# Autonomous Cars – Predicting Future Customers

## Satish Kumar Boguda [1], Arsid Shailaja[2]

[1]Software Engineer – Data Scientist, California, USA

[2] IT Manager – Hyderabad, India

-----------------------------------------------------------------------***-----------------------------------------------------------------------

**Abstract -** *Innovation is something that can never end. The automotive industry has been more influenced than ever by technological advances over the past decade. The driver-less vehicle that has long been a high-tech dream is now happening, it's genuine, not science fiction, yet it's a reality. With the regulatory and infrastructural challenges still remain uncertain, self-driving innovation has pulled in various stakeholders ranging from diminutive startups to leading car manufacturers around the world. Pioneers, for example Tesla, Nissan, BMW, Ford, General Motors Toyota and so on have already spent billions of dollars on research and development, digital platforms, the integration of artificial intelligence and machine learning capabilities to process terabytes of vehicle-generated data to convey real-time information to other vehicles that dodge latency in order to operate safely. If the technical and regulatory milestones are reached, this innovative development is expected to reach the consumer market by 2021. In this paper, we will build an intelligent machine learning model based on the historic car acquisition data of previous customers, that will predict if the future customers who will be able to buy autonomous cars or not. Predicting future customer sales is one of the main pillars of growth to generate more income, which will eventually bring big business value to the organization. In order to develop a predicative data model for future clients with a view to buy autonomous vehicles in the near future, we will utilize the frameworks and libraries of Python.*

*Key Words***: Autonomous Cars, Data Science, Machine Learning, Python, Python Libraries, Pandas, Numpy, Seaborn, Matplotlib, TensorFlow, Scikit-Learn, Regression, Logistic Regression, Data Analysis, Data Visualization.**

## 1.INTRODUCTION

This section provides an overview of the below terms

- Autonomous Cars
- Data Science
- Machine Learning
- Python Libraries
- Data Analysis
- Data Visualization
- Regression
- Logistic Regression

## 1.1 Autonomous Cars

Technologies of the 21st century have the greatest impact on our lives with the digital transformation journey going far beyond the Internet of Everything, the biggest challenge organizations face today is the ability to generate real-time actionable insights for solving the most critical business problems that drives the best value for the organization. However, with the introduction of innovative technologies such as Artificial Intelligence, Machine Learning, Edge computing and Analytics, the transportation sector is redefining the way the vehicles are connected, operated and communicated, which is the game-changer in today's automotive industry.

Self-driven or autonomous cars are robotic vehicles that are designed to travel from source to destination completely without the help of a human operator. All together for a vehicle to be completely autonomous, a vehicle must be able to navigate over roads to a target destination with no human intercession. If you see, as part of traditional cars, human drivers watch for pedestrians around the road, speed signs, surrounding cars, traffic instructions, other obstacles, and so on. If anything goes on, you use brakes to avoid collision with them, and if the way is clear, you accelerate and proceed onward and ultimately influence the vehicle to operate safely and accident-free.

Autonomous vehicle replaces the human driver and performs all the above operations on its own through the below features embedded in the car.

**Camera** – To see the objects, signs of traffic, bridges, pedestrians and so on.

**GPS** – Guide routes by specifying the exact location of the vehicle

**Control Units** – All vehicle operations are performed, controlled and monitored

**Sensors** – To see, tune in and process the ongoing information produced from different sensors, for example, Lidar, Radar and Ultrasonic

The automotive industry is rapidly changing itself into a smart ecosystem of mobility. The foundation of this innovation is connectivity, opening up numerous opportunities, but also a series of cyber security risks that have never existed before.

Society of Automotive Engineers (SAE) International, a US-based automotive standardization body, has characterized 6 different levels of automation classification ranging from fully manual to fully automated systems.

**Level 0** – the vehicle is fully monitored by the human driver, and there is no vehicle control from the automated system of the vehicle.

**Level 1** – One or more vehicle controls such as Automatic Braking, Parking Assistance, Cruise Control and so on are automated. Human driver has the ability to regain control at any time.

