



---

# **Prediction of Dublin Bike Sharing System**

---

*A thesis submitted in fulfilment of the requirements for the degree  
of  
**Master of Science (Data Analytics)***

*by*

**Poojitha Vishukumar**  
**(Student ID: 18231013)**

*in the  
Discipline of Information Technology*

Supervised by

**Dr. Edward Curry**

*National University of Ireland, Galway*

*August 2019*

## ***Acknowledgement***

I would like to thank Dr. Enda Howley for allocating me this project and giving me the opportunity to work on a growing and applied area of research. I wish to acknowledge the support of my enthusiastic supervisor Dr. Edward Curry whose dedicated support ,expertise and timely feedback helped me understand the direction of the project and without whom this thesis would not have been possible. I also would like to express my deepest gratitude to God almighty. Lastly, I would like to thank my family and friends, who have been there for me over the course, providing advice and encouragement along the way. I hereby, certify that this material, which I now submit for assessment on the programme of study leading to the award of Masters in Science in Computer Science (Data Analytics) is entirely my own work, that I have exercised care to ensure the work is original, and does not to the best of my knowledge breach any law of copyright, and has not been taken from the work of others save and to the extent that such work has been cited and acknowledged within the text of my work.

**Signed:**

**Dated: August 28, 2019**

# *Abstract*

Bike sharing systems are rapidly being introduced in European cities for daily mobility leading to sustainable transportation. These systems add to city's public transport system providing fast and easy services to the public. The bike sharing is increasing dramatically due to its unique feature of being convenient, low-cost and environmentally friendly, giving rise to problems to both users and operators. Amongst the European countries, Dublin city is considered in this thesis for the study and research. In this city, there are two types of bikes sharing systems, docked and dockless bike sharing systems. Dockless are newly being introduced but docked bikes are increasing tremendously since few years. In docked based bike sharing system, the bikes are rented by user by signing up to the services, checking out and returning the bike at the dock. When the user signs up and checks out, the basic information such as user's information, station check-in, check-out, bike usage duration and many more is stored.

Nowadays, the primary issue of bike sharing is the uneven distribution of bicycles caused by the ever-changing usage and supply. The uneven distribution of bikes can be due to station area, timing, population, weather, events and many more. These external factors affect the availability of the bikes and gives rise to necessary efficient model for prediction of bike's availability. Efficient prediction model will help in future prediction of availability of bikes with respect to the factors helping the users to use the bikes also, guiding the operators for redistribution of bikes accordingly.

This thesis aims at art of bike sharing systems in Dublin city proposing spatio-temporal bike mobility prediction model based on historical bike sharing data and real data on a per - station with sub-hour granularity. With the time series multivariate prediction model, advancement in re-balancing of bikes could be seen paving the way for rapid deployment and adoption of bike sharing system across the world. To analyse the bike, weather and event data and build the prediction model, deep learning models such as regression model and recurrent neural network model are studied and implemented. In particular, Autoregressive Integrated Moving Average model (ARIMA) and Long Short-Term Memory model (LSTM) have been implemented to predict the availability of the bikes at the given timestamp. The root mean square error is calculated and depicted on trend curve deviation graph for both the models where LSTM proves to work best with defined large series data whereas ARIMA shows best results with comparatively less series data.

# Table of Contents

<b>Acknowledgement</b>	<b>ii</b>
<b>Abstract</b>	<b>iii</b>
<b>Table of contents</b>	<b>iv</b>
<b>List of figures</b>	<b>vi</b>
<b>List of tables</b>	<b>vii</b>
<b>Abbreviations</b>	<b>viii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 General Overview	1
1.2 Problem Statement	3
1.3 Core Research Questions	4
1.4 Existing Approaches	4
1.5 Proposed Approaches	5
1.6 System Architecture	6
1.7 Research Methodology	7
1.8 Contributions	8
1.9 Thesis Outline	9
<b>2 State of the Art</b>	<b>10</b>
2.1 Introduction	10
2.2 Time series model	10
2.3 Review of Regression model related work	11
2.4 Review on Long Short-Term model related work	13
2.5 Auto Regressive Integrated Moving Average Model related work	15
2.6 Conclusion	19
<b>3 Dataset description</b>	<b>21</b>
3.1 Introduction	21
3.2 Dublin Bike dataset	21
3.3 Weather dataset:	22
3.4 Data column's description:	23
<b>4 Method and Implementation</b>	<b>24</b>
4.1 Introduction:	24
4.2 Data Preparation pipeline	25
4.2.1 Data collection	25

