



Hyejoon Lee & Alexis Balayre

Multiple Linear Regression

School of Aerospace, Transport and Manufacturing
Computational Software of Techniques Engineering

MSc
Academic Year: 2023 - 2024

Supervisors: Dr Stuart Barnes & Dr Peter Sherar
23th October 2023

Table of Contents

Table of Contents	ii
List of Figures	iv
List of Tables	v
1 Introduction	1
2 Datasets	2
3 Methodologies	4
3.1 Multicollinearity	4
3.1.1 Data Pre-processing	4
3.1.1.1 Advertising.csv	4
3.1.1.2 Auto.csv	5
3.2 Optimal splitting ratio	6
3.3 Polynomial Regression	7
3.4 Gradient Descent	7
3.4.1 Normalised Dataset	7
3.4.2 Gradient Clipping	8
3.5 Regression criteria	8
3.5.1 R squared value	8
4 Results & Discussion	9
4.1 Advertising dataset	9
4.1.1 Multivariate linear regression	9
4.1.2 TV + radio + newspaper	10
4.1.3 Overall	12
4.2 Auto dataset	13
4.3 Conclusion	15
References	16
A Documentation	17
B Linear Regression Custom Model 1 Source Code	19
C Linear Regression Custom Model 2 Source Code	22

D	Quadratic Regression Custom Model Source Code	24
E	Degree 5 Regression Custom Model Source Code	25
F	Utility Functions Source Code	27

List of Figures

3.1	Heat map of advertising dataset	5
3.2	Heat map of auto dataset	6
3.3	Optimal splitting ratio against the number of parameters	6
4.1	Fitting the data with custom Model	9
4.2	Fitting the data with sklearn Model	10
4.3	TV+radio+TV*radio	10
4.4	P-value of variables	11
4.5	Correlation between Budget and Sales	12
4.6	Linear Regression	13
4.7	Quadratic Regression	13
4.8	Degree 5 Regression	14

List of Tables

2.1	Advertising Data Sample	2
2.2	Advertising Data Statistics	2
2.3	Auto Data Sample	3
2.4	Auto Data Statistics	3
3.1	VIF factor	5
4.1	TV+radio+newspaper data	11
4.2	Polynomial regression results for custom model	14
4.3	Polynomial regression results for sklearn library	14

Chapter 1

Introduction

Linear regression is a powerful tool for predicting a single quantitative (real-valued) variable based on any number of real-valued predictors or features. A regression model that has only one independent variable is considered to be the simplest one. However, many real-world scenarios involve situations where changes in the dependent variable, denoted as 'y,' are influenced by multiple independent variables. These linear regression models, which use several independent variables to explain variations in the dependent variable, are called numerous regression models. As mentioned earlier, a multiple regression model involves two or more independent variables, as shown in the equation below:

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik} + \varepsilon_i \quad (1.1)$$

By default, linear regression assumes that the relationship between independent and dependent variables is linear and determines the best-fitting line with the least Mean Squared Error (MSE), which measures the difference between the actual data and the predicted values. However, there are instances where the inherent pattern in the data is not sufficiently captured by a linear relationship. In such cases, polynomial regression becomes a more suitable alternative. It can be defined using the following formula:

$$y = \omega_0 + \omega_1 x + \dots + \omega_n x^n \quad (1.2)$$

As the order of the polynomial increases, so does the complexity of the model. The upcoming analysis will focus on exploring datasets that pose a challenge to predict using linear regression. The aim is to design an optimal model that works well in these scenarios. The analysis will also involve an examination of the relationship between multiple independent and dependent variables, with the goal of creating a polynomial regression model.

Chapter 2

Datasets

Advertising Dataset

The dataset contains advertising and sales data seen in class. To improve sales of a particular product some data consisting of the sales of that product in 200 different markets, along with advertising budgets for the product in each of those markets for three various media: TV, radio, and newspaper.

	TV	radio	newspaper	sales
1	230.1	37.8	69.2	22.1
2	44.5	39.3	45.1	10.4
3	17.2	45.9	69.3	9.3
4	151.5	41.3	58.5	18.5
5	180.8	10.8	58.4	12.9

Table 2.1: Advertising Data Sample

	TV	radio	newspaper	sales
count	200.000000	200.000000	200.000000	200.000000
mean	147.042500	23.264000	30.554000	14.022500
std	85.854236	14.846809	21.778621	5.217457
min	0.700000	0.000000	0.300000	1.600000
25%	74.375000	9.975000	12.750000	10.375000
50%	149.750000	22.900000	25.750000	12.900000
75%	218.825000	36.525000	45.100000	17.400000
max	296.400000	49.600000	114.000000	27.000000

Table 2.2: Advertising Data Statistics

Auto Dataset

The auto dataset gives some information about cars, auto gas mileage, horse-power and other information. In particular, the variables are mpg, cylinders, displacement, horsepower, weight, acceleration, year origin, and name.

	mpg	cyl.	disp.	hp	wt.	accel.	yr	org	name
0	18.0	8	307.0	130.0	3504	12.0	70	1	chevrolet chevelle malibu
1	15.0	8	350.0	165.0	3693	11.5	70	1	buick skylark 320
2	18.0	8	318.0	150.0	3436	11.0	70	1	Plymouth satellite
3	16.0	8	304.0	150.0	3433	12.0	70	1	amc rebel sst
4	17.0	8	302.0	140.0	3449	10.5	70	1	ford torino

Table 2.3: Auto Data Sample

	mpg	cyl.	disp.	hp	wt.	accel.	yr	org
count	392	392	392	392	392	392	392	392
mean	23.45	5.47	194.41	104.47	2977.58	15.54	76	1.58
std	7.81	1.71	104.64	38.49	849.40	2.76	3.68	0.81
min	9.0	3	68.0	46.0	1613	8.0	70	1
25%	17.0	4	105.0	75.0	2225.25	13.78	73	1
50%	22.75	4	151.0	93.5	2803.5	15.5	76	1
75%	29.0	8	275.75	126.0	3614.75	17.03	79	2
max	46.6	8	455.0	230.0	5140	24.8	82	3

Table 2.4: Auto Data Statistics

Chapter 3

Methodologies

3.1 Multicollinearity

Linear regression is based on four fundamental assumptions. Firstly, the relationship between the dependent variable and the independent variables should be linear. Secondly, the dependent variables should be normally distributed and have the same variance. Thirdly, the independent variables should be independent of each other.

