder

What impacts ice cream

sales?

Linear Regression for continuous outcomes



#

Assignment Project Exam Help

= + (*temperature) + (*day) + (*price) + error

<https://powcoder.com>

Add WeChat powcoder

Some linear combination of temperature values, day of the week dummy variables and price estimate the number of cones that will sell.

The linear combination equation captures information

Assignment Project Exam Help

outcome

coefficients

$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p + \epsilon,$

constant

predictors

error (noise)

<https://powcoder.com>


Add WeChat powcoder

The diagram shows the linear combination equation $Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p + \epsilon$. Arrows point from labels to parts of the equation: 'outcome' points to Y , 'coefficients' points to the β terms, 'constant' points to β_0 , 'predictors' points to the x terms, and 'error (noise)' points to ϵ . Overlaid on the equation is red text: 'Assignment Project Exam Help' at the top, 'https://powcoder.com' in the middle, and 'Add WeChat powcoder' at the bottom.

The linear combination equation captures information

Outcome:
The “dependent”, “y”
or “target”.
Number of Ice Cream Cones

Assignment Project Exam Help


$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p + \epsilon,$$

<https://powcoder.com>
Add WeChat powcoder

The linear combination equation captures information

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \epsilon$$

<https://powcoder.com>

Constant:

The “intercept” or “beta-naught” has no predictor associated with it.

Avg. Number of Ice Cream Cones expected to sell if predictors were all 0.

Add WeChat powcoder

The linear combination equation captures information

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$$

<https://powcoder.com>

Add WeChat powcoder

Predictors:

The “informative features”, “x” or “independent” variables.
Variables affecting sales in the data, temp, day, & price.

The linear combination equation captures information

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \epsilon$$

<https://powcoder.com>

Add WeChat powcoder

Coefficients:

The “weight”, “betas” or
“coefficients” multiplied
with the specific “x”
variable value.

The linear combination equation captures information

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_n x_n + \epsilon$$

<https://powcoder.com>

Add WeChat powcoder

Error:

The “error” or “noise” represents the value the equation is wrong compared to the actual Y.

The linear combination equation captures information

outcome

coefficients

constant

error (noise)

predictors

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$$

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

The diagram shows the linear combination equation $Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$. Arrows point from labels to parts of the equation: 'outcome' to Y , 'coefficients' to the β terms, 'constant' to β_0 , 'error (noise)' to ϵ , and 'predictors' to the x terms. Overlaid on the equation are three red text elements: 'Assignment Project Exam Help', the URL 'https://powcoder.com', and 'Add WeChat powcoder'.

The combinations of beta coefficients seeks to minimize the squared errors between the actual Y values and the equation. **This combination manifests as the “best fit line”**



#

$$= -0.05 * \text{temperature} + 0.001 * \text{saturday_dummy} - 0.5 * \text{price}$$

| Beta-Naught | Temperature | Saturday_dummy | Price |
|-------------|-------------------|----------------|----------|
| 6 | 0.25 * 80 degrees | 3 * 1 | -0.5 * 5 |
| 5 | 0.25 * 88 | 3 * 0 | -0.5 * 2 |



#

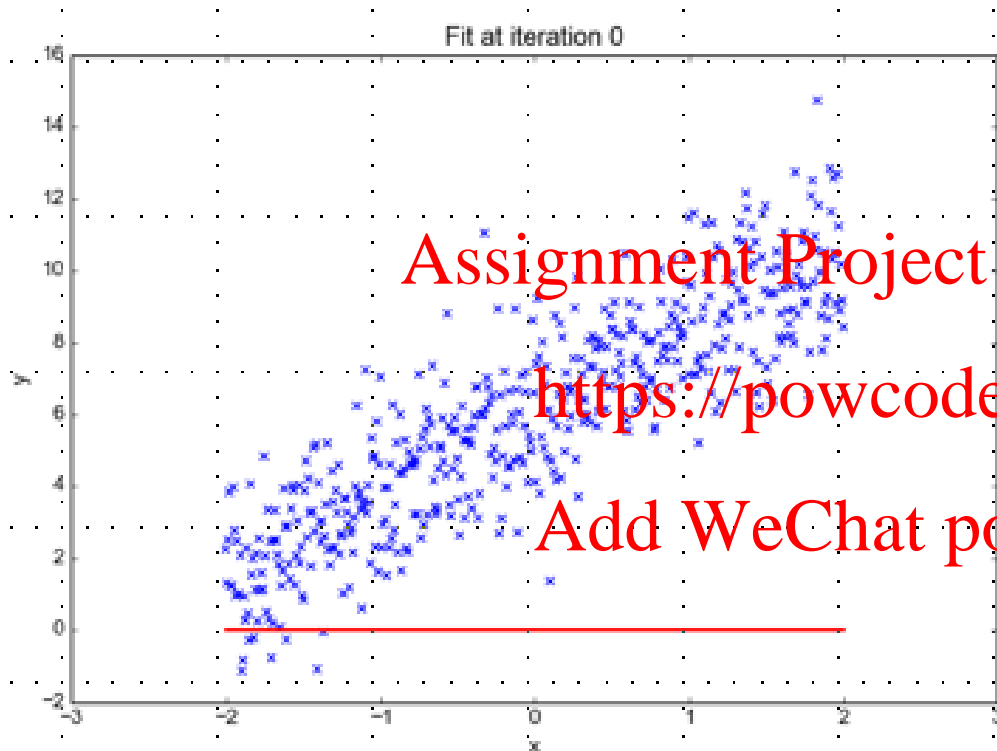
$$= -0.25 * \text{temperature} + 1 * \text{saturday_dummy} + -0.5 * \text{price}$$

<https://powcoder.com>

Add WeChat powcoder

| Beta-Naught | Temperature | Saturday_dummy | Price | Best Fit Prediction |
|-------------|-------------------|----------------|----------|---------------------|
| 6 | 0.25 * 80 degrees | 3 * 1 | -0.5 * 5 | 26.5 |
| 6 | 0.25 * 88 | 3 * 0 | -0.5 * 2 | 27 |