1.1.2	Data Pre-processing: .....	26
4.2.2	Data Cleaning .....	26
4.2.3	Data Transformation .....	27
4.3	Modelling and Implementation of Machine Learning Algorithms.....	29
4.3.1	Overview of Long Short Term Memory: .....	29
4.3.2	Mechanism and Implementation of LSTM:.....	30
4.3.3	Overview of ARIMA Model: .....	34
4.3.4	Mechanism and implementation of ARIMA model .....	35
4.3.5	Conclusion.....	39
<b>5</b>	<b>Results and Discussions.....</b>	<b>40</b>
5.1	Plots showing correlation of attributes .....	40
5.2	Results of LSTM and ARIMA model.....	45
5.3	Model comparisons .....	47
5.4	Conclusion.....	49
<b>6</b>	<b>Conclusions.....</b>	<b>50</b>
6.1	Limitations.....	51
6.2	Future work.....	52
<b>7</b>	<b>References .....</b>	<b>53</b>

## List of Figures:

1.1 Dublin Docking Terminal	2
1.2 Bike sharing working	2
1.3 Working of Rental and Return of bikes	2
1.4 System Architecture	6
1.5 Methodology Flowchart	7
2.1 Average number of bikes per day using three weeks moving average(L) and yearly moving average(R)	12
2.2 Monthly volume indices, seasonal index curve, and long-term trend estimate used in regression problem	12
2.3 Daily bike counts, temperature and precipitation	16
2.4 Plot of actual and predicted counts for test data	18
3.1 Snippet of Static Bike dataset (Data.gov.ie. 2019)	21
3.2 Snippet of dynamic bike's responses from JCDecaux API	22
3.3 Snippet of weather responses from weather API saved as csv file	22
4.1 Standard machine learning CRISP-DM Pipeline	24
4.2 Data Preparation Pipeline	25
4.3 LSTM network	29
4.4 Mechanism of LSTM	30
4.5 Function of forget gate layer	31
4.6 Function of input gate	31
4.7 Operations by input gate	32
4.8 Function of output gate	32
4.9 Flowchart of ARIMA model	34
4.10 Plot of series data hour-wise for stationarity check	35
4.11 Plot of series data month-wise for trend and stationarity check	35
4.12 Plot of differenced hourly series for stationarity	36
4.13 Plot of differenced monthly series for stationarity	36
4.14 Rolling mean and standard deviation	37
4.15 ADF statistics	37
4.16 ACF plot showing correlation between bike time series and time lag	38
4.17 ACF plot showing correlation between differenced bike time series and time lag	38
4.18 Combination of p,d,q with AIC values	39
5.1 Interactive map showing availability of Dublin Bikes	41
5.2 Hourly weather and availability of Barrow Station	42
5.3 Weather and availability plot	42
5.4 Monthly bike-weather trend for Georges Quay station	44
5.5 Train loss vs Validation loss	45
5.6 Forecasting in LSTM considering all stations	45

<i>5.7 Predicted and expected values in ARIMA</i>	<i>46</i>
<i>5.8 Forecasting through ARIMA model considering one station</i>	<i>46</i>

## List of Tables:

<i>Table 2.1 Experimental parameters.....</i>	<i>14</i>
<i>Table 2.2 Average RMSEs for each method.....</i>	<i>14</i>
<i>Table 2.3 Estimation Results of Base and ARIMA models.....</i>	<i>17</i>
<i>Table 2.4 Summary Statistics for forecast .....</i>	<i>17</i>
<i>Table 2.5 Summary table of Proposed approaches .....</i>	<i>20</i>
<i>Table 3.1: Data column's description.....</i>	<i>23</i>
<i>Table 4.1 Bike and weather attributes used as features in the model .....</i>	<i>28</i>
<i>Table 5.1 Summary table comparing all the approaches .....</i>	<i>48</i>