To create an accurate model, dependent and independent variables in `advertising.csv` and `auto.csv` will be analysed in these four ways. When there is multicollinearity in the data, it becomes difficult to determine the exact impact of each variable on the dependent variable. By removing the columns with high multicollinearity from the base model using heatmaps and the VIF library, a model with less error can be created.

1. Heatmap

It is a visualisation technique that shows the correlation between data using the degree of colour. If you input the table value that calculates the correlation index using the `corr()` function as the input of the heatmap, you can visually check the one-to-one correlation between each variable.

2. VIF(Variance Inflation Factor)

It shows the performance of linear regression of one independent variable to other independent variables, and if the correlation between the dependent variable and other variables is large, the r-squared value of the regression equation is high, and the vif index is inflated.

3.1.1 Data Pre-processing

3.1.1.1 Advertising.csv

The dataset's heatmap and VIF index were analysed to detect the presence of multicollinearity. The results indicate that there was no evidence of multicollinearity between the variables as the correlation between each predictor was low, as shown in Figure 3.1 and Table 3.1.



Figure 3.1: Heat map of advertising dataset

	VIF Factor	Feature
0	2.486772	TV
1	3.285462	radio
2	3.055245	newspaper

Table 3.1: VIF factor

Due to the robust correlation observed between 'sales' and the 'TV' variable and the limited insights offered by the 'newspaper' data, a comparative analysis was conducted to assess the impact of including or excluding the 'newspaper' variable. This investigation encompassed the development of two distinct multiple regression models.

The first model, referred to as a multiple linear regression model, considered all three independent variables individually, as delineated in Equations (3.1) and (3.2). In contrast, the second model employed a multiple regression model with interaction, where the product of two variables was considered as a nonlinear term to account for the interaction effects of each independent variable, as illustrated in Equation (3.3).

$$sales = \theta_0 + \theta_1 * TV + \theta_2 * radio \quad (3.1)$$

$$sales = \theta_0 + \theta_1 * TV + \theta_2 * radio + \theta_3 * newspaper \quad (3.2)$$

$$sales = \theta_0 + \theta_1 * TV + \theta_2 * radio + \theta_3 * TV * radio \quad (3.3)$$

3.1.1.2 Auto.csv

Figure 3.2 shows a correlation between horsepower and mpg of 0.78, which is generally considered a correlation above the absolute value of 0.4. The main aim of this paper is to analyse the relationship between the dependent variable mpg and the independent variable horsepower while excluding other variables from the dataset. To accomplish this, two variables were extracted from the entire dataset, horsepower and mpg, and were utilised as the independent and dependent variables, respectively. The objective of this paper is to investigate the correlation between these two variables.

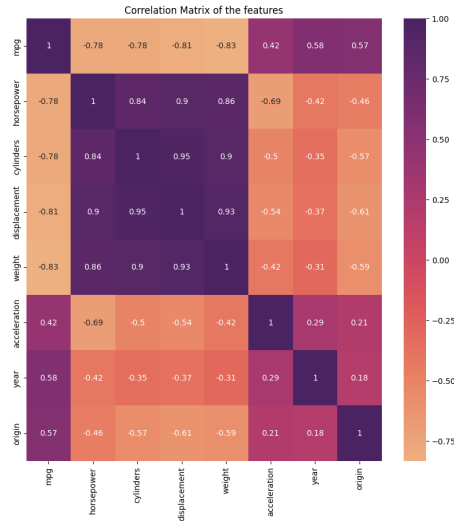


Figure 3.2: Heat map of auto dataset

3.2 Optimal splitting ratio

In the development of the model, a random sample was drawn from the dataset, with the remaining data earmarked for validation. In the context of determining the appropriate allocation of data for training and testing, the optimal train/test split ratio is calculated using the formula $\sqrt{p}:1$, where 'p' denotes the number of parameters requisite for an accurate representation of the linear regression model that encapsulates the underlying data (3.3). The calculation of the ideal training dataset size entails the substitution of the parameter count into the aforementioned formula. (1)

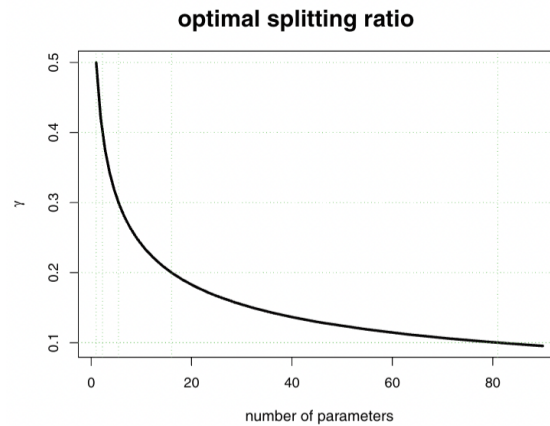


Figure 3.3: Optimal splitting ratio against the number of parameters

For example, for the two variables model, the optimal train size is 0.63, while for the three variables model, it is 0.66.

