

**A RESEARCH PROJECT ON A MODEL THAT PREDICTS CAR PURCHASE
AMOUNT USING ARTIFICIAL NEURAL NETWORKS.**



TECHNICAL UNIVERSITY OF KENYA

Education and Training for the Real World

THE TECHNICAL UNIVERSITY OF KENYA

FACULTY OF APPLIED SCIENCES AND TECHNOLOGY

SCHOOL OF MATHEMATICS AND ACTUARIAL SCIENCE

DEPARTMENT OF PURE AND APPLIED MATHEMATICS

NAME: BONIFACE KARANJA MACHARIA

REG NO. SMPQ/00849/2018

**A Research Project Submitted in Partial Fulfillment of the Requirements for
the Award of the Degree in Bachelor of Science Mathematics.**

DECLARATION

I declare that this project is my original work and has not been presented for a degree award in any other university for the award of a BSc Mathematics degree

SIGNATURE..... DATE.....

NAME: BONIFACE MACHARIA

REG. No. SMPQ/00849/2018

This project has been submitted for examination with my approval as a university supervisor.

SIGNATURE..... DATE.....

NAME: DR. KIMUNGUYI JOSEPHS

DEPARTMENT OF PURE AND APPLIED MATHEMATICS

THE TECHNICAL UNIVERSITY OF KENYA



TECHNICAL UNIVERSITY OF KENYA

DEDICATION

I dedicate this work to my parents Mr. Stephen Macharia and Mrs. Virginia Wambui, my siblings, my supervisor Dr. Kimunguyi josephs, and my friends Reuben Mwangi, Abraham Musembi, Maureen Githaiga and Xaviour Aluku for their continued support, both morally as well as financially. I will forever be grateful for your undying love and support. This project could not have been a success without your most valued inputs. I will never let you down.



ACKNOWLEDGEMENT

First, I want to acknowledge God for His continued blessings, health, strength, and grace even as I maneuvered and navigated through the project and the entire coursework. Secondly, I wish to appreciate my entire family for great support and encouragement, my friends for being there for me, and for guidance and advice as I pursued the project.

I would also like to acknowledge Xavier Aluku (CEO AFRICDSA) for allowing me to gain additional skills in data science and machine learning with python which has helped me maneuver through the project.

Finally, I acknowledge my supervisor DR. Kimungunyi for walking with me all through the project correcting, guiding, and encouraging me and the entire teaching staff of the department of pure and applied mathematics at The Technical University of Kenya.



Table of Contents

A Research Project Submitted in Partial Fulfillment of the Requirements for the Award of the Degree in Bachelor of Science Mathematics.	0
DECLARATION.....	1
DEDICATION	2
ACKNOWLEDGEMENT	3
List of figures.....	6
LIST OF LINKS.....	6
List of Abbreviations	6
ABSTRACT.....	7
CHAPTER 1.....	8
INTRODUCTION.....	8
Background of study	8
Statement of problem	9
Significance of study.....	9
CHAPTER 2.....	10
2.0 LITERATURE REVIEW	10
CHAPTER 3.....	12
3.1 MACHINE LEARNING	12
3.2 LINEAR REGRESSION.....	12
3.3 ARTIFICIAL NEURAL NETWORKS(ANNs)	13
3.3.1 Neuron mathematical model	13
3.3.2 multi neuron model.....	14
3.4 Scaling the data	16
CHAPTER 4.....	17
4.1 DATA COLLECTION.....	17
4.2 IMPORTING THE LIBRARIES	17
4.3 DATA VIEW.....	17
4.4 DATA PREPROCESSING.....	18
4.5 EXPLORATORY DATA ANALYSIS.....	19
4.6 LINEAR REGRESSION	20
4.7 MODEL TRAINING-.....	20
4.7.1 DATA NORMALIZATION	21
4.7.2 TRAIN, TEST AND SPLIT TECHNIQUE	21
4.7.3 MODEL EVALUATION.....	22

4.74 MODEL RESULTS	23
CHAPTER 5.....	24
5.1 CONCLUSION	24
5.2 Appendices.....	25
REFERENCE.....	29



TECHNICAL UNIVERSITY OF KENYA

List of figures

- 1.** Fig 1.0-A graph showing the trendlines of car prices as of Sept 2021
- 2.** Fig 2.0- figure showing the 9 columns and classification of content data types
- 3.** Fig 3.0- figure showing the descriptive statistics of our data and its quartiles
- 4.** Fig 4.0- figure showing the correlation heatmap of different variables
- 5.** Fig 5.0 – figure showing the regression scatter plot between annual salary and car purchase amount
- 6.** Fig 6.0- figure showing the table of rows used as x-variable
- 7.** Fig 7.0- figure showing the graph of normal distribution
- 8.** Fig 8.0 figure showing the graph of model loss during training