**Level 2** - Two or more interconnected automated features, e.g. cruise control in combination with lane changing capability, work in parallel. However, when control is taken over by the human driver if necessary, the automated system is deactivated.

**Level 3** - The driver can turn his attention away safely, while automated technology makes it possible to drive without human intervention for a long time. However, during an emergency the car identifies the situation and cautions the human driver to take over providing ample time to take control of it.

**Level 4** - All critical safety functions are performed without driving intervention throughout the journey by automated technology. Infact the driver can watch his favorite Netflix Series  or to go to sleep or relax. However, in some extreme emergency situations such as severe weather, traffic jams and so on, the human driver control is required otherwise if the driver does not respond, the car will abort the journey accordingly.

**Level 5** – There is no need for human intervention at all. The car is fully controlled by automated technology and is capable of performing all operations and can drive to any destination provided it is legal to operate a fully autonomous car.

## 1.2 Data Science

Data science is a combination of various areas such as statistics, programming, algorithms, analytics, processes and so on that enables you to generate strategic information and drives organizational decisions.

Data has turned into the most copious thing that we have today. Be it an industry giant like Walmart or a mid-size shop retailer, every company has been storing the data in some or the other form which results to enormous data explosion. One of the biggest challenge enterprises face today is what can they do with this data? Companies like Amazon, Target, General Motors, Tesla and so on have access to large quantities of consumer statistics containing purchase information, past search history, customer location, age, gender, etc. and apply various data science models on this data to segment their customers, identify different patterns, generate actionable insights and recommend the right product to the right customer at the right time.

Below are some of the benefits businesses can gain by applying various data science process and systems:

- Identify risks ahead of time
- Reduced operational cost
- Faster and accurate decision-making
- Deliver customers what they want in real-time
- Prevent Fraudulent Transactions

## 1.3 Machine Learning

"Machine learning is the subset of Artificial Intelligence that enables the computer to act and make data driven decisions to carry out certain tasks. Machine learning algorithms (programs) are designed in a way that they can learn and improve over time when exposed to new data. In other words, machine learning is a type of artificial intelligence that provides computers with the ability to learn without being explicitly programmed.

One methodology where the machine learning algorithm is trained first is using a set of labelled or unlabelled training data to create the model. In the next step, the machine learning algorithm is presented with new input data which makes the prediction based on the training model and that forecast is evaluated for precision and if the accuracy is worthy, then the machine learning model is deployed. Now if the accuracy is not satisfactory, the Machine Learning algorithm is trained repeatedly with an expanded training data-set until the accuracy is satisfactory.
Machine Learning is used in various industry sectors which includes

- Travel Industry
- Marketing
- Healthcare

- Social Media
- Sales
- Automation
- Credit & Insurance"

## 1.4 Python

Python is a high-level open source programming language which was created in 1989 by a Dutch Programmer named "Guido Rossum" and it is supported by many platforms such as Windows, Linux etc. Python is the most popular language used today, whether it's a government agency or a defence system or NASA or a Silicon Valley giant, that almost every company today uses python.

Using Python libraries, a wide variety of projects can be developed, including:

- Big Data
- Web Development
- Data Analysis
- Data Visualization
- Machine Learning
- Computer Vision
- Game Developments
- Web Scraping
- Scripting Automation
- Browser Automation
- GUI Development

Some of the advantages of Python Programming Language includes

- Simple and Easy to Learn
- Portable and Extensible
- Supports integration with various programming languages
- Free and open Source
- Huge collection of libraries

## 1.5 Python Libraries

Python library is a collection of standard functions and methods that enables you to perform many actions without explicitly having to write your own code. When we install the Python package, the standard libraries come as default and any python programmer has access to its in-built modules and functions.

For example, if you want to find square root of number. You need not explicitly write code to find it, however with the help of built in library's and functions you can find it with a single line of code like using "math.sqrt(x)"

Example:

math.sqrt(15)  → Square root of 15

In this case we make use of the inbuilt library *math* which provides access to the mathematical functions.

Below are some of the most commonly used libraries in python.