3.3 Polynomial Regression

Polynomial regression works in the same way as multiple linear regression, except that it preprocesses the data to create new combinations of variables before running the regression. However, the curve being fitted has a quadratic nature.

$$\begin{aligned} Y &= \beta_0 * X^2 + \beta_1 * X + \beta_2 \\ X_1 &= X^2 \\ X_2 &= X \\ Y &= \beta_0 * X_1 + \beta_1 * X_2 + \beta_2 \end{aligned}$$

Therefore, after transforming the training set to include polynomial terms up to a degree, the model was predicted by inheriting the previously implemented model. (See Appendix D and E)

3.4 Gradient Descent

Gradient descent is an algorithm for finding the minimum of the cost function.

The learning rate α in (3.4) and the gradient are multiplied and subtracted to keep moving the parameter towards the cost function's minimum. The learning rate is a constant, and its purpose is to control the amount of change in motion. It is usually between 0 and 1. The optimal model was determined by systematically assessing various combinations of learning rates and iterations. This was done by organising learning rate and iteration values into arrays. (see Appendix F)

The slope is proportional to the distance from the minima, meaning that if you are farther away from the minima, you can move more because the slope is steeper on the graph of the cost function, and if you are closer to the minima, you can move less because the slope is gentler. Subtract the product of the learning rate and the slope because the sign of the derivative is opposite to the direction of travel to the minima.

$$\theta_i = \theta_i - \alpha \frac{\partial J(\theta_0, \theta_i)}{\partial \theta_i} \quad (3.4)$$

The Gradient Descent algorithm is used to optimise multivariable linear regression by finding the optimal theta values that minimise the cost function. An algorithm for the Multivariate Regression Cost Function can be implemented by continuously updating the theta values after each model fits using the Gradient Descent algorithm, as shown in Equation 3.5.

$$\theta_j := \theta_j - \alpha \frac{2}{m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i) x_j^i \quad (3.5)$$

3.4.1 Normalised Dataset

Each of the three variables in the given Advertising dataset has a different range and density distribution. To address the issue of heavily influenced term creation by larger

values when multiplying two independent variables, standardization of the three independent and dependent variables using their mean and standard deviation was necessary. This is because the cost function graph becomes very scattered in three dimensions, and the gradient descent method takes a lot of time when dealing with variables that have a wide range.

(2)

3.4.2 Gradient Clipping

When performing gradient descent, it's crucial to prevent erratic movements of the gradient. These movements often occur when the gradient takes on substantial values. The gradient can then point in misleading directions, which can hinder convergence towards the global minimum.

A prominent solution to this challenge is the method of *gradient clipping* (3).

The core idea behind gradient clipping is to limit the magnitude of the gradient. Formally, given a gradient \mathbf{g} , the condition to be satisfied is:

$$\|\mathbf{g}\| \leq \text{threshold} \quad (3.6)$$

The term *threshold* denotes a predefined maximum magnitude for the gradient. If the magnitude $\|\mathbf{g}\|$ surpasses this threshold, the gradient is scaled such that its magnitude is confined to the threshold, but its original direction is retained.

The use of gradient clipping in a gradient descent procedure effectively minimises the potential risks posed by disproportionately high gradient values, thus avoiding potential overflow scenarios.

3.5 Regression criteria

3.5.1 R squared value

In regression analysis, the value of the coefficient of determination R^2 is a measure of the magnitude of the explanatory power of the independent variables on the total variation of the dependent variable and ranges from 0 to 1. The closer the value of R^2 is to 1, the greater the explanatory power of the dependent variable by the independent variables and the better the fit of the regression equation.

It is imperative to scrutinise the adjusted R-squared value to ascertain whether it has been substantively enhanced if an additional independent variable has been integrated into the model.

$$R_a^2 = 1 - (n - 1) \frac{(1 - R^2)}{n - p - 1} \quad (3.7)$$

In equation 3.7, where p means the number of independent variables, p is in the denominator, and as p increases, the R-squared value in the numerator also increases, somewhat offsetting the impact of the increase.

Chapter 4

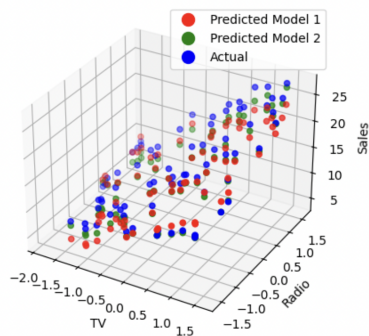
Results & Discussion

4.1 Advertising dataset

4.1.1 Multivariate linear regression

A regression model that considers the relationship between the first two independent variables and the dependent variable is shown in Model 1 (3.1). The first term on the left-hand side of the proposed model is constant, and the next two terms consider one of each independent variable. Model 2(3.3) includes an interaction term formed by multiplying the two independent variables. Considering all the interaction terms may adversely affect the efficiency of the model due to multicollinearity, the VIF and hitmap for all the terms were performed to derive the final multiple regression model function. The results are shown in Figure 4.1 The R-squared value for the multiple linear regression model without interactions is 0.816, but the R-squared value for the multiple regression model with interactions is 0.890, a significant improvement.

3D Scatter plot of Sales depending on TV and Radio budgets - Model 1 and :

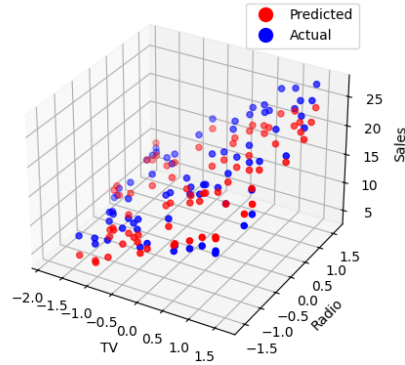


	Model1	Model2
R^2 score	0.816	0.890
MSE	5.34	2.79
RMSE	2.31	1.67

Figure 4.1: Fitting the data with custom Model

To fit a linear regression model using Sklearn, a three-step process was followed. Firstly, a linear regression class instance was created. Subsequently, the model was trained by invoking the fit() function. Finally, the predict() function was called to generate the predictions of the trained model.