LIST OF LINKS

- 1.** https://drive.google.com/file/d/1z8TodtkUYcbyaaxH_5TFgcoOjoFuEEQk/view?usp=sharing- link to the car purchase dataset at Kaggle.com
- 2.** https://en.wikipedia.org/wiki/Sigmoid_function



List of Abbreviations

1. CEO- Chief Executive officer.
2. AFRICDSA- African Center for Data Science and Analytics.
3. ANN- Artificial Neural Network.
4. KNN- K-nearest neighborhood

ABSTRACT

With the increased improvement in technology and infrastructure such as roads, the need for vehicles has increasingly grown bigger over time. People have different reasons for buying a car some of them including but not limited to; comfort, simplicity of movement, entertainment, sporting, as well as business. There are so many factors that determine the choice of a car some of them being price, purpose, availability in the market, and size of the car. A customer must consider all these factors before buying a car. The price of cars has also exponentially gone high with the rise in demand. With the internet nowadays it is easy to purchase vehicles from the comfort of your home.

This research is focused on how we can assist a vehicle dealer company estimate or predict the amount of money a customer might be willing to pay for a car based on their historical data . This data tells us more about the financial credibility of a customer. With this data a dealer might be able to decide if the customer perhaps can be entrusted to pay for the car in installments. The research has also integrated some skills in data science and machine learning with python to help in coming up with the model, training the data, and testing the model.

From this research, we came up with a conclusion that the model is capable of predicting the purchase amount of a vehicle by putting into consideration the customers' historical and financial data. This model can then be fed into a dealer's website where customers can access it in form of a questionnaire where they can fill in their details and the model predict the credit score and prints the amount of money they can have at their disposal before proceeding to check for the brand or type of vehicle that fits the amount. Dealers can also use it as a source of their data when making decisions on what type of cars to restock or modify.

CHAPTER 1

INTRODUCTION

Improvement in technology and infrastructure has led to a great increase in demand for vehicles in day-to-day activities. The price of cars going exponentially lower and quite affordable. In simple terms we are moving into an era where having a vehicle will be a basic need. Vehicles are one of the greatest inventions that have made human life quite flexible and also improved mobility on land.

A study has shown that as of July 2021, the price of cars has undergone a peak with a margin of about 40% mainly due to covid 19 pandemic.

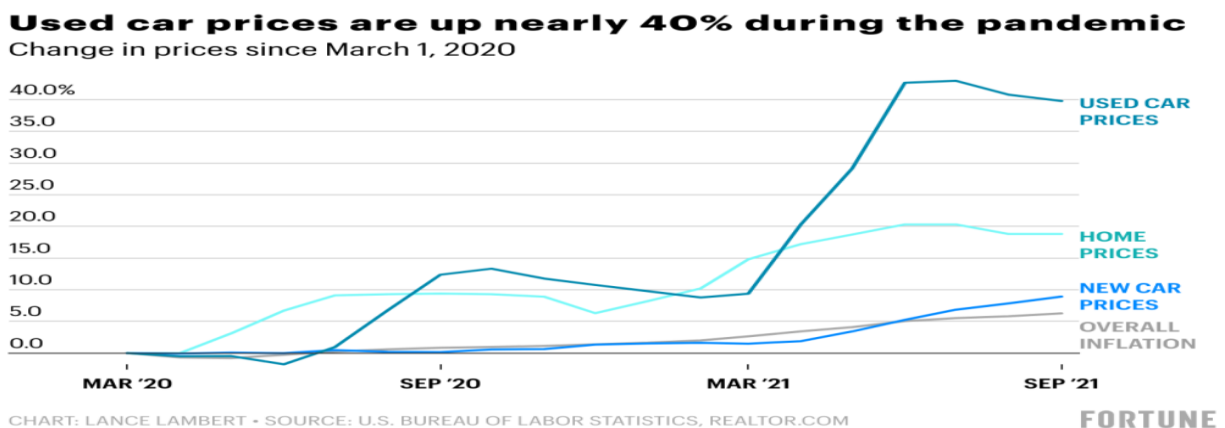


Fig 1.0-A graph showing the trendlines of car prices as at Sept 2021

This has then sparked the idea of narrowing down the search further by introducing a model that can be able to use patterns in data to predict vehicle purchase amounts with similar characteristics and then displaying them to customers from the comfort of their homes

Background of study

The increased demand for cars has led to the emergence of many car dealers. The hunt for a suitable and affordable car has therefore become one of the most tedious activities before you could get one. The Internet has however played a key role to help make the situation

easier for the customers to locate the buyers, access their locations, or even engage them to help them identify a car that fits their budget.

This model further simplifies the hunt by making sure that customers know their credit score thus avoiding overspending. Vehicle dealers also have access to customer data and are able to make informed decisions from the data.