### Pandas

- Pandas is a python library consisting of high-level data structures and tools used for data manipulation and data analysis. Pandas ' main purpose is to help us identify data intelligence.
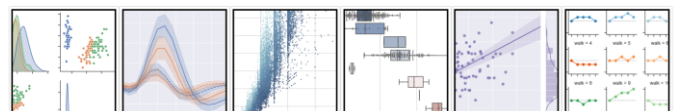


### Numpy

- Numpy stands for Numerical Python, an open source library used for scientific computing such as statistical, mathematical, scientific, engineering, and data science programming. Numpy works very well with multi-dimensional arrays and matrices. Most of the Data Science associated packages are built on top of NumPy.
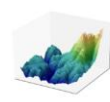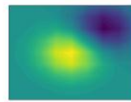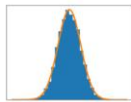
### Seaborn

- Seaborn is a statistical data visualization python library primarily based on matplotlib. It provides a high-level interface for drawing desirable and informative statistical graphics and works very properly with Pandas.



### Matplotlib

- Matplotlib is one of the prevalent Python 2D plotting library that provides various functions to plot excellent and appealing graphs. Utilizing Matplotlib you can create distinctive plots that include histograms, bar charts, error charts, scatterplots and so on with only a couple of fundamental lines of code

- **TensorFlow**

  TensorFlow, otherwise called Google TensorFlow, is a standout amongst the most mainstream open source libraries initially created by the Google Brain group inside Google's Machine Intelligence research organization for machine learning and profound neural networks research. You can design, build and train deep learning models utilizing TensorFlow.

- **Scikit-Learn**

  Scikit-learn is an open source reusable machine learning library for the Python programming language used in data mining and statistical analysis. It provides various machine learning features that include both supervised and unsupervised learning algorithms such as classification, clustering, regression, and so on.

## 1.6 Data Analysis

As you already know, data is collected everywhere around us. Whenever you browse any website on the internet or access a video on your mobile device, data is either collected manually or digitally by scientists. Data doesn't mean information, you need to unlock information and insights from raw data to answer your questions and that is where we use data analysis.

Data Analysis is a process of identifying, inspecting, cleaning, transforming and modelling the data with the aim of discovering meaningful information, making conclusions and aided decision-making in business and scientific areas. With the analysis of data, businesses can identify both risks and opportunities head of time before the event occurs. Some of the key benefits of data analysis are as follows:

- Get the right information at the right time for your business.
- Comprehend the key business process.
- Serve the customer better
- Generate more business value for the organization.

Data Analysis process comprises of the various phases that are iterative in nature –

- Data Collection
- Data Processing
- Data Cleaning
- Data Modelling
- Data Visualization

## 1.7 Data Visualization

Data visualization is an augmenting area in the business world that offers with the illustration of information in the structure of chart, diagram, image and so on and turns massive and small information units into visuals that are less difficult for the human genius to recognize and process. These are created as the visible illustration of information.
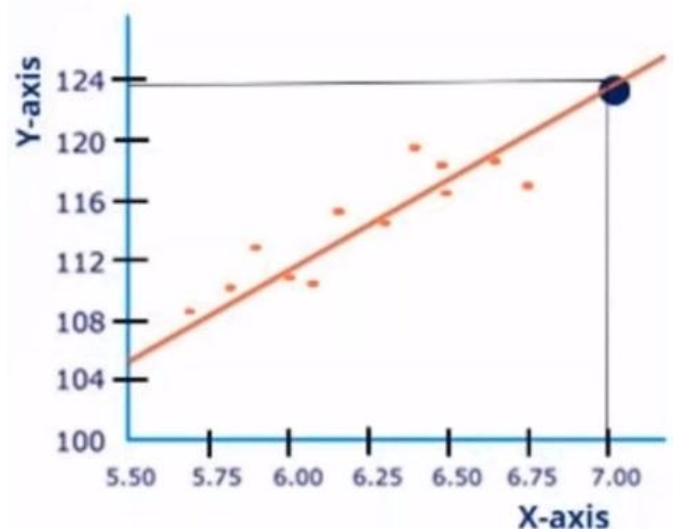
One of the key reasons we use data visualization is that, Human Beings process visual contribution to the exclusion of everything else and faster than any other strategy. By nature, people are visual, so it helps to see data in a visual format. Finding patterns, digesting, and making decisions with visual data is easier. At the point when numerous individuals look at data, they only perceive a sea of meaningless numbers. Data visualization is a pattern that seeks to counteract the natural reaction many people have to raw numbers.