3D Scatter plot of Sales depending on TV and Radio budgets - Scikit-Learn

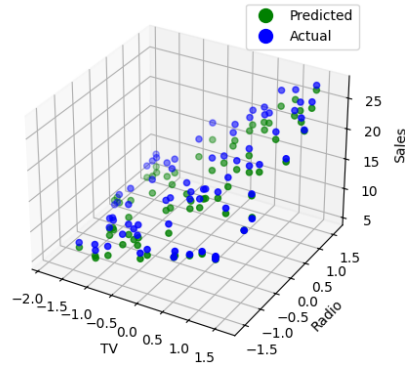


	Sklearn
R^2 score	0.83
MSE	4.38
RMSE	2.09

Figure 4.2: Fitting the data with sklearn Model

In Model2 equation 3.3, θ_3 represents the coefficient for this interaction term. The magnitude and sign of θ_3 will determine how the interaction between TV and radio affects the predicted values of sales. If θ_3 is positive, the interaction has a positive effect on the radio, and if it's negative, the interaction has a negative effect. The model yields a value of approximately 1.37 for θ_3 , indicating the mutual influence of the two independent variables.

3D Scatter plot of Sales depending on TV and Radio budgets - Model 2



	Coefficient
θ_0	13.41
θ_1	3.92
θ_2	2.67
θ_3	1.37

Figure 4.3: TV+radio+TV*radio

4.1.2 TV + radio + newspaper

The model's predictive capability was extended through the introduction of an additional independent variable, 'newspaper,' into the existing framework, yielding the following regression equation:

$$sales = 13.5851 + 3.5680 * TV + 2.8118 * radio + 0.2626 * newspaper \quad (4.1)$$

It is noteworthy that the values of the coefficients in Equation (5.1) may exhibit slight variability in response to fluctuations in the training and test data. Nonetheless, the model can be deemed valid, particularly in consideration of the R-squared value. It is discernible from this model that 'TV' exerts the most substantial influence on the sales outcome, while 'newspaper' demonstrates a relatively minor effect. Despite the increment in the R-squared value, it is imperative to scrutinise the adjusted R-squared value to ascertain whether it has been substantively enhanced, as this additional independent variable has been integrated into the model.

The adjusted R-square value is 0.86, meaning that the predicted value calculated from the regression can explain 86% of the actual y-value.

To conduct a comprehensive analysis, advertising budgets for the product and data for three different media formats, namely TV, radio, and newspaper, were collected from 200 markets. These variables were utilised as independent factors in the analysis, enabling the identification of the impact of media formats on the product's performance in each market. This approach provided a deeper understanding of the relationship between advertising expenditure and media platforms. To account for the interaction between the independent variables, a multiple regression model was built that additionally considers a term consisting of the product of the two independent variables. The analysis showed that significant p-values were obtained for TV and radio, indicating that the null hypothesis about the association between these two variables and sales can be rejected. However, the null hypothesis for newspaper cannot be rejected, suggesting that there is insufficient evidence to conclude that this variable is significantly associated with sales.

	coef	std err	t	P> t	[0.025	0.975]
Intercept	2.9389	0.312	9.422	0.000	2.324	3.554
TV	0.0458	0.001	32.809	0.000	0.043	0.049
radio	0.1885	0.009	21.893	0.000	0.172	0.206
newspaper	-0.0010	0.006	-0.177	0.860	-0.013	0.011

Figure 4.4: P-value of variables

When it comes to advertising, all expenditures are positively correlated with sales. However, it is highly unlikely that newspaper ad expenditures have any correlation with sales. When regressing multiple variables, overlapping variances explaining sales are often not accounted for, resulting in less variance and a lower regression coefficient for newspapers.

The results of Model 3, which includes newspaper data as an independent variable in Model 1, are shown in Table 4.3.

	model	sklearn
RMSE	1.98	2.09
R^2	0.87	0.83

Table 4.1: TV+radio+newspaper data

4.1.3 Overall

Figure 4.5 shows the relationship between sales and ad budget, as plotted using the Spearmanr library. The relationship appears to be roughly linear, indicating that increasing the ad budget leads to a corresponding increase in sales. Furthermore, the parameter values for TV and radio are larger than those for newspapers, suggesting that increasing the budget for these two media channels is more effective in boosting sales than for newspapers.

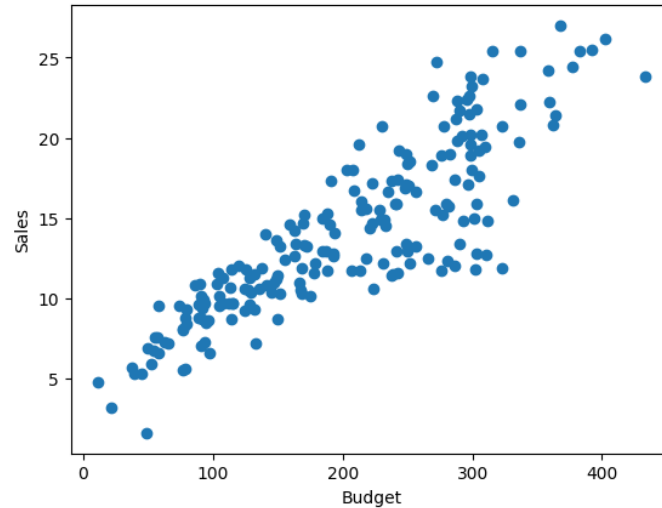


Figure 4.5: Correlation between Budget and Sales

After comparing the obtained model with the model obtained using the Sklearn library (Figure 4.2 and Table 4.1), it was found that the model produces similar results to the Sklearn function.