Statement of problem

Suppose I am working as a car salesman and I would like to develop a model to predict the total dollar amount that customers are willing to pay given the following attributes:

- Customer Name
- Customer e-mail
- Country
- Gender
- Age
- Annual Salary
- Credit Card Debt
- Net Worth



TECHNICAL UNIVERSITY OF KENYA

The model should predict:

- Car Purchase Amount

Significance of study

This research focuses entirely on making a model that can be able to predict car purchase amount using linear regression and customer historical data. The model can then be fed into a dealer's website so that dealers can be able to predict how much a customer might be able to pay for a car based on their credit card worthiness, annual salary and gender.

CHAPTER 2

2.0 LITERATURE REVIEW

Wu et al (2009) performed an analysis of car price estimation utilizing a knowledge-based neuro-fuzzy method. The model took into consideration the following characteristics: model, year of production, and engine size. The model of projection had comparable results to the simplistic model of regression. A regression model focused on a neighboring KNN machine learning algorithm was used to predict a car's speed. This program seems to be remarkably effective as it exchanged more than 2 million vehicles.

Richardson (2009) found a more specific approach. In his thesis research, he expected that more robust cars should be made by automakers. He implemented multiple regression analyses and found out that electric cars have maintained their worth longer than regular cars. This was due to urban warming effects and offers greater fuel efficiency.

Gonggie (2011) suggested a model that would be developed using ANN (Artificial Neural Networks) to predict car prices. He considered several attributes such as millage, estimated car life, and mark. The model was developed to cope with non-linear data interactions, which was not the case with prior models using standard linear regression techniques. This model was able to predict car prices exhibiting better accuracy.

Noor and Jan (2017) suggested the use of multiple linear regression for forecasting car prices. Notably, a standard approach to machine learning algorithms did not produce impressive predictive outcomes. This could be improved by combining multiple machine learning methods into an ensemble.

The production of cars has exponentially increased over the past decade and over 20 million passenger cars being produced every year (2021) has given rise to the growth of the used cars market and also with the introduction of trade-in which has become a booming industry. The recent advent of online portals has facilitated the need for both the customer and the seller to be informed about the trends and patterns in determining the values of cars in the market. Using machine learning techniques, we can be able to come up with a model that can predict the price of a car based on consumer data and a given set of features

In this particular research project, a considerable number of distinct objects are examined to come up with an accurate prediction. The data used in this model were collected from the web, and respective algorithms were also compared to find the best suit for the available project. The model was then evaluated using test data and the percentage of accuracy was obtained. Overfitting and overfitting come into the picture when we are creating a statistical model. The model might be too biased toward the training data and might not perform well on the test data. This is called overfitting. Similarly, the model might not take into consideration all the variance in the data and hence perform poorly on the test data set. This is called underfitting. A perfect balance needs to be achieved in both regression and classification.



CHAPTER 3

METHODOLOGY AND PROBLEM SOLVING

3.1 MACHINE LEARNING

Machine learning is the use and development of computer systems that can learn and adapt without following explicit instructions, by using algorithms, mathematical and statistical models to analyze and draw inferences from patterns in data. Customer needs for a car exhibit similar patterns that a computer can adapt and learn and be able to use these inferences in the future to predict the price of a vehicle without having to go through all the iterative processes. Similarly, customers' historical data also exhibit similar trends which can be adapted by a machine learning model to predict and to draw different inferences about a particular customer.

Machine learning algorithms can be used to come up with a model that can be able to study patterns in customer data and predict the amount of money in USD that a customer might be able to dispose to quench his or her need for a car.

3.2 LINEAR REGRESSION

Our machine learning model particularly uses simple linear regression to predict the car purchase amount. We used the following equation

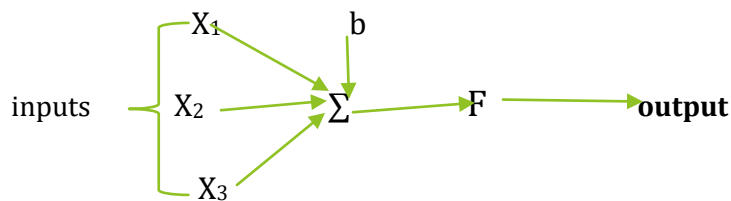
$$Y = aX + b$$

- Where Y is the dependent variable which also signifies the car purchase amount (predictor variable)
- X is the independent variable
- 'a' is the slope
- 'b' is the intercept

3.3 ARTIFICIAL NEURAL NETWORKS(ANNs)

3.3.1 Neuron mathematical model

An artificial neuron perception is a simple mathematical model that involves biological neurons. The neuron collects information from an input channel, process it and predict an output. it can be represented as follows. It contains three subsections. The input section, The summation function, and the activation function.