Minimizing the Sum of Ordinary Least Squared Errors



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Equation:

$$Y = 0 + (0 * x)$$

Beta "Naught" = 0

Intercept is 0

X beta coefficient = 0

No slope

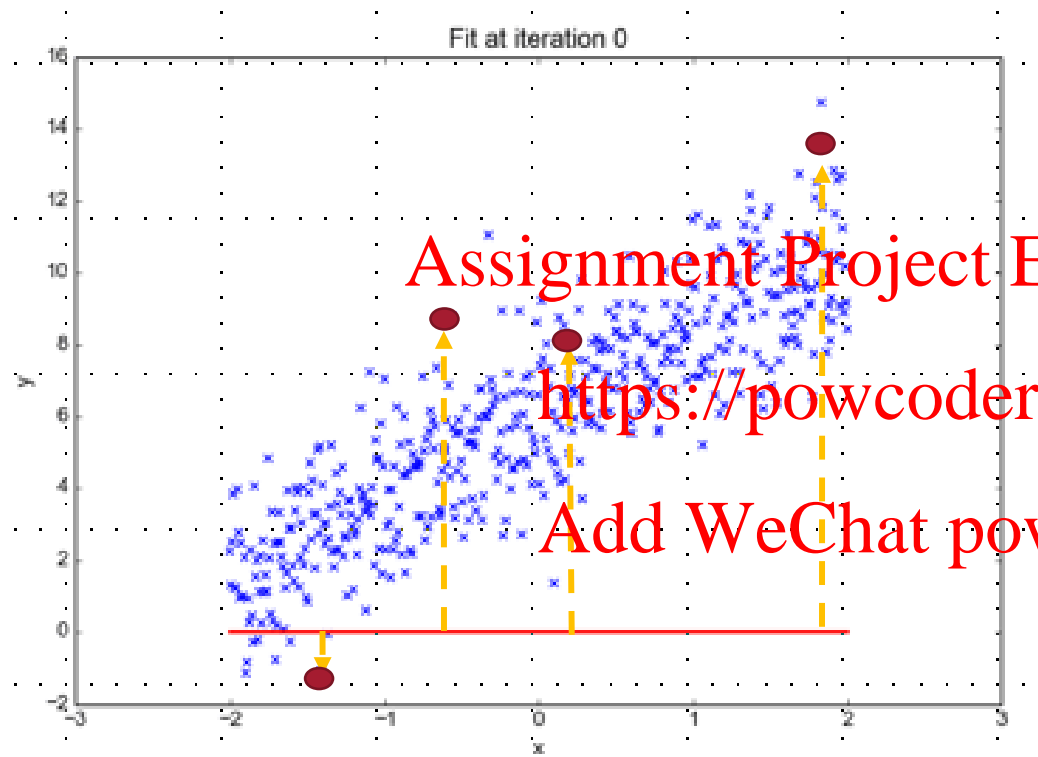
Blue points Y Values represent
actual outcome.

MINUS

Red line is the predicted outcome

Equals the Error

Big Errors



Equation:
 $Y = 0 + (0 * x)$
Beta "Naught" = 0
Intercept is 0
X beta coefficient = 0
No slope

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

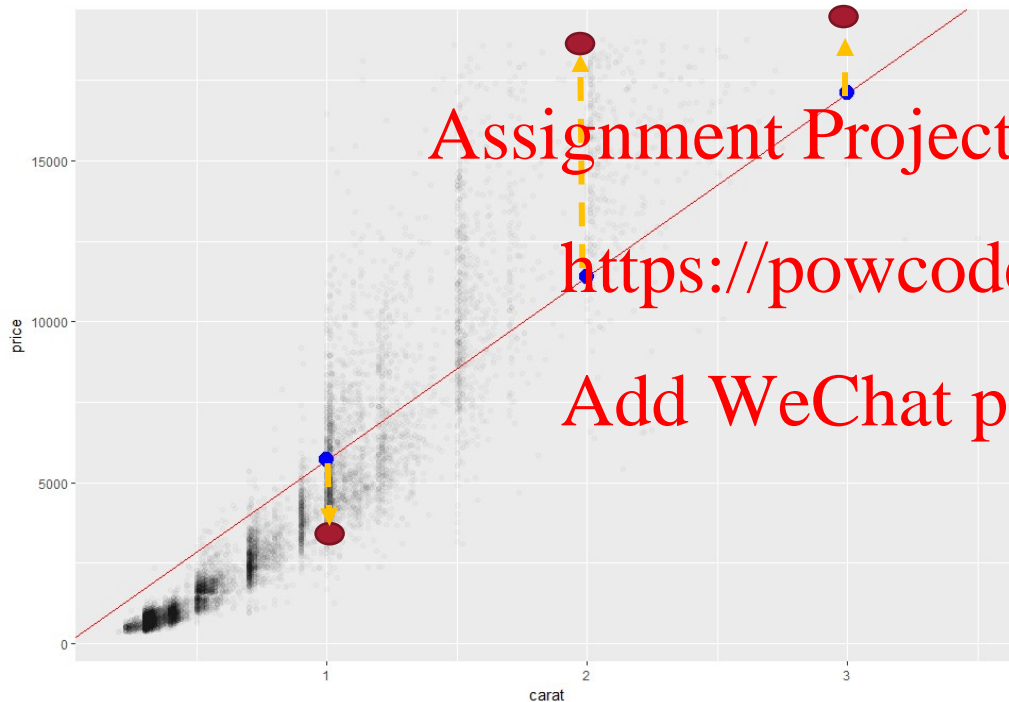
Blue points Y Values represent actual outcome.

MINUS

Red line is the predicted outcome

Equals the Error

What's really going on?



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

- Errors between a prediction and actual.

- Notice some are negative and some are positive

Why Squared Error?

Blue points Y Values represent actual outcome.

MINUS

Red line is the predicted outcome

Equals the Error

| Cones | | Prediction | | Error |
|-------|--|------------|--|-------|
| 85 | | 67 | | 18 |
| 48 | | 42 | | 6 |
| 45 | | 54 | | -9 |
| 27 | | 95 | | -68 |
| 32 | | 1 | | 31 |
| 30 | | 48 | | -18 |
| 69 | | 51 | | 18 |
| 80 | | 95 | | -15 |
| 15 | | 20 | | -5 |
| 61 | | 22 | | 39 |

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

$$(18 + 6 + -9 + -68 + 31 + -18 + 18 + -15 + -5 + 39) = -3$$

Why Squared Error?

Without squaring the errors positive & negative prediction errors cancel each other out.

(18² + 6² + -9² + -68² + 31² + -18² + 18² + -15² + -5² + 39²)