4.2 Auto dataset

Assuming a linear relationship between horsepower and mpg, the model was predicted: $Y = \beta_0 + \beta_1 X_1$.

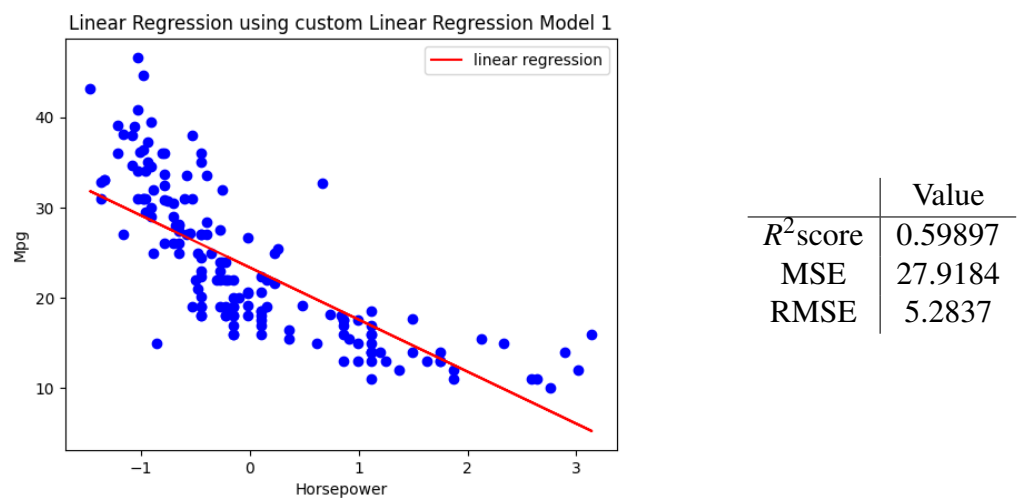


Figure 4.6: Linear Regression

In **Fig 5.5**, the straight line is unable to capture the pattern in the data. This is categorised as an under-fitting. The R squared value is 0.59, which indicates that linear regression is inadequate to explain the model.

A higher-order equation can be generated by adding powers of the original features as new features. The linear model can then be transformed to $Y = \beta_0 X^2 + \beta_1 X + \beta_2$.

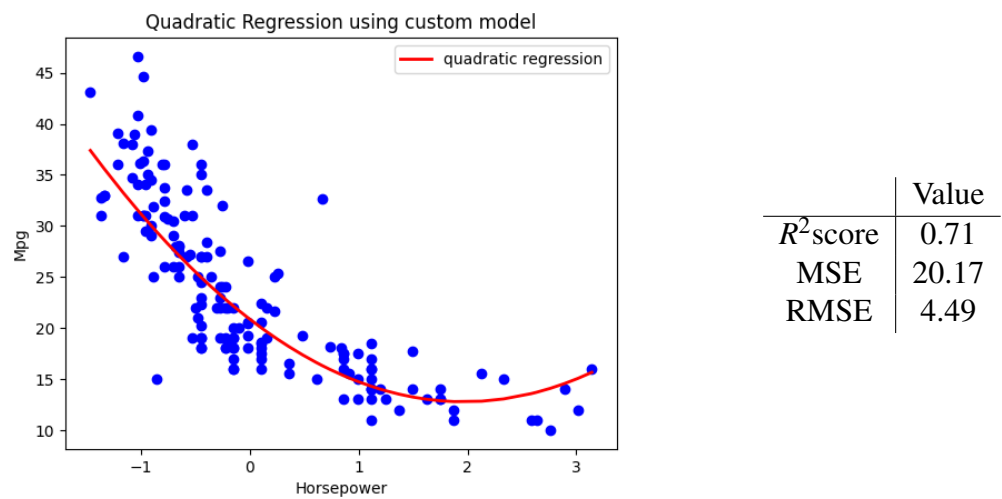
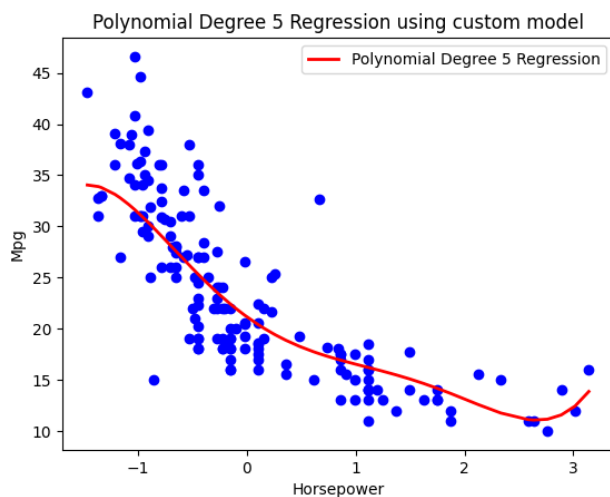


Figure 4.7: Quadratic Regression

Figure 5.6 illustrates the graph and numbers obtained through the quadratic regression model, which fits the data better than the linear curve. Comparing the R^2 values, an increase in the model’s explained probability to 0.71 is evident. Additionally, the RMSE has decreased, and the R^2 score has increased when compared to the linear line.

When the degree value is set to 5, **Figure 5.7** shows that the model passes through more points and finds a better fit. However, it tends to overfit and performs worse than polynomial regression when data other than the training data is included.