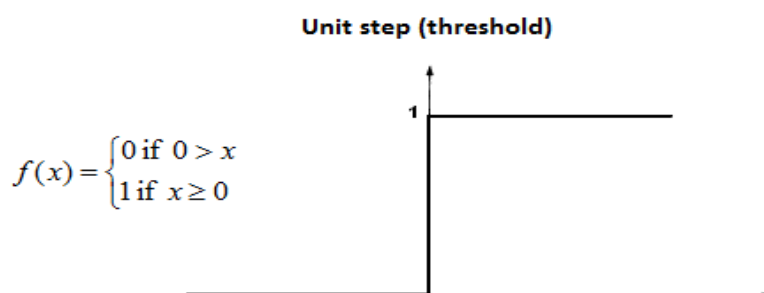


Summation function $y = f(X_1W_1 + X_2W_2 + X_3W_3 + b)$, Where, w_1, w_2, \dots Are the edge weights attached to the input variables, b is the bias.

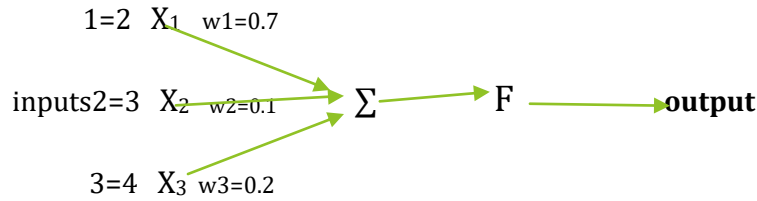
activation function for single-layered neuron $= f^{(1)}(x) = a \left[\omega_0^{(1)} + \sum_{n=1}^N \omega_n^{(1)} x_n \right]$ where, the powers on f and w show that this is a single layered neuron.

For example,

Lets assume that the activation function is a unit step function and is used to map the input between 0 and 1 such that



$$y = f(X_1W_1 + X_2W_2 + X_3W_3 + b),$$



$$y = f(1 * 0.7 + 3 * 0.1 + 4 * 0.2) = f(1.8)$$

$$y = 1 \text{ (because } 2.2 > 0 \text{)}$$

3.3.2 multi neuron model

Now that we have already discussed the single neuron mathematical model, we consider if the inputs were going through multiple layers as explained below. Note that our goal was to subject our single layer unit above to multiple functions, we have to express our mathematical notation in a more explicit notation. This will help us be able to work with arrays in the Numpy library.

The notation will be as follows:

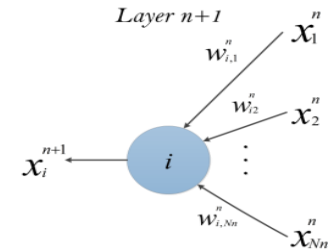
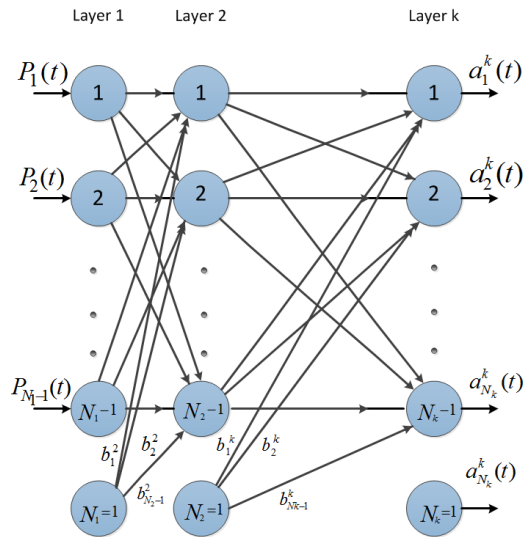
$$P = \begin{bmatrix} P_1 \\ P_2 \\ \vdots \\ P_{N_1} \end{bmatrix}$$

we will use this to work out the following matrix.

$$\begin{pmatrix} W_{11} & W_{12} & \dots & W_{1, N_1} \\ W_{21} & W_{22} & \dots & W_{2, N_1} \\ \vdots & \vdots & \ddots & \vdots \\ W_{m-1,1} & W_{m-1,2} & \dots & W_{m-1, N_1} \\ W_{m,1} & W_{m,2} & \dots & W_{m, N_1} \end{pmatrix}$$

Where m- the number of hidden neurons in the hidden layer.

N_1 - is the number of inputs



$$x_i^{n+1}(t) = \varphi\left(\sum_{j=1}^{N_n} w_{i,j}^n x_j^n(t)\right)$$

Node (n+1, i) representation

Non-Linear Sigmoid Activation function

$$\varphi(w) = \frac{1}{1 + e^{-w}}$$

This is how we will use the formula the formular.

$$\varphi(\omega) = \frac{1}{1 + e^{-\omega}}$$

$$\frac{d}{dw}(\varphi(w)) = \frac{d}{dx} \frac{1}{1 + e^{-\omega}}$$

Using quotient rule to find the derivative

$$\frac{vU' - uv'}{v^2}$$

$$\frac{(1 + e^{-\omega})(0) - (1)(-e^{-\omega})}{(1 + e^{-\omega})^2}$$

$$\frac{e^{-\omega}}{(1 + e^{-\omega})^2}$$

$$\frac{1 - 1 + e^{-\omega}}{(1 + e^{-\omega})^2} = \frac{[1 + e^{-\omega}]}{(1 + e^{-\omega})^2} - \frac{1}{(1 + e^{-\omega})^2}$$

$$= \frac{1}{(1 + e^{-\omega})} - \frac{1}{(1 + e^{-\omega})^2}$$

$$= \frac{1}{1 + e^{-\omega}} \left(1 - \frac{1}{1 + e^{-\omega}}\right)$$

$$= \varphi(\omega)(1 - \varphi(\omega))$$

this is stored as forward and backward propagation.