Some of the key benefits of Data Visualization.

- Exploratory Data Analysis
- Communicate Data Clearly
- Less space and more information.
- Visually appealing reports
- Share unbiased representation of data
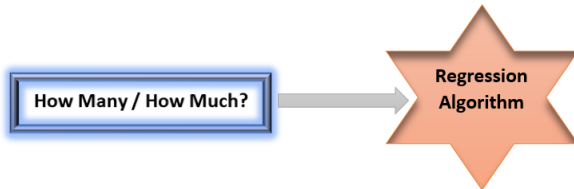- Can be accessed quickly by wider audience.

## 1.8 Regression

Regression, one of the most important technique used for modelling and analyzing the data is a supervised learning algorithm which is a predictive modelling technique that helps in identifying the relationship between a dependent variable usually referred as "target" and independent variables knows as "predictor".



$$y = m x + c$$

Regression Algorithms are used to calculate the numeric value for example whenever the problem demands that you have consisting of values of a target variable and of one or more explanatory variables to get a mathematical value, then in such case, you go with Regression Algorithm.
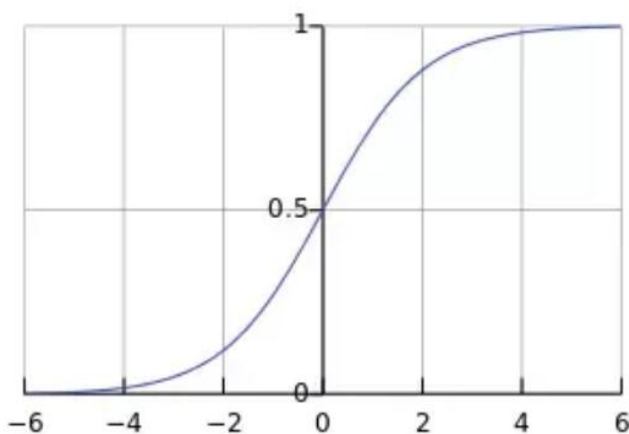


Using this regression algorithm, we fit a curve / line to the data points in such a way that the variations between the distances of data points from the curve or line are minimized. In this example we will predict the number of customers who will be interested to buy the autonomous cars based on the features from the purchasing history of previous customers using customer data set to predict the future sales of the company using the current and historical information.

Regression algorithms are classified into 3 types

- Linear Regression
- Logistic Regression
- Polynomial Regression

## 1.9 Logistic Regression

Logistic regression falls underneath the category of supervised Machine learning algorithm, which is basically used when dividing our data into two or more groups and output should be discrete or categorical. The relationship between the categorical dependent variable and one or more independent variables is measured by evaluating probabilities using logistics and sigmoid functions.



$y = logistic (c + x1*w1 + x2*w2 +x3*w3 +........+xn*wn)$

$y = 1 / 1 + e [- (c + x1*w1 + x2*w2 +x3*w3 +........+xn*wn)]$

The end result of the Logistic Regression is a sigmoid curve, likewise recognized as the S-curve where the value of the independent variable on Horizontal - X axis would determine the dependent variable on the Vertical - Y axis.

The sigmoid function is given by the following equation.

$$S(x) = \frac{1}{1+e^{-x}}$$

The Sigmoid function converts any value from -∞ infinity to ∞ infinity to your discrete values which is required in the logistic regression.

Wikipedia - Logistic Regression

The continuous nature of the input variable can provide more granularity in logistic regression and therefore provide more adequate responses. For example: You are leading an employee engagement workshop as a HR Manager for one of the fortune 500 company and you want to predict which employees will depart the organization in the future. You have access to the historical information of previous employees who left the company and you discover some of the below factors that precipitated the previous employees to leave the organization.