## ***Abbreviations:***

<b>CRISP-DM</b>	<b>C</b> ross <b>I</b> ndustry <b>S</b> tandard <b>P</b> rocess for <b>D</b> ata <b>M</b> ining
<b>ARIMA</b>	<b>A</b> uto <b>R</b> egressive <b>I</b> ntegrated <b>M</b> oving <b>A</b> verage
<b>ETL</b>	<b>E</b> xtract, <b>T</b> ransform, <b>L</b> oad
<b>WEKA</b>	<b>W</b> aikato <b>E</b> nvironment for <b>K</b> nowledge <b>A</b> nalysis
<b>LSTM</b>	<b>L</b> ong <b>S</b> hort- <b>T</b> erm <b>M</b> emory
<b>RNN</b>	<b>R</b> ecurrent <b>N</b> eural <b>N</b> etwork
<b>RMSE</b>	<b>R</b> oot <b>M</b> ean <b>S</b> quare <b>E</b> rror
<b>DNN</b>	<b>D</b> eep <b>N</b> eural <b>N</b> etwork
<b>AR</b>	<b>A</b> uto <b>R</b> egression
<b>MA</b>	<b>M</b> oving <b>A</b> verage
<b>ARMA</b>	<b>A</b> uto <b>R</b> egressive <b>M</b> oving <b>A</b> verage
<b>ARIMA</b>	<b>A</b> uto <b>R</b> egressive <b>I</b> ntegrated <b>M</b> oving <b>A</b> verage
<b>ACF</b>	<b>A</b> uto <b>C</b> orrelation <b>F</b> unction
<b>AIC</b>	<b>A</b> kaike <b>I</b> nformation <b>C</b> riteria
<b>BSS</b>	<b>B</b> ike <b>S</b> haring <b>S</b> ystem
<b>DCC</b>	<b>D</b> ublin <b>C</b> ity <b>C</b> ouncil
<b>SBRP</b>	<b>S</b> tatic <b>B</b> ike <b>R</b> epositioning <b>P</b> roblem
<b>DBRP</b>	<b>D</b> ynamic <b>B</b> ike <b>R</b> epositioning <b>P</b> roblem

# 1 Introduction

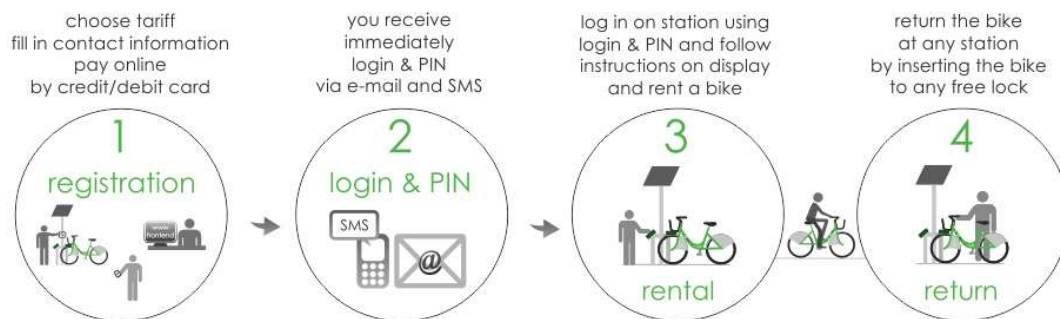
## 1.1 General Overview

Bike sharing systems is a new generation of traditional bike rentals where user is able to easily rent, return and take membership of bikes from a particular place to another automatically without any human resources. According to recent statistics, there are about over 500 bike sharing systems around the world with more than 500 thousand bikes. (Dua, D. and Graff, C.,2019). The BSS program is growing tremendously due to their important role in traffic, environment and health issues. It can be used seamlessly for individual trips consisting of multiple transportation modes which is important component of a modern, sustainable and efficient multimodal transportation network. In general, Distributed bike sharing system is of two types, docked based and dock-less based bike sharing system. In docked based BSS, the bikes are rented and returned from and to the docking stations.