3.4 Scaling the data

We scale the data so as we can work with a sample of the data. Note that our data has over 6000 rows. We scale our data using the follows the following formula.

$$x' = \frac{x - x_{min}}{x_{max} - x_{min}}$$

Information taken
from old data

Where, x' = the new rescaled value



TECHNICAL UNIVERSITY OF KENYA

CHAPTER 4

MODEL ANALYSIS AND DISCUSSION

4.1 DATA COLLECTION

We collected a data set from www.Kaggle.com. Kaggle is a website that provides data analysts and researchers with the necessary data and resources needed for research. We have provided the link to the dataset above. This data was the most suitable dataset to suit our purpose of creating a model that predicts the amount of money a customer is willing to pay for a car given their age, annual income, creditworthiness as well as gender and net worth.

4.2 IMPORTING THE LIBRARIES

We started by importing the libraries into the Jupiter notebook. These libraries helped in setting up the environment in which we carried out data visualization as well as data analysis. We used the following block of codes.

```
import pandas as pd---Pandas libraries are useful when working with tabular data

import NumPy as np---NumPy libraries are useful when working with arrays

import matplotlib.pyplot as plt----these libraries are useful in plotting data presentations

import seaborn as sns---these libraries are useful in data visualization

%matplotlib inline

sns.set()
```

4.3 DATA VIEW

In this step, we loaded our data set into our work frame. This was to ensure that all the values were included in our dataset. the dataset was in tabular form including 500 rows and 9 columns. We, therefore, scaled this down to 10 head rows and 5 tail rows. We also queried the information about these rows to ensure that there were no null values in our dataset. the following is a representation of the output.