	Value
R^2 score	0.71
MSE	20
RMSE	4.47

Figure 4.8: Degree 5 Regression

Table 5.1 and **Table 5.2** present the results obtained from the implemented model and the sklearn library, respectively. It appears that the models produced using the library are generally better than those produced using the implemented model.

	linear	quadratic	degree 5
RMSE	5.28	4.49	4.47
R^2	0.59	0.71	0.71

Table 4.2: Polynomial regression results for custom model

	linear	quadratic	degree 5
RMSE	5.28	4.58	4.5
R^2	0.59	0.69	0.70

Table 4.3: Polynomial regression results for sklearn library

4.3 Conclusion

When working with complex data sets in data science, it can be challenging to pinpoint the most crucial variables. While machine learning techniques can leverage multiple variables to improve predictions, this can result in overfitting. As a result, it is often more effective to focus on the most significant variables. For instance, Multicollinearity was identified in the sales dataset, and only the essential variables were selected for modelling.

Furthermore, we utilised a basic gradient descent method to construct the model, which can be implemented in various ways, including mini-batch and stochastic. It's worth noting that initialising values to zero can significantly impact training time and results since the derivative is a crucial algorithm component. Once the slope value is determined, the algorithm adjusts the point's position by adding or subtracting it to move closer to the minimum.

References

1. Joseph VR. Optimal ratio for data splitting. Statistical Analysis and Data Mining: The ASA Data Science Journal. 2022;15(4). Available at: <https://onlinelibrary.wiley.com/doi/full/10.1002/sam.11583>. © 2022 The Author. Statistical Analysis and Data Mining published by Wiley Periodicals LLC.
2. Watt J, et al. Machine Learning Refined: Foundations, Algorithms, and Applications 2nd Edition. Cambridge University Press; March 12, 2020.
3. Bajaj A. Understanding Gradient Clipping (and How It Can Fix Exploding Gradients Problem); 2023. Available at: <https://neptune.ai/blog/understanding-gradient-clipping-and-how-it-can-fix-exploding-gradients-problem>. Accessed: 11 October 2023.

Appendix A

Documentation

1. Project tree

```
data /
  - Advertising.csv
  - Auto.csv
library /
  - LinearRegressionCustomModel1.py
  - LinearRegressionCustomModel2.py
  - PolynomialDegree5RegressionCustomModel.py
  - QuadraticRegressionCustomModel.py
  - utils.py
assessment_advertising-dataset.ipynb
assessment_auto-dataset.ipynb
requirements.txt
```

2. Getting Started

To run the program, follow these steps:

1. Create a virtual environment using `python3 -m venv venv`.
2. Activate the virtual environment using `source venv/bin/activate`.
3. Install the required dependencies using `pip3 install -r requirements.txt`.
4. Use the jupyter notebook `assessment_advertising-dataset.ipynb` to run the program on the advertising dataset.
5. Use the jupyter notebook `assessment_auto-dataset.ipynb` to run the program on the auto dataset.

3. Detailed Features of Classes and Functions

Classes

`LinearRegressionCustomModel1`: Class to perform linear regression with the gradient descent algorithm, using the model:

$$y = \theta_0 + \theta_1 \times x_1 + \theta_2 \times x_2 + \dots + \theta_n \times x_n$$

`LinearRegressionCustomModel2`: Class to perform linear regression with the gradient descent algorithm, using the model:

$$y = \theta_0 + \theta_1 \times x_1 + \theta_2 \times x_2 + \theta_3 \times (x_1 \times x_2)$$

`QuadraticRegressionCustomModel`: Extends the `LinearRegressionCustomModel1` Class to perform quadratic regression.

`PolynomialDegree5RegressionCustomModel`: Extends the `LinearRegressionCustomModel1` Class to perform polynomial regression of degree 5.

Each class has the following methods:

- `fit(self, x_train, y_train)`: Method to fit the model to the training data.
- `predict(self, x_test)`: Method to predict the output for a given input.
- `metrics(self, x_test, y_test)`: Method to compute the model's metrics.

Functions in `utils.py`

`computeOptimalTrainSize(r)`: Function to compute the optimal training dataset size based on the number `r` of predictors.

`computeOptimalParameters(x_train, x_test, y_train, y_test, model_type)`: Function to compute the best hyperparameters (learning rate and number of steps).

`prepare_data(x, y, train_size)`: Function to prepare the data (split the data into training and testing sets, and standardise the data).

Appendix B

Linear Regression Custom Model 1 Source Code

```
1 import numpy as np
2 from sklearn import metrics
3
4
5 # Class to perform linear regression with gradient descent
  algorithm - Model 1
6 #  $y = \_0 + \_1 * x_1 + \_2 * x_2 + \dots + \_n * x_n$ 
7 class LinearRegressionCustomModel1:
8     def __init__(self, learning_rate, iterations_nb):
9         self.learning_rate = learning_rate # Learning rate
10        self.iterations_nb = iterations_nb # Number of iterations
11        self.theta = None
12
13        # Method to fit the model to the training data
14        def fit(self, x_train, y_train):
15            # Add column of ones to x_train - Bias / Intercept
16            x_train = np.column_stack((np.ones(x_train.shape[0]),
17 x_train))
18
19            # Number of training samples and number of features
20            n, r = x_train.shape
21
22            # Initialise theta with a vector of zeros
23            self.theta = np.zeros(r)
24
25            # Initialise variables for convergence check
26            prev_cost = float("inf") # Set initial previous cost to
infinity
27            tolerance = 1e-6 # Convergence tolerance
28            gradient_threshold = 1e5 # Gradient clipping threshold
29
30            # Gradient descent algorithm
31            for _ in range(self.iterations_nb):
32                # Compute predictions ( $y_{pred} = x_{train} * \theta$ )
33                y_pred = x_train.dot(self.theta)
34
35                # Compute error (difference between predictions and
actual values)
```



```

35         error = y_pred - y_train
36
37         # Compute gradient
38         gradient = (2 / n) * error.dot(x_train)
39
40         # Gradient clipping to prevent divergence (if gradient
is too large)
41         gradient = np.clip(gradient, -gradient_threshold,
gradient_threshold)
42
43         # Update parameters
44         self.theta -= self.learning_rate * gradient
45
46         # Compute cost (mean squared error)
47         cost = (1 / n) * np.sum((y_pred - y_train) ** 2)
48
49         # Convergence check
50         if abs(prev_cost - cost) < tolerance:
51             break
52
53         # Update previous cost
54         prev_cost = cost
55
56         # Return None if training was not successful
57         if np.isnan(cost):
58             return None, float("inf") # Return None for theta and
infinity for cost
59
60         return self.theta, cost
61
62     # Method to predict the output for a given input
63     def predict(self, x_test):
64         # Add column of ones to x_test - Bias / Intercept
65         x_test = np.column_stack((np.ones(x_test.shape[0]), x_test)
)
66
67         # Return predictions with interaction term
68         y_pred = x_test.dot(self.theta) # Compute predictions
69         assert (
70             y_pred.shape[0] == x_test.shape[0]
71             ), "Mismatch in prediction shape" # Check if the shape of
the prediction is correct
72
73         return y_pred
74
75     # Method to compute the metrics of the model
76     def metrics(self, x_test, y_test):
77         # Compute predictions
78         y_pred = self.predict(x_test)
79
80         # Compute metrics
81         rSquare = metrics.r2_score(y_test, y_pred) # R squared
82         meanAbErr = metrics.mean_absolute_error(y_test, y_pred) #
Mean absolute error
83         meanSqErr = metrics.mean_squared_error(y_test, y_pred) #
Mean squared error