Assignment Project Exam Help

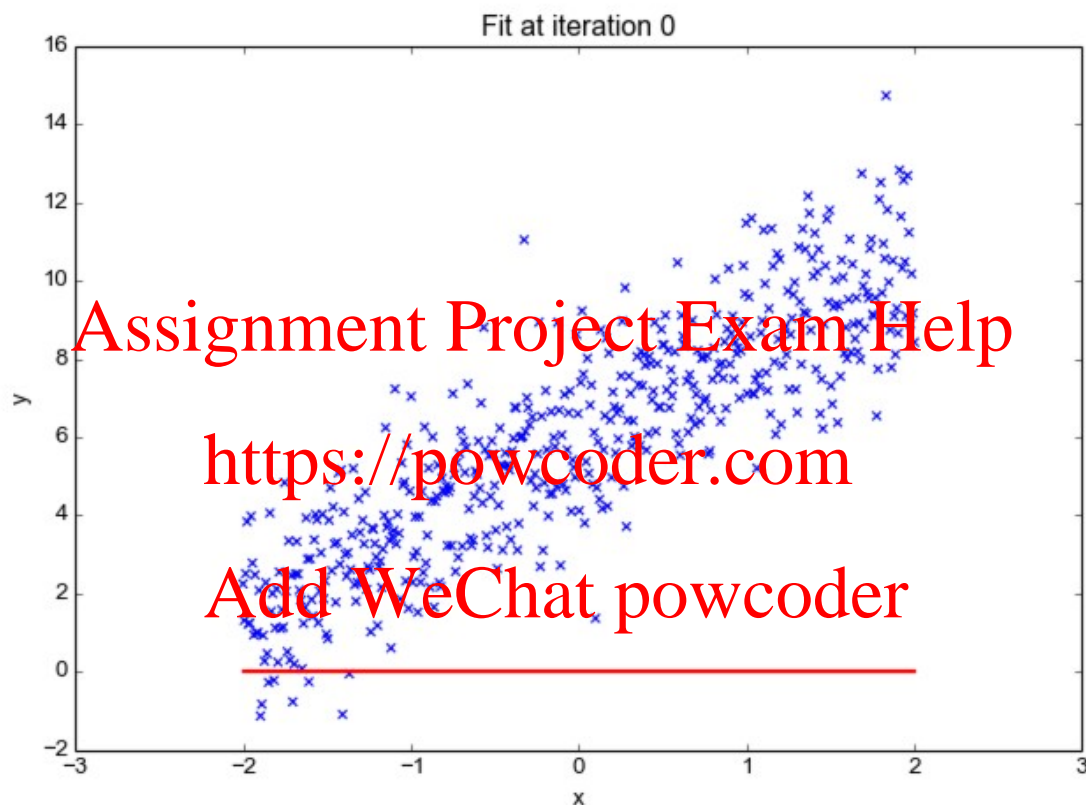
 <https://powcoder.com>

Add WeChat powcoder

$$324 + 36 + 81 + 4624 + 961 + 324 + 324 + 225 + 25 + 1521 = 8445$$

Squaring the error means all errors have the same impact on the optimization function.

So what is really going on?



The algorithm is optimizing the inputs and weights (beta's) to **minimize the sum of squared errors.**

This is called "ordinary least squares (OLS)."

Let's Practice

Open C Regression v1.R

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Topics

- Explanatory vs. predictive modeling with regression
- Example: prices of Toyota Corollas
- Fitting a predictive model
- Assessing predictive accuracy
- Selecting a subset of predictors

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Explanatory Vs Predictive

Reviewing beta coefficients can explain relationships

- There is a positive relationship between number of rooms and housing prices.

```
> fit2
```

```
Call:
lm(formula = MEDV ~ RM, data = trainSet)
```

```
Coefficients:
(Intercept)    -38.704
```

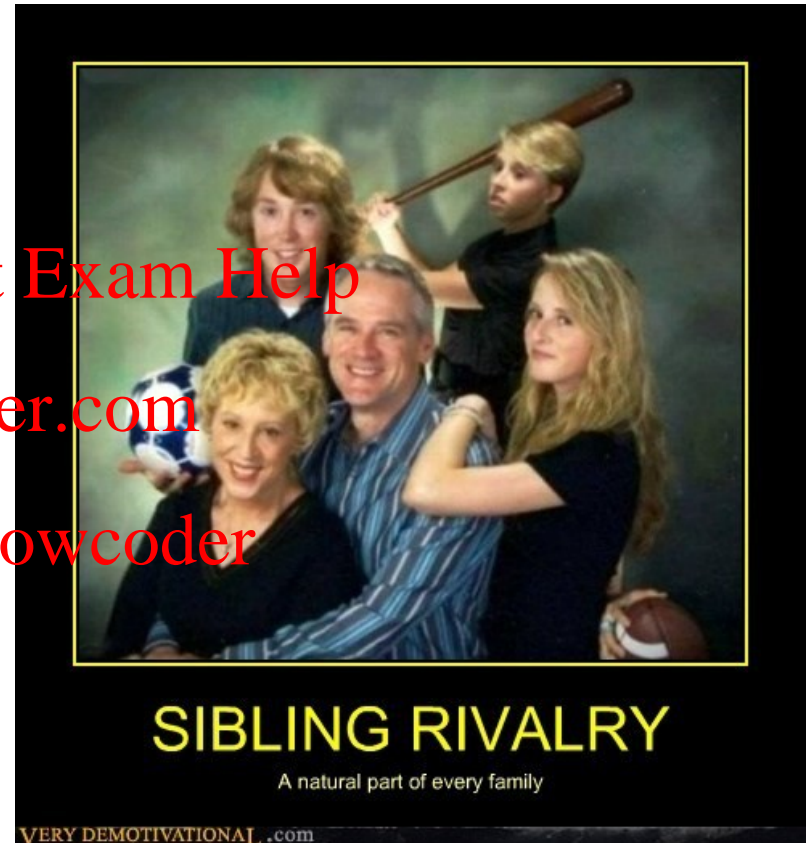
```
RM
 9.534
```

- Holding all other inputs constant the median price would increase 9.534 for each room.*

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Explanatory Modeling

Goal: Explain relationship between predictors (explanatory variables) and target

- Familiar use of regression in data analysis

Assignment Project Exam Help

- Model Goal: Fit the data well and understand the contribution of explanatory variables to the model

<https://powcoder.com>

Add WeChat powcoder

- “goodness-of-fit”: R^2 , residual analysis, p-values



Predictive Modeling

Goal: predict target values in other data where we have predictor values, but not target values

- Classic data mining context
- Model Goal: **Optimize predictive accuracy**
- Train model on training data
- Assess performance on validation (hold-out) data
- Explaining role of predictors is not primary purpose (but useful)

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Explanatory Vs Predictive Modeling

Make sure you understand the point of
your project
(explanatory or predictive)

- Do leaders want to understand a phenomena?
- Do leaders want to make accurate predictions about the future?

This impacts how you evaluate the model
and even what variables you choose.



Sibling rivalry never ends...

Agenda

| Start | End | Item |
|-------|-----|---|
| | | Regression in Action |
| | | Partitioning Data |
| | | Understanding Regression (simple univariate) |
| | | Build & Evaluate a Multiple Linear Regression |
| | | Housekeeping |

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



How does truecar.com know the price is “great?”