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 500 entries, 0 to 499
Data columns (total 9 columns):
 #   Column                Non-Null Count  Dtype
---  -
 0   Customer Name         500 non-null   object
```

1	Customer e-mail	500	non-null	object
2	Country	500	non-null	object
3	Gender	500	non-null	int64
4	Age	500	non-null	float64
5	Annual Salary	500	non-null	float64
6	Credit Card Debt	500	non-null	float64
7	Net Worth	500	non-null	float64
8	Car Purchase Amount	500	non-null	float64

Fig 2.0- figure showing the 9 columns and classification of content data types

The output also classified the data into different data types depending on the contents in the column.

4.4 DATA PREPROCESSING

In this step, we checked the descriptive statistics of the data. we ran a code to obtain the mean, standard deviation, the interquartile as well as the min and max values from the data.

	Gender	Age	Annual Salary	Credit Card Debt	Net Worth	Car Purchase Amount
count	500.000000	500.000000	500.000000	500.000000	500.000000	500.000000
mean	0.506000	46.241674	62127.239608	9607.645049	431475.713625	44209.799218
std	0.500465	7.978862	11703.378228	3489.187973	173536.756340	10773.178744
min	0.000000	20.000000	20000.000000	100.000000	20000.000000	9000.000000
25%	0.000000	40.949969	54391.977195	7397.515792	299824.195900	37629.896040
50%	1.000000	46.049901	62915.497035	9655.035568	426750.120650	43997.783390
75%	1.000000	51.612263	70117.862005	11798.867487	557324.478725	51254.709517
max	1.000000	70.000000	100000.000000	20000.000000	1000000.000000	80000.000000

TECHNICAL UNIVERSITY OF KENYA

Fig 3.0- figure showing the descriptive statistics of our data and its quartiles

4.5 EXPLORATORY DATA ANALYSIS

In this step, we conducted an informed exploratory data analysis. The objective of this step was to obtain the correlation between the different variables in our dataset. this correlation was obtained from the regression equation $Y = aX + b$ where **a and b are correlation coefficients**. We constructed a correlation heatmap to explain the correlation between all the variables in the dataset

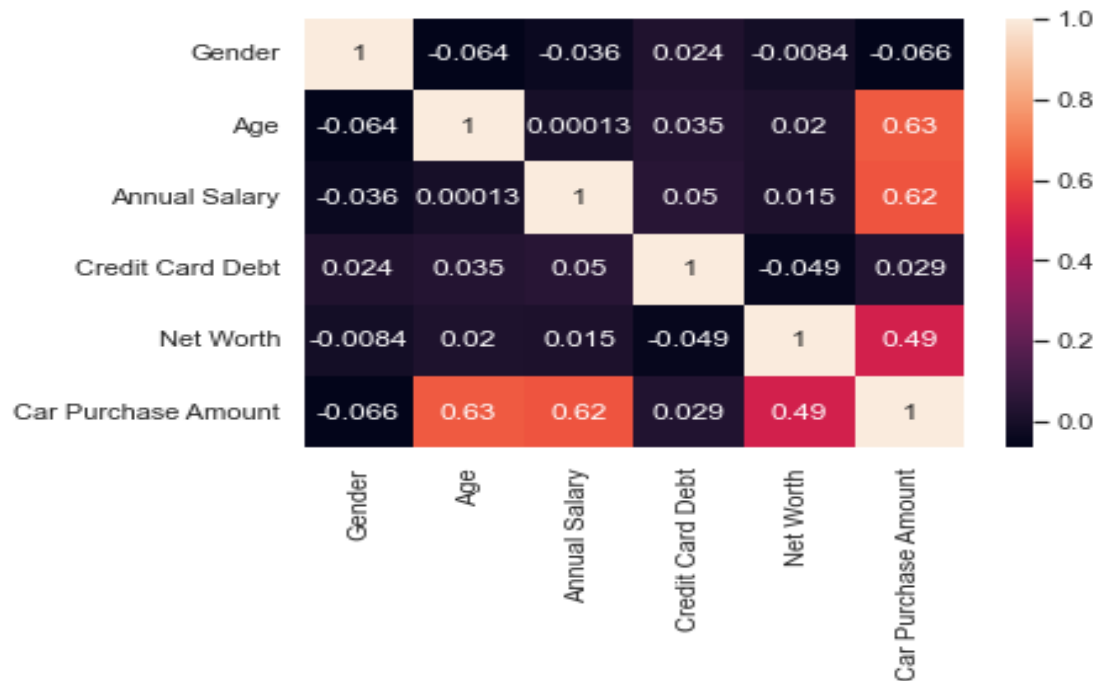


Fig 4..0- figure showing the correlation heatmap of different variables

From the heat map, we noted that this matrix is symmetric. The highest correlation was between age and car purchase amount recording a 0.63 correlation followed by the annual salary of the customer with a correlation of 0.62. the lowest correlation recorded was between gender and net worth meaning that the gender of a customer does not determine their net worth. The diagonal indicates that there is a perfect positive correlation between every variable in the table with itself.

r=1

4.6 LINEAR REGRESSION

Linear regression works by predicting one variable Y based on another variable X. X is called the independent variable while Y is the dependent variable. We were guided by the equation $Y=aX+b$. The goal is to obtain the relationship (model) between the car purchasing amount and annual salary which has the largest correlation.

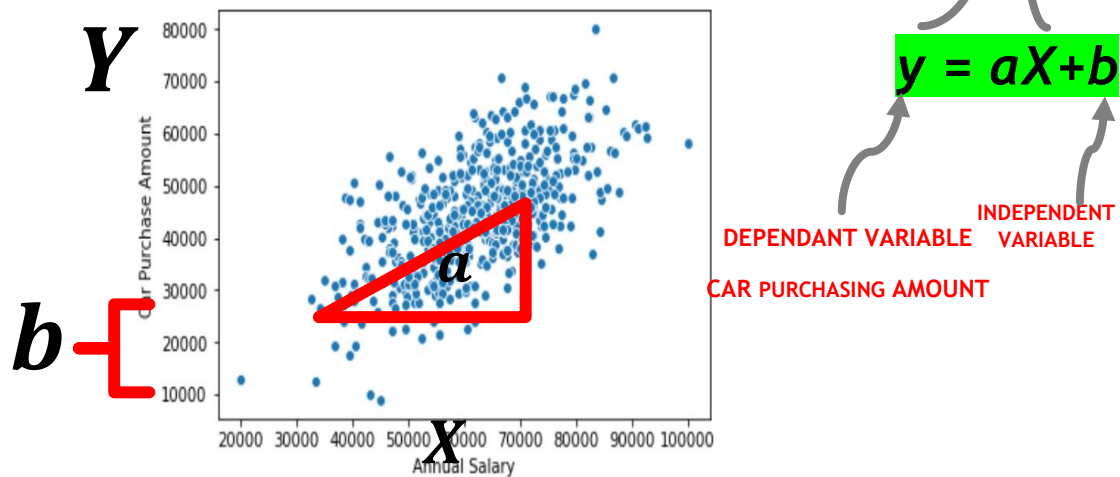


Fig 5.0 – figure showing the regression scatter plot between annual salary and car purchase amount

4.7 MODEL TRAINING-

To train the model, we first defined our independent variable x and dependent variable. We drop the following rows to represent the x-variable: **Gender, Age, Annual Salary, Credit Card Debt, and Net Worth** which can be represented in the following table.

	Gender	Age	Annual Salary	Credit Card Debt	Net Worth
0	0	41.851720	62812.09301	11609.380910	238961.2505
1	0	40.870623	66646.89292	9572.957136	530973.9078
2	1	43.152897	53798.55112	11160.355060	638467.1773
3	1	58.271369	79370.03798	14426.164850	548599.0524
4	1	57.313749	59729.15130	5358.712177	560304.0671
...
495	0	41.462515	71942.40291	6995.902524	541670.1016
496	1	37.642000	56039.49793	12301.456790	360419.0988
497	1	53.943497	68888.77805	10611.606860	764531.3203
498	1	59.160509	49811.99062	14013.034510	337826.6382
499	1	46.731152	61370.67766	9391.341628	462946.4924

500 rows × 5 columns

Fig 6..0- figure showing the table of rows used as x-variable