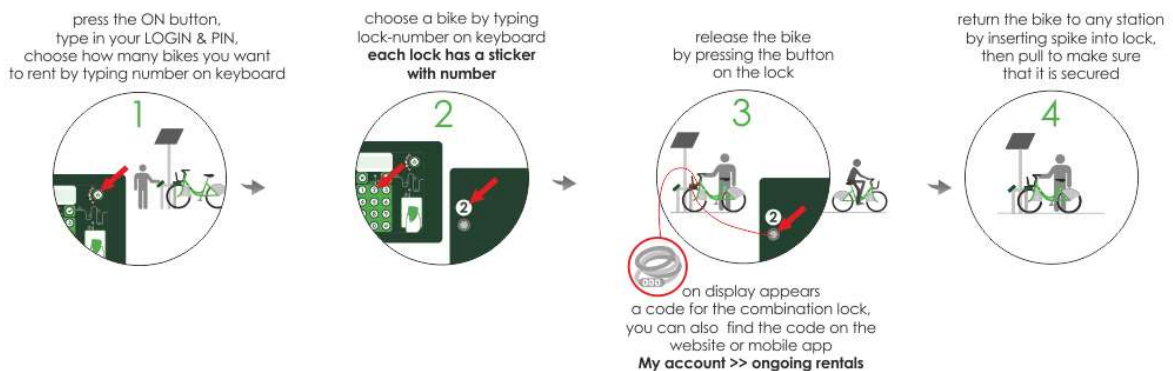
Typically, a rental station includes one terminal stand and several bikes shown in Fig. 1.1. A terminal works on IoT(Internet of Things), where a device is attached to the terminal which communicates with the electronic locker attached to the bikes. (Raviv, T., Tzur, M. and Forma, I. A., 2013) Initially, user needs to register online to the services by entering the basic contact details which proceeds to the payment. After the payment is done, an electronic message is received with PIN number which is used to login at the terminal for rental of bikes as shown in fig 1.2. Internally, when the user rents a bike, signal is sent to the terminal that the bike with respective locker has been vacated. A user can return the bike to destination stand when there is a vacant locker as shown in fig 1.3. All the rental and return transactions are recorded and reported in real time by JCDecaux of Dublin City Council which is the central control facility for Dublin city. Thus, this information of state of system, in terms of number of available bikes and number of vacant lockers available for each station is updated by DCC to the BSS operators and these operators make the information available online for the users.



1.1 Dublin Docking Terminal  
(En.wikipedia.org, 2019)



1.2 Bike sharing working  
(Anon, 2019)



1.3 Working of Rental and Return of bikes  
(Anon, 2019)

Apart from real- world application of bike sharing systems for the accessibility of the users and redistributors, the data generated by these systems can be used for further researches. Antithetical to other transport services such as bus or subway or luas, the total duration of travel, the arrival time and the departure time and position is explicitly recorded in these systems. This feature turns bike sharing system into a virtual sensor network that can be used for sensing mobility in the city (Fanaee-T, Hadi & Gama, João.,2013).

## 1.2 Problem Statement

Bike sharing has significantly been enhancing the urban mobility as a sustainable mode of transportation, but it has few limitations due to the effects of fluctuating spatial and temporal demand. According to previous researches, (Chen et al., 2016, Li et al., 2015, Lin, 2018, Zhou, 2015), it is very common for bike sharing systems with fixed number of terminals that some terminals can be completely empty with no bikes to check out during peak hours of usability while others are full precluding bikes from being returned at those locations. For dock-less bike sharing systems, there is lot of flexibility with the parking system but it even possesses more challenges with respect to ensuring bike availability at some places and also, prevent surplus bikes from blocking sidewalks and parking areas. For both types of BSSs, accurate bike-sharing demand predictions can help operators to generate optimal routes and schedules for rebalancing to improve the efficiency of BSSs.

Docked bike sharing service is widely being used in Dublin bringing much of convenience for the short trips, reducing greenhouse gas emission and congestions but the management of the bikes due to its fluctuating spatial and temporal demand , limiting the capacity usage of bikes and docking stations has increased the concern. Empty or saturated bike stations at certain places at certain time causes problems to the commuters. In order to address this issue, bikes are manually rebalanced by the trucks with an arbitrary timetable. Due to the unbalanced distribution in the number of bikes at docking stations, these resources cannot be fully utilized (Wang and Kim, 2018). This problem of predicting the frequency usage of bike and manually loading of bikes accordingly can be effectively dealt by predicting the availability of the bikes with respect to time at respective stations.