- Employee John was searching for a pay hike.
- Employee Victor felt there was no potential group with the organization.
- Employee Harry was not content with the organizational changes.
- Employee Maria looking for an opportunity on site

This way all employees will have one variable in the data that says left or still working. Therefore, variable that has values left and still working will be our dependent variable and all different variables will be our unbiased variable due to the fact the worker will leave the organization because they are now not blissful with the job.

With actionable insights generated from the logistical regression output, companies can streamline their corporate strategies to accomplish their business goals like minimizing expenditures or losses, augmenting rate of profitability, ROI and so on.

## 2. Case Study

The rapid evolution of the shared versatility advertisement, especially in the automotive industry, for example, Uber, Lyft has shown not only the upsides of advantage lights action plans, but also their ability to transmit quickly across the globe. Autonomous vehicles are among the top innovations and no one could have imagined that such vehicles would ever hit the road about a decade ago. However, with the innovation of technology advances such as artificial intelligence and machine learning models, business concepts and models are revolutionized, ultimately leading to a faster, more convenient and more efficient way of processing business transactions that helps to understand customers in real time by offering the right service/product at the right time and at the right place that ultimately leads to high business value. Historical-based machine learning models can provide us with hidden insights, patterns or tendency to predict the future before it actually takes place. In this approach, we will predict the future sales of autonomous cars that are still in the very early stages of consumer reality by using the historical sales data of the customers who purchased the traditional cars through evaluating different attributes and identifying different features from the existing dataset that gives us the ability to predict the chances of a customer buying the driverless car or not.

**Note**: No evidence concerning regulatory approvals in the federal regulations governing autonomous vehicles has been described or disclosed in this document and there is no confidential information used about the research conducted by any automotive manufacturer worldwide.

## 3. Proposed Methodology

We will develop a predictive logistic regression model with the Python environment in this paper. Accessing the open source libraries such as Pandas, Numpy, Seaborn, Matplotlib and so on allows us to use the existing methods and functions which provide us with nice integrated features without having to write our own code explicitly. In order to build the accuracy of the model, we will access the Jupyter Notebook, a web-based, open source interactive computational environment that helps us to create an accurate predictive model.

Now in order to predict the category of customers who may be interested in purchasing autonomous cars, we need to access the historical sales information dataset about existing customers. Also, we must identify the factors that influence customers when buying autonomous cars with the help of the logistic regression.

Note: No customer sales data related to any auto manufacture has been used in the paper. The data set below represents sample records of sales information.

| Sno | Customer Number | Customer Name | Location | Age | Gender | Salary | High End Model | Category |
|---|---|---|---|---|---|---|---|---|
| 1 | 100023456 | John | Los Angles | 45 | Male | 120000 | 1 | B |
| 2 | 100023786 | Maria | Irvine | 34 | Female | 100000 | 0 | B |
| 3 | 100023894 | Sneha | Los Angles | 32 | Female | 145000 | 1 | B |
| 4 | 100046578 | Harry | Santa Ana | 27 | Male | 132000 | 1 | L |
| 5 | 100090876 | Thomas | Irvine | 26 | Male | 125000 | 1 | B |
| 6 | 100034501 | Jerry | Sandiego | 36 | Male | 160000 | 1 | B |
| 7 | 100009876 | Kate | Bakersfield | 41 | Female | 90000 | 0 | B |
| 8 | 100034976 | William | Beverly Hills | 20 | Male | 225000 | 1 | B |
| 9 | 100020986 | Hudson | Sandiego | 25 | Male | 155000 | 1 | L |
| 10 | 100012785 | Mike | Carmel Valley | 33 | Male | 135000 | 0 | L |
| 11 | 100028734 | Angelo | Escondido | 28 | Male | 124000 | 0 | B |
| 12 | 100000678 | Daniel | Bakersfield | 30 | Male | 110000 | 0 | L |
| 13 | 100023567 | Nick | Irvine | 19 | Male | 150000 | 1 | L |
| 14 | 100028765 | Pavan | Los Angles | 39 | Male | 180000 | 0 | B |
| 15 | 100098790 | Charles | Long Beach | 24 | Male | 130000 | 1 | B |