TRUECar. Sign In

Used Cars for Sale / Used Cars Search / Toyota / Maynard, MA

Used Toyota for Sale in Maynard, MA

Showing 1 – 30 of 88 Used Toyota Listings

Sort By: Best Match

Your Used Car Search

[Clear All](#) [Save Search](#)

Toyota x

01754 x

DISTANCE FROM 01754

75 miles

MODELS

Models

BODY STYLES

Body Styles

PRICE

Min to Max [OK](#)

YEARS

Min to Max [OK](#)

2015 Toyota Corolla LE

Mileage: 17,183 miles
Location: Maynard, MA
Exterior: Barcelona Red Metallic
Interior: Ash
VIN: 2T1BURHE6FC312663

TRUECar Rating: Great Price

\$13,595

\$1,601 below market

[View Details](#)

☐ Compare

2014 Toyota Corolla S

Mileage: 20,490 miles
Location: Auburndale, MA
Exterior: Black
Interior: Black
VIN: 2T1BURHE3EC166480

TRUECar Rating: Great Price

\$14,400

\$1,028 below market

[View Details](#)

☐ Compare

Assignment Project Exam Help
<https://powcoder.com>
Add WeChat powcoder

Yeah! I got a new job!

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Welcome Dale to TrueCar's Competitor: OldCar

Dale, can you predict used car prices? Then we will know if a car is priced above/below expected value

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Let's help Dale again.

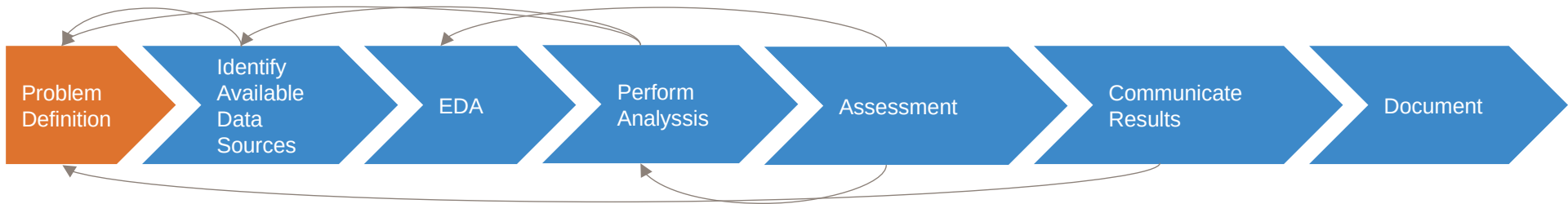
But of course. This sounds like a predictive modeling project.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Remember the Workflow?



Assignment Project Exam Help

1. Problem Formulation >> Predict Toyota Prices
2. Define data requirements >> Use ~1400 cars & car attributes
3. Explore the data >> in script
4. Perform Analysis & Create Project Artifacts >> fit a linear regression
5. Asses/Adjust the Project Artifacts >> Adj. R-Squared, P-Values etc.
6. Communicate Results >> examine the coefficients and readout Adj R²
7. Document to make it repeatable >> Keep notes in script

Let's practice.

Open D_oldCar.R

Using the preprocessing and partitioning code in your toolbox, create a linear model of Toyota prices.

Price in Euros

Tires binary 0/1 are new tires on car

Age in months as of 8/04

KM (kilometers)

Fuel Type (diesel, petrol, CNG)

HP (horsepower)

Metallic color (1=yes, 0=no)

Automatic transmission (1=yes, 0=no)

CC (cylinder volume)

Doors

Quarterly_Tax (road tax)

Weight (in kg)

Assignment Project Exam Help
<https://powcoder.com>
Add WeChat powcoder



Regression is susceptible to **multi-collinearity**

Math

The presence of two or more predictor variables sharing the same linear relationship with the outcome variable.

In other words

Two+ informative features are measuring essentially the same thing.

Example

When predicting ice cream sales you include Fahrenheit and Celsius temperatures as two separate informative features.

Don't shock the cat by having multi-collinearity.



Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Regression is susceptible to **multi-collinearity**

What happens with multi-collinearity?

The effect is exaggerated e.g. double counted.

The good news

Assignment Project Exam Help

- The problem is so prevalent/impactful that R's linear regression function handles it.
- Other algorithms are not affected by double counting

<https://powcoder.com>

Best Practice

Add WeChat powcoder

Even though R removes it, do not rely on the function. Understanding the data you put in will avoid “garbage in, garbage out” scenarios.

Make happy cats by
knowing your data



Back to R, script D!

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



The Summary of the Fit

```
> summary(fit)
```

```
Call:
lm(formula = Price ~ ., data = treatedTrain)
```

```
Residuals:
```

```
    Min       1Q   Median       3Q      Max
-11141.3  -774.3    -19.9    763.0   6653.1
```

```
Coefficients: (3 not defined because of singularities)
```

| | Estimate | Std. Error | t value | Pr(> t) | |
|------------------------|--------------|-------------|---------|----------------------|-----|
| (Intercept) | -6476.090893 | 1363.877191 | -4.748 | 0.00000231286938 | *** |
| Weight_clean | 18.866071 | 1.291646 | 14.606 | < 0.0000000000000002 | *** |
| Quarterly_Tax_clean | 11.299234 | 1.895942 | 5.960 | 0.00000000336699 | *** |
| Doors_clean | -69.703648 | 44.544760 | -1.565 | 0.118 | |
| CC_clean | -0.037332 | 0.091886 | -0.406 | 0.685 | |
| Automatic_clean | 185.878083 | 177.884776 | 1.042 | 0.298 | |
| Met_Color_clean | 106.772102 | 84.118619 | 1.269 | 0.205 | |
| HP_clean | 26.793641 | 3.812888 | 7.027 | 0.00000000000363 | *** |
| Fuel_Type_catP | 2132.597511 | 358.370269 | 5.951 | 0.00000000354884 | *** |
| Fuel_Type_catN | 0.319050 | 0.211948 | 1.505 | 0.133 | |
| Fuel_Type_catD | NA | NA | NA | NA | |
| KM_clean | -0.017617 | 0.001488 | -11.840 | < 0.0000000000000002 | *** |
| Age_08_04_clean | -123.703581 | 2.940716 | -42.066 | < 0.0000000000000002 | *** |
| Fuel_Type_lev_x_Diesel | NA | NA | NA | NA | |
| Fuel_Type_lev_x_Petrol | NA | NA | NA | NA | |

```
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 1328 on 1137 degrees of freedom
```

```
Multiple R-squared:  0.8701,    Adjusted R-squared:  0.8688
```

```
F-statistic: 692.4 on 11 and 1137 DF,  p-value: < 0.00000000000000022
```