our y-variable is the car purchase amount. Notice that the gender row is represented in 0s and 1s to represent female and male respectively.

4.7.1 DATA NORMALIZATION

Data normalization is the process of transforming our data into unit spere or rather representing our data in the Euclidean space. Meaning in the range [0,1]

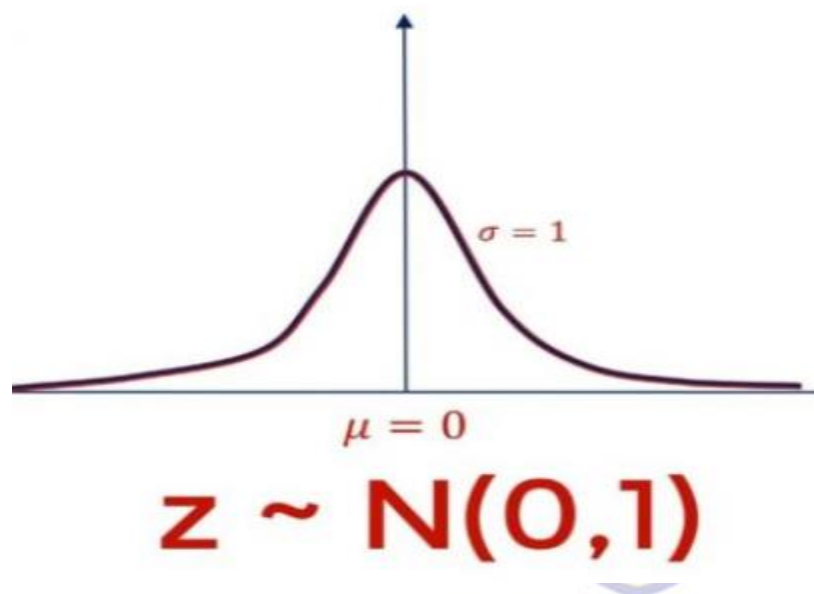


Fig 7.0- figure showing the graph of normal distribution.

This normalized data is now uniformly distributed therefore training the model is easy now. However, we do not normalize all the data, we normalize only the sample that we are about to use.

4.7.2 TRAIN, TEST AND SPLIT TECHNIQUE

After normalizing our dataset, we then continue to prepare our data for predictions. We split our data into two portions, that is, the train sample and the remainder. We train about 80% of our dataset. This ensures that we maximize the inputs and minimize the errors during training. However, from this training set, we test a set of 25% of the training set. This helps to reduce the time taken during testing of the trained data. The trained data is as random as possible. This means that every time we run the trained sample, it gives us a different value.

The model is trained sequentially in phases named epochs. An epoch indicates the number of passes of the entire training the dataset that our machine learning algorithm has completed. This technique is found in the tensor flow library. We define our epoch as 100 meaning that we want the model training to be done in 100 instances. This ensures accuracy because data is worked on in different small portions and then combined into one unit. We also use the verbose algorithm to ensure that the computer displays every epoch instance and all the information in the background including the value loss.

4.7.3 MODEL EVALUATION

After training and testing our data, we determined the model value loss. This helped us check if our model was appropriate. The model had a higher loss at the beginning of our training but gradually reduced the loss to values approaching 0. We came up with this graph from the following equation of root mean squared error

$$RMSE = \sqrt{\sum_{i=1}^n (predicted - Actual)^2 / N}$$



Fig 8.0 figure showing the graph of model loss during training

4.74 MODEL RESULTS

After training our model and testing our model, we finally have an output value. This value is variant every time we train a new set of data. This makes our model suitable to predict the car purchase amount because different data is trained every time a new customer arrives.



CHAPTER 5

5.1 CONCLUSION

In conclusion, we created a machine learning model that was able to predict the amount of money a customer might be willing to pay for a car provided their financial records as data. Using this data, we can be able to identify patterns with other past car sales. This model used linear regression techniques to be able to identify and match patterns in the data. We were able to exclusively evaluate all the assumptions such as all customers willing to buy a car from our dealer or manufacturer company will also be willing to provide their financial information so that a car worth their money can be suggested. This model can also be fed into a company's website as an extranet so that customers can be able to interact with it. The model does not need to train the data at all times, the model should be able to detect patterns matching past sales and be able to suggest car models to different customers on the website so that they can be able to choose from a collection of cars ranging from the price indicated. In addition, our model was able to close the gap in the car sales market of embracing technology and appreciating the inputs made possible by the development of machine learning over time to be able to fit in the industry. Now customers can be able to buy or order cars from the comfort of their homes once this model is deployed in a car dealership website. This invention will help customers easy their hunt for cars, they can even start saving up to increase their worth so as they can be able to buy their dream cars. This way we will create a society that is mindful of their financial life.

TECHNICAL UNIVERSITY OF KENYA

5.2 Appendices