### Step 1 - Importing the Dataset into Python Environment

In this step we will access the Jupyter Notebook and import the dataset into the python environment by making using of the standard python libraries Pandas, Numpy, Seaborn, Matplotlib

**Loading the Historical Sales Data of Previous Customers**

```
import seaborn as sb
import math
import pandas as pds
import numpy as npy
import matplotlib.pyplot as mplp
%matpalotlib inline
```

```
autocar_data= pd.read_csv("Autonomous_Cars.csv")
```

### Step 2 - Analyzing the Data

Use the Jupyter Notebook to break down the loaded data set and analyze the distinctive relationships by asking some of the following questions.

1) Category of customers with high-end car, buy or lease.
2) Which location has the major high-end car purchases?
3) What is the age of customers who are under high salary?
4) Male customers vs female customers who leased the car.

By accessing data from different columns, we will create different plots such as a distribution graph, a correlation graph, and so on to determine how one variable affects the other variable.

Presently we will draw up a check track for the quantity of clients who leased the vehicle and other people who purchased the vehicle by utilizing the Seaborn library.

```
sb.countplot(x="leased", data=autocar_data)
```

The number of males vs. females who rented the car will be compared using the below code.

```
sb.countplot(x="leased", hue="Gender", data=autocar_data)
```

Next, the age of individual customers are analyzed using the code below.

```
autocar_data["Age"].plot.hist()
```

### Step 3 - Data Cleaning

During the data science model building process, collecting complete data is one of the biggest challenges. You need to handle big data sets in most scenarios that generally contain null values or incomplete documents. It is the primary responsibility to clean the data and remove all unnecessary data values that are part of the dataset, ultimately resulting in greater accuracy.

First, we will check if there are null values on the dataset. The code below tells us if the values in the data set are null then a Boolean result is written down to us, i.e. it essentially checks your missing data and it is either true or false

```
autocar_data.isnull()
```

Now the null values can be superseded by correct alternative values or the column can be removed. Using the below code we can drop the columns

```
autocar_data.drop("Sno", axis=1, inplace=True)
```

Using the below code, you can check if the null values are removed from the data set or not.

```
autocar_data.isnull().sum()
```

Note: Using the above approach, the different columns can be accessed accordingly based on your dataset.

### Step 4 - Build the model using training data set

In this step, we will create the model with the train data set. The first thing we are doing here is to divide the data set into 2 parts, data set for training and testing. The dependent variable and independent variable are defined.

We now will use the columns "customer age" and "salary" of two independent variable data to predict whether the independent "high-end model" variable is purchased or not.

Independent Variables

```
X = autocar_data.iloc[:,[3,5]].values
```

Dependent Variable

```
Y = autocar_data.iloc[:,6].values
```

Using the below code, we will split the dataset into training and testing

```
from sklearn.cross_validation import train_test_split

X_train, X_test, Y_train, Y_test = train_test_split(x, y test_size=0.30, random_state=0)
```

Now we will apply the logistic regression

```
classifier=LogisticRegression(random_state=0)
classifier.fit(X_train, Y_train)
```

### Step 5 - Test the Predictive Model using the test data set

In this step we use the test dataset and predict the model with the code below.

```
y_pred=classifier.predict(X_test)
```

### Step 6 - Evaluate the accuracy of the Machine learning Model

In this step, we will evaluate the predictive model's accuracy using the code below. We will utilize the standard "accuracy_score" function to calculate the accuracy of the model.

```
from sklearn.metrics import accuracy_score
```

```
accuracy_score(y_test,y_pred)
```

## 4. Future Scope

The majority of accidents that are happening today are due to human error in recognition, decision or performance which can be avoided by driving your car on a computer that will be statistically much safer than driving a car yourself. Also driving is probably one of the challenging and most time-consuming tasks we do on a daily basis, and this can be avoided by self-driving cars that will bring us hours back to

our day, preventing us from keeping our eyes on the road. Autonomous cars are the future of transportation where companies like Google, Tesla, GM, Toyota, Nissan etc. are already testing their vehicles on the road and most industry leaders predict that self-driving vehicles will become commonplace on our roads in 5 to 10 years. We are in the early stages of this driverless revolution that will disrupt in many ways and transform human lives that no one has ever thought about.