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

The Summary of the Fit

```
> summary(fit)
```

```
Call:
```

```
lm(formula = Price ~ ., data = treatedTrain)
```

```
Residuals:
```

```
      Min       1Q   Median       3Q      Max
-11141.3  -774.3   -19.9    763.0   6653.1
```

```
Coefficients: (Not defined because of singularities)
              Estimate Std. Error t value Pr(>|t|)
(Intercept)  -6476.090893  1363.877191  -4.748  0.00000231286938 ***
weight_clean    18.866071    1.291646  14.606 < 0.0000000000000002 ***
Quarterly_Tax_clean  10.299334    1.895917   5.950  0.00000000336699 ***
Doors_clean    -69.703648    44.544760  -1.565    0.118
CC_clean       -0.037332     0.091886  -0.406    0.685
Automatic_clean  185.378082   177.884776   1.042    0.298
Met_Color_clean  106.772762    84.113619   1.269    0.205
HP_clean       16.798641    38.123840   0.620  0.000000000000363 ***
Fuel_Type_catP  2132.597511   358.370269   5.951  0.00000000354884 ***
Fuel_Type_catN    0.319050    0.211948   1.505    0.133
Fuel_Type_catD           NA           NA           NA           NA
KM_clean       -0.017617    0.001488 -11.840 < 0.0000000000000002 ***
Age_08_04_clean -123.703581    2.940716 -42.066 < 0.0000000000000002 ***
Fuel_Type_lev_x_Diesel           NA           NA           NA           NA
Fuel_Type_lev_x_Petrol           NA           NA           NA           NA
---
```

```
signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 1328 on 1137 degrees of freedom
```

```
Multiple R-squared:  0.8701,    Adjusted R-squared:  0.8688
```

```
F-statistic: 692.4 on 11 and 1137 DF,  p-value: < 0.00000000000000022
```

Auto ID of multi
colinearity

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

The Summary of the Fit

```
> summary(fit)
```

```
Call:
lm(formula = Price ~ ., data = treatedTrain)
```

```
Residuals:
```

| Min | 1Q | Median | 3Q | Max |
|----------|--------|--------|-------|--------|
| -11141.3 | -774.3 | -19.9 | 763.0 | 6653.1 |

```
Coefficients: (5 not defined because of singularities)
```

| | Estimate | Std. Error | t value | Pr(> t) | |
|------------------------|--------------|-------------|---------|----------------------|-----|
| (Intercept) | -6476.090893 | 1363.877191 | -4.748 | 0.00000231286938 | *** |
| weight_clean | 18.866071 | 1.291646 | 14.606 | < 0.0000000000000002 | *** |
| Quarterly_Tax_clean | 11.309234 | 1.305943 | 8.660 | 0.00000000336699 | *** |
| Doors_clean | -69.703648 | 44.344760 | -1.563 | 0.118 | |
| CC_clean | -0.037332 | 0.091886 | -0.406 | 0.685 | |
| Automatic_clean | 185.378082 | 177.884776 | 1.042 | 0.298 | |
| Met_Color_clean | 106.772762 | 64.113619 | 1.269 | 0.205 | |
| HP_clean | 16.793641 | 3.812888 | 4.407 | 0.0000000000363 | *** |
| Fuel_Type_catP | 2132.597511 | 358.370269 | 5.951 | 0.00000000354884 | *** |
| Fuel_Type_catN | 0.319050 | 0.211948 | 1.505 | 0.133 | |
| Fuel_Type_catD | NA | NA | NA | NA | |
| KM_clean | -0.017617 | 0.001488 | -11.840 | < 0.0000000000000002 | *** |
| Age_08_04_clean | -123.703581 | 2.940716 | -42.066 | < 0.0000000000000002 | *** |
| Fuel_Type_lev_x_Diesel | NA | NA | NA | NA | |
| Fuel_Type_lev_x_Petrol | NA | NA | NA | NA | |

```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 1328 on 1137 degrees of freedom
Multiple R-squared:  0.8701,    Adjusted R-squared:  0.8688
F-statistic: 692.4 on 11 and 1137 DF,  p-value: < 0.00000000000000022
```

Another name for errors.
Summary stats for the errors.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

The Summary of the Fit

```
> summary(fit)
```

```
Call:
lm(formula = Price ~ ., data = treatedTrain)
```

```
Residuals:
```

```
      Min       1Q   Median       3Q      Max
-11141.3  -774.3   -19.9    763.0   6653.1
```

```
Coefficients: (1 not defined because of singularities)
```

| | Estimate | Std. Error | t value | Pr(> t) | |
|------------------------|--------------|-------------|---------|----------------------|-----|
| (Intercept) | -6476.090893 | 1363.877191 | -4.748 | 0.00000231286938 | *** |
| Weight_clean | 18.866071 | 1.291646 | 14.606 | < 0.0000000000000002 | *** |
| Quarterly_Tax_clean | 11.299284 | 1.895942 | 5.960 | 0.00000000336699 | *** |
| Doors_clean | 162.703648 | 44.54480 | 3.655 | 0.000340118 | |
| CC_clean | -0.037332 | 0.091886 | -0.406 | 0.685 | |
| Automatic_clean | 185.378082 | 177.884776 | 1.042 | 0.298 | |
| Met_Color_clean | 106.772762 | 84.113619 | 1.269 | 0.205 | |
| HP_clean | 23.793741 | 3.31188 | 7.197 | 0.0000000000363 | *** |
| Fuel_Type_catP | 2132.597511 | 358.370269 | 5.951 | 0.00000000354884 | *** |
| Fuel_Type_catN | 0.319050 | 0.211948 | 1.505 | 0.133 | |
| Fuel_Type_catD | NA | NA | NA | NA | |
| KM_clean | -0.017617 | 0.001488 | -11.840 | < 0.0000000000000002 | *** |
| Age_08_04_clean | -123.703581 | 2.940716 | -42.066 | < 0.0000000000000002 | *** |
| Fuel_Type_lev_x_Diesel | NA | NA | NA | NA | |
| Fuel_Type_lev_x_Petrol | NA | NA | NA | NA | |

```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 1328 on 1137 degrees of freedom
```

```
Multiple R-squared:  0.8701,    Adjusted R-squared:  0.8688
```

```
F-statistic: 692.4 on 11 and 1137 DF,  p-value: < 0.00000000000000022
```

Treated Variable Names
i.e. informative features

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

The Summary of the Fit

```
> summary(fit)
```

```
Call:
```

```
lm(formula = Price ~ ., data = treatedTrain)
```

```
Residuals:
```

```
      Min       1Q   Median       3Q      Max
-11141.3  -774.3   -19.9    763.0   6653.1
```

```
Coefficients: (3 not defined because of singularities)
```

| | Estimate | Std. Error | t value | Pr(> t) |
|------------------------|--------------|------------|---------|--------------------------|
| (Intercept) | 16476.090893 | 1363.87191 | -4.748 | 0.0000231286938 *** |
| weight_clean | 18.866071 | 1.291646 | 14.606 | < 0.0000000000000002 *** |
| Quarterly_Tax_clean | 11.299234 | 1.895942 | 5.960 | 0.00000000336699 *** |
| Doors_clean | -69.703648 | 44.544760 | -1.565 | 0.118 |
| CC_clean | -9.037130 | 0.091836 | -0.406 | 0.685 |
| Automatic_clean | 185.378082 | 177.884776 | 1.042 | 0.298 |
| Met_Color_clean | 106.772762 | 84.113619 | 1.269 | 0.205 |
| HP_clean | 26.793641 | 3.812888 | 7.027 | 0.00000000000363 *** |
| Fuel_Type_catP | 112.37511 | 158.370269 | 0.705 | 0.00000000354884 *** |
| Fuel_Type_catN | -0.319050 | 0.211148 | -1.505 | 0.133 |
| Fuel_Type_catD | NA | NA | NA | NA |
| KM_clean | -0.017617 | 0.001488 | -11.840 | < 0.0000000000000002 *** |
| Age_08_04_clean | -123.703581 | 2.940716 | -42.066 | < 0.0000000000000002 *** |
| Fuel_Type_lev_x_Diesel | NA | NA | NA | NA |
| Fuel_Type_lev_x_Petrol | NA | NA | NA | NA |

```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 1328 on 1137 degrees of freedom
```

```
Multiple R-squared:  0.8701,    Adjusted R-squared:  0.8688
```

```
F-statistic: 692.4 on 11 and 1137 DF,  p-value: < 0.00000000000000022
```