```

```
84     rootMeanSqErr = np.sqrt(  
85         metrics.mean_squared_error(y_test, y_pred)  
86     ) # Root mean squared error  
87  
88     return rSquare, meanAbErr, meanSqErr, rootMeanSqErr
```

Appendix C

Linear Regression Custom Model 2 Source Code

```
1 import numpy as np
2 from sklearn import metrics
3
4
5 # Class to perform linear regression with gradient descent
  algorithm - Model 2
6 #  $y = \_0 + \_1 * x1 + \_2 * x2 + \_3 * (x1 * x2)$ 
7 class LinearRegressionCustomModel2:
8     def __init__(self, learning_rate, iterations_nb):
9         self.learning_rate = learning_rate # Learning rate
10        self.iterations_nb = iterations_nb # Number of iterations
11        self.theta = None
12
13    # Method to fit the model to the training data
14    def fit(self, x_train, y_train):
15        # Convert DataFrame to numpy array
16        x_train = x_train.values
17
18        # Add column of ones to x_train - Bias / Intercept
19        interaction_term = (
20            x_train[:, 0] * x_train[:, 1]
21        ) # Interaction term between x1 and x2
22        x_train = np.column_stack(
23            (np.ones(x_train.shape[0]), x_train, interaction_term)
24        )
25
26        # Number of training samples and number of features
27        n, r = x_train.shape
28
29        # Initialise theta with a vector of zeros
30        self.theta = np.zeros(r)
31
32        # Gradient descent algorithm
33        for _ in range(self.iterations_nb):
34            # Compute predictions
35            y_pred = x_train.dot(self.theta)
36
37            # Compute error (difference between predictions and
```

```

actual values)
38     error = y_pred - y_train
39
40     # Compute gradient
41     gradient = (2 / n) * error.dot(x_train)
42
43     # Update parameters
44     self.theta -= self.learning_rate * gradient
45
46     # Compute cost (mean squared error)
47     cost = (1 / n) * np.sum((y_pred - y_train) ** 2)
48
49     return self.theta, cost
50
51 # Method to predict the output for a given input
52 def predict(self, x_test):
53     # Add column of ones to x_test - Bias / Intercept
54     x_test = np.column_stack((np.ones(x_test.shape[0]), x_test)
55 )
56
57     # Return predictions with interaction term
58     y_pred = (
59         self.theta[0]
60         + self.theta[1] * x_test[:, 1]
61         + self.theta[2] * x_test[:, 2]
62         + self.theta[3] * (x_test[:, 1] * x_test[:, 2])
63     )
64     assert (
65         y_pred.shape[0] == x_test.shape[0]
66     ), "Mismatch in prediction shape" # Check if the shape of
the prediction is correct
67
68     return y_pred
69
70 # Method to compute the metrics of the model
71 def metrics(self, x_test, y_test):
72     # Compute predictions
73     y_pred = self.predict(x_test)
74
75     # Compute metrics
76     rSquare = metrics.r2_score(y_test, y_pred) # R squared
77     meanAbErr = metrics.mean_absolute_error(y_test, y_pred) #
Mean absolute error
78     meanSqErr = metrics.mean_squared_error(y_test, y_pred) #
Mean squared error
79     rootMeanSqErr = np.sqrt(
80         metrics.mean_squared_error(y_test, y_pred)
81     ) # Root mean squared error
82
83     return rSquare, meanAbErr, meanSqErr, rootMeanSqErr

```

Appendix D

Quadratic Regression Custom Model Source Code

```
1 from library.LinearRegressionCustomModel1 import
   LinearRegressionCustomModel1
2 import numpy as np
3
4
5 # Extend the LinearRegressionCustomModel1 class to implement a
   quadratic regression model
6 class QuadraticRegressionCustomModel(LinearRegressionCustomModel1):
7     # Function to add quadratic terms to the dataset
8     def __add_quadratic_terms(self, X):
9         # Square each feature
10         X_squared = np.square(X)
11
12         # Concatenate the original features and their squared
   values
13         return np.hstack((X, X_squared))
14
15     def fit(self, x_train, y_train):
16         # Transform x_train to include quadratic terms
17         x_train_transformed = self.__add_quadratic_terms(x_train)
18
19         # Call the fit method of the parent class (
   LinearRegressionCustomModel1)
20         theta, cost = super().fit(x_train_transformed, y_train)
21
22         return theta, cost
23
24     def predict(self, x_test):
25         # Transform x_test to include quadratic terms
26         x_test_transformed = self.__add_quadratic_terms(x_test)
27
28         # Call the predict method of the parent class
29         return super().predict(x_test_transformed)
```