## 5. Conclusion

Data which is the center of gravity for many organizations today, is redefining the businesses to deliver the customers the state of art solutions across different industry sectors by generating meaningful information through real-time data processing  that helps executives make strategic decision that gives best value for the organization. In analyzing the performance, responses and behavior of the key business process, data scientists play a key role by extracting meaningful insights from complex and large datasets and

predicting an event before it happens in the future. By using the python programming frameworks and libraries to leverage the logistic regression machine learning model, this paper demonstrates the accuracy of predictive model that helps companies predict future sales of customers who have the potential to buy the most awaited autonomous innovation car.

## REFERENCES

[1] Boguda, Satish Kumar. "Magnetic Resonance Imaging (MRI) – Digital Transformation Journey Utilizing Intelligent Technologies." *irjet.net.* IRJET,  March 2019. Web. 21 April 2019. <MRI Digital Transformation Journey Utilizing Intelligent Technologies.>

[2] Dr. Meher Geeta "Magnetic Resonance Imaging (MRI) – Digital Transformation Journey Utilizing Intelligent Technologies." *irjet.net.* IRJET,  March 2019. Web. 21 April 2019. <MRI Digital Transformation Journey Utilizing Intelligent Technologies.>

[3] "Self-Driving Car." *Wikipedia,* Wikipedia, 2019. Web. 15 April 2019. <Self Driving Car.>

[4] "Jupyter", *Jupyter,* Jupyter. 2019. Web. 10 April 2019. <Jupyter.>

[5] "Regression." *Wikipedia,* Wikipedia, 2019. Web. 10 April 2019. <Regression.>

[6] "Logistic Regression." *Wikipedia,* Wikipedia, 2019. Web. 10 April 2019. <Logistic Regression.>

[7] "Data Analysis." *Wikipedia,* Wikipedia, 2019. Web. 14 April 2019. <Data Analysis.>

[8] "Data Visualization." *Wikipedia,* Wikipedia, 2019. Web. 18 April 2019. <Data Visualization.>

[9] "The Python Standard Library." *Python Software Foundation,* Python Software Foundation, 2019. Web. 12 April 2019 <Python Library.>

[10] "Python Data Analysis Library." *Pandas,* Pandas, 2019. Web. 25 March 2019. <Pandas.>

[11] "NumPy." *NumPy Developers*, NumPy Developers, 2019. Web. Mar 29 2019 <NumPy.>

[12] Michael, Waskom, "Seaborn: statistical data visualization." *Michael, Waskom,* Michael, Waskom, 2019. Web. April 05 2019. <Seaborn.>

[13] John Hunter, Darren Dale, Eric Firing , Michael Droettboom, Matplotlib Development Team, "matplotlib." *Matplotlib,* matplotlib, 2019.Web April 07 2019. <matplotlib.>

[14] "TensorFlow." *Wikipedia,* Wikipedia, 2019. Web. 16 April 2019. <TensorFlow.>

[15] "Scikit-learn." *Scikit learn.* Scikit learn. 2019. Web. 20 April 2019. <scikit-learn.>

[16] The Art of Data Science by Elizabeth Matsui, Roger D. Peng.

[17] Data Science and Analytics with Python by Jesus Rogel Salazar.

[18] Human + Machine: Reimagine Work in the Age of AI by James Wilson and Paul Daugherty.

[19] Artificial Intelligence and Neural Networks by F.Acar Savaci.

[20] Python Machine Learning by Sebasthan Raschka.

[21] Python Essential Reference by David. M. Beazley.

## BIOGRAPHIES

Satish Kumar Boguda is a skilled IT software engineer with over 12 years of experience by serving various industry customers in Oil & Gas, Utilities, HealthCare Production, Manufacturing, Supply Chain, Sales, Finance, Transport and Energy sectors.

Shailaja Arsid is currently working as Information Technology Manager at Hyderabad.