Coefficients or Beta values

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

The Summary of the Fit

```
> summary(fit)
```

```
Call:
```

```
lm(formula = Price ~ ., data = treatedTrain)
```

```
Residuals:
```

```
      Min       1Q   Median       3Q      Max
-11141.3  -774.3   -19.9    763.0   6653.1
```

```
Coefficients: (3 not defined because of singularities)
```

| | Estimate | Std. Error | t-value | Pr(> t) |
|------------------------|--------------|------------|---------|--------------------------|
| (Intercept) | 16476.090893 | 1363.87191 | -4.748 | 0.00000231286938 *** |
| weight_clean | 18.866071 | 1.291646 | 14.606 | < 0.0000000000000002 *** |
| Quarterly_Tax_clean | 11.299234 | 1.895942 | 5.960 | 0.00000000336699 *** |
| Doors_clean | -69.703648 | 44.544760 | -1.565 | 0.118 |
| CC_clean | -9.037130 | 0.091836 | -0.406 | 0.685 |
| Automatic_clean | 185.378082 | 177.884776 | 1.042 | 0.298 |
| Met_Color_clean | 106.772762 | 84.113619 | 1.269 | 0.205 |
| HP_clean | 26.793641 | 3.812888 | 7.027 | 0.00000000000363 *** |
| Fuel_Type_catP | 112.377511 | 158.370269 | 0.711 | 0.478 |
| Fuel_Type_catN | -0.319050 | 0.211148 | -1.505 | 0.133 |
| Fuel_Type_catD | NA | NA | NA | NA |
| KM_clean | -0.017617 | 0.001488 | -11.840 | < 0.0000000000000002 *** |
| Age_08_04_clean | -123.703581 | 2.940716 | -42.066 | < 0.0000000000000002 *** |
| Fuel_Type_lev_x_Diesel | NA | NA | NA | NA |
| Fuel_Type_lev_x_Petrol | NA | NA | NA | NA |

```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 1328 on 1137 degrees of freedom
```

```
Multiple R-squared:  0.8701,    Adjusted R-squared:  0.8688
```

```
F-statistic: 692.4 on 11 and 1137 DF,  p-value: < 0.00000000000000022
```

The average amount that the coefficients vary from the actual average value of our response variable.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

The Summary of the Fit

```
> summary(fit)
```

```
Call:
```

```
lm(formula = Price ~ ., data = treatedTrain)
```

```
Residuals:
```

```
      Min       1Q   Median       3Q      Max
-11141.3  -774.3   -19.9    763.0   6653.1
```

```
Coefficients: (3 not defined because of singularities)
```

| | Estimate | Std. Error | t-value | Pr(> t) |
|------------------------|---------------|-------------|---------|------------------------|
| (Intercept) | -16476.090893 | 1363.857191 | -4.748 | 0.00000231286938 *** |
| weight_clean | 18.866071 | 1.291646 | 14.606 | 0.0000000000000002 *** |
| Quarterly_Tax_clean | 11.299234 | 1.895942 | 5.960 | 0.00000000336699 *** |
| Doors_clean | -69.703648 | 44.544760 | -1.565 | 0.118 |
| CC_clean | -9.037130 | 0.094836 | -0.406 | 0.685 |
| Automatic_clean | 185.378082 | 177.884776 | 1.042 | 0.298 |
| Met_Color_clean | 106.772762 | 84.113619 | 1.269 | 0.205 |
| HP_clean | 26.793641 | 3.812888 | 7.027 | 0.00000000000363 *** |
| Fuel_Type_catP | 112.577511 | 158.370269 | 0.711 | 0.00000000354884 *** |
| Fuel_Type_catN | -0.319050 | 0.211148 | -1.505 | 0.133 |
| Fuel_Type_catD | NA | NA | NA | NA |
| KM_clean | -0.017617 | 0.001488 | -11.840 | 0.0000000000000002 *** |
| Age_08_04_clean | -123.703581 | 2.940716 | -42.066 | 0.0000000000000002 *** |
| Fuel_Type_lev_x_Diesel | NA | NA | NA | NA |
| Fuel_Type_lev_x_Petrol | NA | NA | NA | NA |

```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 1328 on 1137 degrees of freedom
```

```
Multiple R-squared:  0.8701,    Adjusted R-squared:  0.8688
```

```
F-statistic: 692.4 on 11 and 1137 DF,  p-value: < 0.00000000000000022
```

The coefficient t-value is a measure of how many standard deviations our coefficient estimate is far away from 0. In another way, values away from 0 indicate a real relationship.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

The Summary of the Fit

```
> summary(fit)
```

```
Call:
```

```
lm(formula = Price ~ ., data = treatedTrain)
```

```
Residuals:
```

```
      Min       1Q   Median       3Q      Max
-11141.3  -774.3   -19.9   763.0  6653.1
```

```
Coefficients: (3 not defined because of singularities)
```

| | Estimate | Std. Error | t value | Pr(> t) |
|------------------------|--------------|------------|---------|--------------------------|
| (Intercept) | 16476.090893 | 1363.87191 | -4.748 | 0.0000231286938 *** |
| weight_clean | 18.866071 | 1.291646 | 14.606 | < 0.0000000000000002 *** |
| Quarterly_Tax_clean | 11.299234 | 1.895942 | 5.960 | 0.00000000336699 *** |
| Doors_clean | -69.703648 | 44.544760 | -1.565 | 0.118 |
| CC_clean | -9.037130 | 0.091836 | -0.406 | 0.685 |
| Automatic_clean | 185.378082 | 177.884776 | 1.042 | 0.298 |
| Met_Color_clean | 106.772762 | 84.113619 | 1.269 | 0.205 |
| HP_clean | 26.793641 | 3.812888 | 7.027 | 0.00000000000363 *** |
| Fuel_Type_catP | 112.87511 | 158.370269 | 0.715 | 0.476 |
| Fuel_Type_catN | -0.319050 | 0.211148 | -1.505 | 0.133 |
| Fuel_Type_catD | NA | NA | NA | NA |
| KM_clean | -0.017617 | 0.001488 | -11.840 | < 0.0000000000000002 *** |
| Age_08_04_clean | -123.703581 | 2.940716 | -42.066 | < 0.0000000000000002 *** |
| Fuel_Type_lev_x_Diesel | NA | NA | NA | NA |
| Fuel_Type_lev_x_Petrol | NA | NA | NA | NA |

```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 1328 on 1137 degrees of freedom
```

```
Multiple R-squared:  0.8701,    Adjusted R-squared:  0.8688
```

```
F-statistic: 692.4 on 11 and 1137 DF,  p-value: < 0.00000000000000022
```