Appendix E

Degree 5 Regression Custom Model Source Code

```
1 from library.LinearRegressionCustomModel1 import
   LinearRegressionCustomModel1
2 import numpy as np
3
4
5 # Extend the LinearRegressionCustomModel1 class to implement a
   polynomial regression model of degree 5
6 class PolynomialDegree5RegressionCustomModel(
   LinearRegressionCustomModel1):
7     # Add polynomial terms up to degree 5 to the dataset
8     def __add_degree_5_terms(self, X):
9         # Square each feature
10         X_deg_2 = np.power(X, 2)
11         X_deg_3 = np.power(X, 3)
12         X_deg_4 = np.power(X, 4)
13         X_deg_5 = np.power(X, 5)
14
15         # Concatenate the original features and their polynomial
   values
16         return np.hstack((X, X_deg_2, X_deg_3, X_deg_4, X_deg_5))
17
18     def fit(self, x_train, y_train):
19         # Transform x_train to include polynomial terms up to
   degree 5
20         x_train_transformed = self.__add_degree_5_terms(x_train)
21
22         # Call the fit method of the parent class (
   LinearRegressionCustomModel1)
23         theta, cost = super().fit(x_train_transformed, y_train)
24
25         return theta, cost
26
27     def predict(self, x_test):
28         # Transform x_test to include polynomial terms up to degree
   5
29         x_test_transformed = self.__add_degree_5_terms(x_test)
30
31         # Call the predict method of the parent class
```

```
32     return super().predict(x_test_transformed)
```

Appendix F

Utility Functions Source Code

```
1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import train_test_split
4
5 from library.LinearRegressionCustomModel1 import
   LinearRegressionCustomModel1
6 from library.LinearRegressionCustomModel2 import
   LinearRegressionCustomModel2
7 from library.PolynomialDegree5RegressionCustomModel import (
8     PolynomialDegree5RegressionCustomModel,
9 )
10 from library.QuadraticRegressionCustomModel import
   QuadraticRegressionCustomModel
11
12
13 # Function to compute the optimal train size
14 def computeOptimalTrainSize(r):
15     # Compute optimal train size
16     optimalTestSize = 1 / (np.sqrt(r) + 1)
17     optimalTrainSize = 1 - optimalTestSize
18
19     return optimalTrainSize
20
21
22 # Function to compute the optimal parameters depending on the model
   type
23 def computeOptimalParameters(x_train, x_test, y_train, y_test,
   model_type):
24     # Set hyperparameters
25     iterations_numbers = [
26         200,
27         300,
28         400,
29         500,
30         1000,
31         2000,
32         3000,
33         5000,
34         7000,
35         10000,
```



```

36         20000,
37         50000,
38     ] # Number of iterations
39     learning_rates = [0.1, 0.01, 0.05, 0.001, 0.005, 0.0001,
40                       0.0005] # Learning rate
41
42     # Initialise an empty array to store the results
43     results = []
44
45     # Set seed for reproducibility
46     np.random.seed(0)
47
48     # Computes linear regression for each hyperparameters
49     for iterations_nb in iterations_numbers:
50         for learning_rate in learning_rates:
51             # Initialise model
52             if model_type == "linear_2":
53                 # Linear Model 2
54                 model = LinearRegressionCustomModel2(learning_rate,
55                                                       iterations_nb)
56             elif model_type == "quadratic":
57                 # Quadratic model
58                 model = QuadraticRegressionCustomModel(
59                     learning_rate, iterations_nb)
60             elif model_type == "degree_5":
61                 # Polynomial degree 5 model
62                 model = PolynomialDegree5RegressionCustomModel(
63                     learning_rate, iterations_nb
64                 )
65             else:
66                 # Model 1
67                 model = LinearRegressionCustomModel1(learning_rate,
68                                                       iterations_nb)
69             # Fit model to training data
70             theta, cost = model.fit(x_train, y_train)
71             # Skip if training was not successful
72             if theta is None:
73                 continue
74             # Compute metrics
75             rSquare, meanAbErr, meanSqErr, rootMeanSqErr = model.
76             metrics(x_test, y_test)
77             # Save results
78             results.append(
79                 {
80                     "iterations_nb": iterations_nb,
81                     "learning_rate": learning_rate,
82                     "cost": cost,
83                     "rSquare": rSquare,
84                     "meanAbErr": meanAbErr,
85                     "meanSqErr": meanSqErr,
86                     "rootMeanSqErr": rootMeanSqErr,
87                 }
88             )
89
90     # Convert results to dataframe
91     results = pd.DataFrame(results)

```

```
87
88     # Sort results by MSE
89     results = results.sort_values("meanSqErr", ascending=True)
90
91     # Return best hyperparameters
92     return (
93         int(results.iloc[0]["iterations_nb"]),
94         results.iloc[0]["learning_rate"],
95         results,
96     )
97
98
99 # Function to prepare data (split data into train and test sets,
    and standardise data)
100 def prepare_data(x, y, train_size):
101     # Split data into train and test sets
102     x_train, x_test, y_train, y_test = train_test_split(x, y,
103                                                         train_size=train_size)
104
105     # standardise data
106     x_train = (x_train - x_train.mean()) / x_train.std() #
    standardise train set
107     x_test = (x_test - x_test.mean()) / x_test.std() # standardise
    test set
108
109     # Return train and test sets
110     return x_train, x_test, y_train, y_test
```