```
# importing our libraries

import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

%matplotlib inline
sns.set()

#loading our data set into our workspace. note that the data is in .csv format.
car_df=pd.read_csv("Car_Purchasing_Data.csv",encoding = "ISO-8859-1")

car_df

# extract the information in the first 10 rows in our dataset,and the last 10 rows.
# make sure that we dont have any null values
car_df.head(10)
car_df.tail(5)
car_df.info()

# describing our data so as to determine the mean, standard deviation, min&max values as
well as the quatiles
car_df.describe()

# constructing a pairplot to determine the relationship between the different variables in our
dataset
sns.pairplot(car_df)

sns.distplot(car_df['Car Purchase Amount'],
hist_kws=dict(edgecolor='red',linewidth=1,color='green'))

# plotting our scatter plot and the line of best fit
```

```

plt.figure(figsize=(10,8))
ax = plt.axes()
plt.scatter(x=car_df['Annual Salary'],y=car_df['Car Purchase Amount'],color='red')
ax.plot([20000, 40000, 60000, 80000,90000], [0, 20000, 40000, 60000,70000])
ax.set_xlabel('Annual Salary')
ax.set_ylabel('Car Purchase Amount')
plt.show()#correlation
car_df.corr()

# constructing a heat map to explain our correlation.
sns.heatmap(car_df.corr(), annot=True)

plt.rcParams["figure.figsize"]=(12,8)
car_df.plot.line()
plt.show()

# dropping the various columns that we need in our as our X values
X = car_df.drop(["Customer Name", "Customer e-mail","Country","Car Purchase Amount"],axis
= 1)

# dropping the item column that we need as our Y values
y = car_df["Car Purchase Amount"]

# normalizing our data in arrays
from sklearn.preprocessing import MinMaxScaler
scaler = MinMaxScaler()
X_scaled = scaler.fit_transform(X)
X_scaled

# getting the maximum values in arrays

```

```

scaler.data_max_
# getting minimum values
scaler.data_min_
# reshaping our data in the purchase amount column
y = y.values.reshape(-1,1)
# fitting our data into the model
y_scaled = scaler.fit_transform(y)
y_scaled
# engaging the train test and split technique test size 25%
from sklearn.model_selection import train_test_split
X_train,X_test,y_train,y_test = train_test_split(X_scaled,y_scaled,test_size =0.25)
# checking the shape of our x_train data set
X_train.shape
# checking the shape of our y_train dataset
y_train.shape
# transforming values in the y_scaled dataset to fit n our model
y_scaled = scaler.fit_transform(y)

#checking the shape of our X_test dataset
X_test.shape
# training the model sequentially
#checking for the output shape
# using the activation fnction to produce an output
# engaging tensorflow
#checking model summary
# cross checking model parameters
import tensorflow.keras
from keras.models import Sequential
from keras.layers import Dense

```



TECHNICAL UNIVERSITY OF KENYA

```

model=Sequential()
model.add(Dense(5,input_dim = 5,activation="relu"))
model.add(Dense(5,activation = "relu"))
model.add(Dense(1,activation = "linear"))
model.summary()

# compiling the model to optimize model training
# using epochs to train our model in a sequential manner
# using verbose to check model progreesion and loss
# calculating model loss using mean squared error
model.compile(optimizer= "adam" ,loss = "mean_squared_error")

epochs_hist = model.fit(X_train,y_train,epochs = 100,batch_size=
75,verbose=1,validation_split=0.2)

# model evaluation
epochs_hist.history.keys()

# plotting the model value loss progrss graph during training and validation
plt.plot(epochs_hist.history["loss"])
plt.plot(epochs_hist.history["val_loss"])
plt.title("Model loss Progress During Training")
plt.ylabel("Training and Validation loss")
plt.legend(["Training loss", "Validation loss"])

# comparing the model predictions and crosschecking the time taken to train the model
X_test = np.array([[1,50,50000,10000,600000]])
y_predict= model.predict(X_test)
# printing the predicted value
print("Expected Purchase Amount",y_predict)

```

REFERENCE

1. Introduction to statistical learning and data science by Gareth James
2. Understanding machine learning from theory to algorithms by Shai-ben David and Shai Shalev-Schwarz
3. Towards data science
4. Machine learning for Dummies by John Mueller and Luca Massaron
5. https://en.wikipedia.org/wiki/Sigmoid_function wikipedia
6. Kaggle.com
7. Github