P-Values

The probability of seeing a value larger than the t value.
Small p-values mean its less likely due to chance.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

The Summary of the Fit

```
> summary(fit)
```

```
Call:
```

```
lm(formula = Price ~ ., data = treatedTrain)
```

```
Residuals:
```

```
      Min       1Q   Median       3Q      Max
-11141.3  -774.3   -19.9    763.0   6653.1
```

```
Coefficients: (3 not defined because of singularities)
```

| | Estimate | Std. Error | t value | Pr(> t) |
|------------------------|--------------|------------|---------|--------------------------|
| (Intercept) | 16476.090893 | 1363.87191 | -4.748 | 0.0000231286938 *** |
| weight_clean | 18.866071 | 1.291646 | 14.606 | < 0.0000000000000002 *** |
| Quarterly_Tax_clean | 11.299234 | 1.895942 | 5.960 | 0.00000000336699 *** |
| Doors_clean | -69.703648 | 44.544760 | -1.565 | 0.118 |
| CC_clean | -9.037130 | 0.091836 | -0.406 | 0.685 |
| Automatic_clean | 185.378082 | 177.884776 | 1.042 | 0.298 |
| Met_Color_clean | 106.772762 | 84.113619 | 1.269 | 0.205 |
| HP_clean | 26.793641 | 3.812888 | 7.027 | 0.00000000000363 *** |
| Fuel_Type_catP | 112.377511 | 158.370269 | 0.705 | 0.00000000354884 *** |
| Fuel_Type_catN | -0.319050 | 0.211148 | -1.505 | 0.133 |
| Fuel_Type_catD | NA | NA | NA | NA |
| KM_clean | -0.017617 | 0.001488 | -11.840 | < 0.0000000000000002 *** |
| Age_08_04_clean | -123.703581 | 2.940716 | -42.066 | < 0.0000000000000002 *** |
| Fuel_Type_lev_x_Diesel | NA | NA | NA | NA |
| Fuel_Type_lev_x_Petrol | NA | NA | NA | NA |

```
---
```

```
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 1328 on 1137 degrees of freedom
```

```
Multiple R-squared:  0.8701,    Adjusted R-squared:  0.8688
```

```
F-statistic: 692.4 on 11 and 1137 DF,  p-value: < 0.00000000000000022
```

P-Values

In stats $p < 0.05$ is good but in business I have seen $p < 0.2$.

It's a good idea to rebuild a model without variables that do not meet the cutoff.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

The Summary of the Fit

```
> summary(fit)
```

```
Call:
lm(formula = Price ~ ., data = treatedTrain)
```

```
Residuals:
    Min       1Q   Median       3Q      Max
-11141.3  -774.3   -19.9   763.0  6653.1
```

```
Coefficients: (3 not defined because of singularities)
```

| | Estimate | Std. Error | t value | Pr(> t) |
|------------------------|--------------|------------|---------|--------------------------|
| (Intercept) | 16476.090893 | 1363.87191 | -4.748 | 0.0000231286938 *** |
| weight_clean | 18.866071 | 1.291646 | 14.606 | < 0.0000000000000002 *** |
| Quarterly_Tax_clean | 11.299234 | 1.895942 | 5.960 | 0.00000000336699 *** |
| Doors_clean | -69.703648 | 44.544760 | -1.565 | 0.118 |
| CC_clean | -9.037130 | 0.091836 | -0.406 | 0.685 |
| Automatic_clean | 185.378082 | 177.884776 | 1.042 | 0.298 |
| Met_Color_clean | 106.772762 | 84.113619 | 1.269 | 0.205 |
| HP_clean | 26.793641 | 3.812888 | 7.027 | 0.00000000000363 *** |
| Fuel_Type_catP | 112.577511 | 158.370269 | 0.711 | 0.00000000354884 *** |
| Fuel_Type_catN | -0.319050 | 0.211148 | -1.505 | 0.133 |
| Fuel_Type_catD | NA | NA | NA | NA |
| KM_clean | -0.017617 | 0.001488 | -11.840 | < 0.0000000000000002 *** |
| Age_08_04_clean | -123.703581 | 2.940716 | -42.066 | < 0.0000000000000002 *** |
| Fuel_Type_lev_x_Diesel | NA | NA | NA | NA |
| Fuel_Type_lev_x_Petrol | NA | NA | NA | NA |

```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
Residual standard error: 1328 on 1137 degrees of freedom
```

```
Multiple R-squared:  0.8701,    Adjusted R-squared:  0.8688
```

```
F-statistic: 692.4 on 11 and 1137 DF,  p-value: < 0.00000000000000022
```

P-Values
We could drop ??? Based on p values?

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

What variables should we drop with a p-value ≥ 0.2 ?

Back to the script D

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Parsimonious Model

```
> summary(fit2)

Call:
lm(formula = Price ~ ., data = treatedTrainParsimony)

Residuals:
    Min       1Q   Median       3Q      Max
-11231.1   -766.9    -24.4    769.7   6665.8

Coefficients: (3 not defined because of singularities)
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  6655.418241    1327.407929   4.984 0.00000075614045 ***
Weight_clean    19.086315     1.263358   15.108 < 0.0000000000000002 ***
Quarterly_Tax_clean  11.289780     1.895568    5.956 0.00000000344140 ***
Doors_clean   -71.134917    44.179384   -1.610    0.108
HP_clean      35.230442     3.716592    9.483 0.00000000000293 ***
Fuel_Type_catP 2169.457810    354.025248    6.128 0.00000000122404 ***
Fuel_Type_catN    0.282825     0.208225    1.358    0.175
Fuel_Type_catD      NA         NA         NA      NA
KM_clean      -0.011808     0.001482   -12.014 < 0.0000000000000002 ***
Age_08_04_clean  123.520949     2.008408   61.470 < 0.0000000000000002 ***
Fuel_Type_lev_x_Diesel      NA         NA         NA      NA
Fuel_Type_lev_x_Petrol      NA         NA         NA      NA
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1328 on 1140 degrees of freedom
Multiple R-squared:  0.8698,    Adjusted R-squared:  0.8689
F-statistic:  952 on 8 and 1140 DF,  p-value: < 0.0000000000000022
```

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Parsimony or compactness is desirable in models. The more features in a model, the more complexity we introduce, data integrity, data interactions, time to score and time to predict are all impacted.

Parsimonious Model

```
> summary(fit2)

Call:
lm(formula = Price ~ ., data = treatedTrainParsimony)

Residuals:
    Min       1Q   Median       3Q      Max
-11231.1  -766.9   -24.4    769.7   6665.8

Coefficients: (3 not defined because of singularities)
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  6655.41824    1137.40799   5.854 0.00000075614045 ***
Weight_clean    19.086315     1.263358  15.108 < 0.0000000000000002 ***
Quarterly_Tax_clean  11.289780     1.895568   5.956 0.000000000344140 ***
Doors_clean   -71.134917    44.179384  -1.610    0.108
HP_clean      35.230442     3.716572   9.483 0.000000000000293 ***
Fuel_Type_catP 2169.457810    354.025248   6.128 0.00000000122404 ***
Fuel_Type_catN    0.282825     0.208225   1.358    0.175
Fuel_Type_catD      NA         NA         NA      NA
KM_clean      -0.017808     0.001482  -12.014 < 0.0000000000000002 ***
Age_08_04_clean  123.520949     2.908408  42.470 < 0.0000000000000002 ***
Fuel_Type_lev_x_Diesel      NA         NA         NA      NA
Fuel_Type_lev_x_Petrol      NA         NA         NA      NA
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1328 on 1140 degrees of freedom
Multiple R-squared:  0.8698,    Adjusted R-squared:  0.8689
F-statistic:  952 on 8 and 1140 DF,  p-value: < 0.00000000000000022
```

R-Sq: how much of the variation are the model is fitting. R-Sq measures the linear relationship between Price & features It always lies between 0 and 1

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Compare the two models

fit1

Residual standard error: 1328 on 1137 degrees of freedom
Multiple R-squared: 0.8701, Adjusted R-squared: 0.8688
F-statistic: 692.4 on 11 and 1137 DF, p-value: < 0.000000000000000022

<https://powcoder.com>

fit2

Residual standard error: 1328 on 1140 degrees of freedom
Multiple R-squared: 0.8698, Adjusted R-squared: 0.8689
F-statistic: 952 on 8 and 1140 DF, p-value: < 0.000000000000000022

It can be said that both models explain ~87% of the variation in car prices. Dropping the variables improved accuracy and reinforces the fact that the variables didn't add value.

Evaluating a Prediction Model

RMSE- Root Mean Squared Error

MAPE- Mean Absolute Percentage Error

Assignment Project Exam Help

| Actual Values | Predicted/Forecasted |
|---------------|----------------------|
| 10 | 6 |
| 12 | 8 |
| 20 | 18 |
| 36 | 20 |

Besides P-Values which is a variable level KPI, and adjusted R-Sq there are two popular KPI for evaluating continuous model predictions.

RMSE

RMSE- Root Mean Squared Error

| y | y-hat | ERROR | ERROR-SQ |
|---------------|----------------------|--------|----------------|
| Actual Values | Predicted/Forecasted | Errors | Squared Errors |
| 10 | 16 | -6 | 36 |
| 12 | 8 | -4 | 16 |
| 20 | 17 | 3 | 9 |
| 36 | 34 | 2 | 4 |

Mean

$$\frac{36+16+9+4}{4}$$

To manually calculate RMSE, work the acronym backwards.



RMSE

RMSE- Root Mean Squared Error

| Actual Values | Predicted/Forecasted | Errors | Squared Errors |
|---------------|----------------------|--------|----------------|
| 10 | 16 | -6 | 36 |
| 12 | 8 | -4 | 16 |
| 20 | 17 | 3 | 9 |
| 36 | 34 | 2 | 4 |

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Mean

$$\frac{36+16+9+4}{4}$$

Square Root

$$=4.03$$

In the same units being measured, tells you +/- the prediction error.

MAPE

MAPE- Mean Absolute Percentage Error

| Actual Values | Predicted/Forecasted | Errors | Absolute | As % of Forecast |
|---------------|----------------------|--------|----------|------------------|
| 10 | 16 | -6 | 6 | =6/16 or 37% |
| 12 | 8 | 4 | 4 | =4/8 or 50% |
| 20 | 17 | 3 | 3 | =3/17 or 17% |
| 36 | 34 | 2 | 2 | =2/34 or 5% |

Mean of Percentages

$$\frac{37\% + 50\% + 17\% + 5\%}{4}$$

$$=27.7\%$$

Instead of squaring error, take the absolute error. Then divide that by the forecast value. Lastly calculate a mean average of all the percentage errors.

Back to script D

How good is Dale's model?

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Your Data Mining Toolbox

Previous Lessons

- Some R Programming
- Knowledge of Data Preparation
- Exploratory Data Analysis
- Basic Visualization

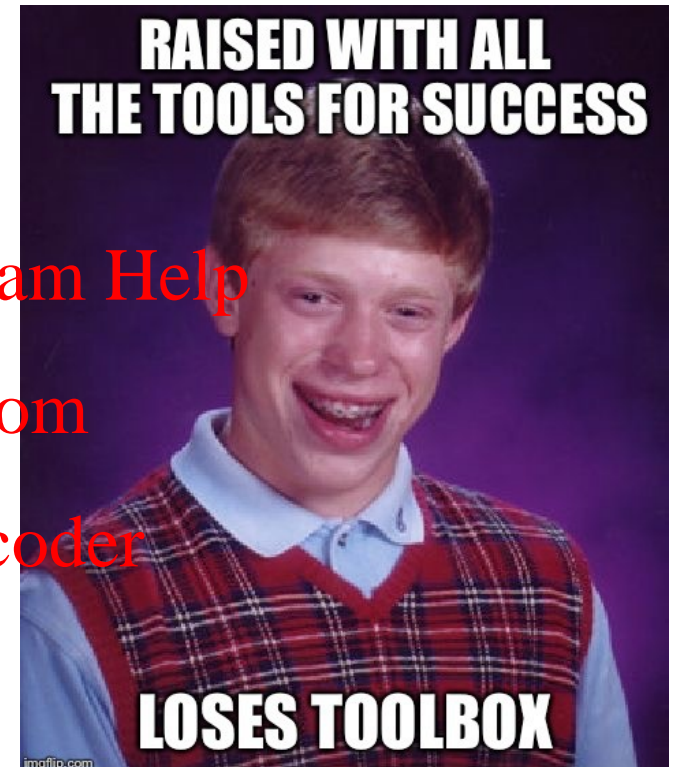
After today

- You can predict continuous business outcomes simplistically

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder



Regression is an initial starting algorithm. It puts you on a path to more complex machine learning but more importantly you can start to frame business problems in terms algorithms can understand.

Housekeeping , Reading & Homework

- Next Week is Logistic Regression, KNN

Assignment Project Exam Help

- Chapter 7
 - Chapter 10
- <https://powcoder.com>

Add WeChat powcoder

- Homework – check syllabus!