### 1.3 Core Research Questions

The core goal of this research is to answer the questions related to rebalancing problem. For instance, what if everyone gets off from the Dublin bus at the major stop like City Centre and then they all grab the bikes from the nearest station, those bikes will be depleted before long resulting in unequal riding patterns. (Dell'Amico, M., Hadjicostantinou, E., Iori, M. and Novellani, S., 2014).

The major challenging questions that have been focused on and is studied are:

- Is it possible to classify bike sharing stations according to the bike availability?
- What if the station is empty when the user wants to rent a bike?
- What if the station is saturated when the user wants to return the bike?
- How to build a complete ETL framework to analyse the data and report results?
- Can one-commodity pickup-and-delivery capacitated vehicle routing problem be solved?

### 1.4 Existing Approaches

For both types of Bike sharing systems, accurate bike-sharing demand predictions can help operators to generate optimal routes and schedules for rebalancing to improve the efficiency of BSSs. The affordability and availability of bikes influences the variability of users. The ability to predict the number of users can allow the sponsored businesses and governments that oversee these systems to manage them in a more efficient and cost-effective manner.

The problem of repositioning operations which consists of removing and reloading of bikes from overloaded stations to shortage stations using dedicated fleet of trucks where the decisions regarding the routes the truck needs to follow and the number of bikes that needs to be removed or loaded at each stations at each visit of truck is solved by dividing the problem into two categories, Static bike repositioning problem(SBRP) and dynamic bike repositioning problem (DBRP). The repositioning activities carried at the night when the rate of usage of bikes is negligible falls under SBRP and the operations carried during day when the bike systems usage is rapidly changing is considered as DBRP. The Static problem is solved by estimating the demand pattern on past demand data on similar days. Solution for SBRP fails to work efficiently during the unpredictable events such as traffic slowdown, vehicle breakdown, heavy rain, and many more unfortunate events( Raviv, T., Tzur, M. and Forma, I. A., 2013).

The demand faced by the BSSs is also sensitive towards the partial predictable effects such as weather conditions, public events in the city and many more. This kind of situations affecting the usability of the bikes requires dynamic repositioning and which is dealt by statistical time series model which analyses the historical data, the traffic data, arrival and departure rates on hourly, daily and weekly basis. To make the dynamic approach to predict the demand pattern more efficiently, the current traffic trends and weather data is considered along with the historical data.

### 1.5 Proposed Approaches

The use of Internet of Things technologies having embedded sensors to obtain the real-time data has allowed the cities to provide a better service, where the data is used to study, plan and manage the bike sharing docking station and adapt to the need of the commuters. As discussed in the previous sections, the efficient way of rebalancing scheme is by predicting the demand of bikes in different locations covered by different services and selecting the way to relocate bikes in order to satisfy user's demand.

Rather than following the conventional static method for rebalancing, this thesis proposes a dynamic rebalancing schemes which aims at providing user's satisfaction by minimizing probability of system failures. (Chiariotti, F., Pielli, C., Zanella, A. and Zorzi, M., 2018). There are many other factors that alters the usability of the bike sharing patterns, such as time of day, population of the city, location of the station, nearby amenities, weather and climatic changes, week of the day, nearby events and bank holidays and many more. Different machine learning time series regression models help to predict the availability of bikes on basis of multivariate attributes. In this, time series model, a sequence is taken at successfully equally spaced points at each hour of time. There are lot of work done with univariate dataset wherein just the availability of the bikes is considered against the time interval. To make the work closer to real life problems, multivariate data set has been chosen. In this thesis, multivariate time analysis consists of simultaneous multiple time series that deals with dependent data such as time stamp which includes hour of the day, weekend or not, season of the year, public holiday or not, temperature, weather and feels like attributes.

The goal of the thesis to use deep learning models in Keras such as Multivariate Linear Regression and Long short-term memory model which effectively predicts the number of bikes shared which are available at any given time period. It can be effectively be done by using available information of bike shared on hourly basis with weather, day of the week, status of

the day and events. The focus on behavioural pattern of the stations depending on the space and usage can be done to categorize them into highly used to least used stations during the period of time.

## 1.6 System Architecture

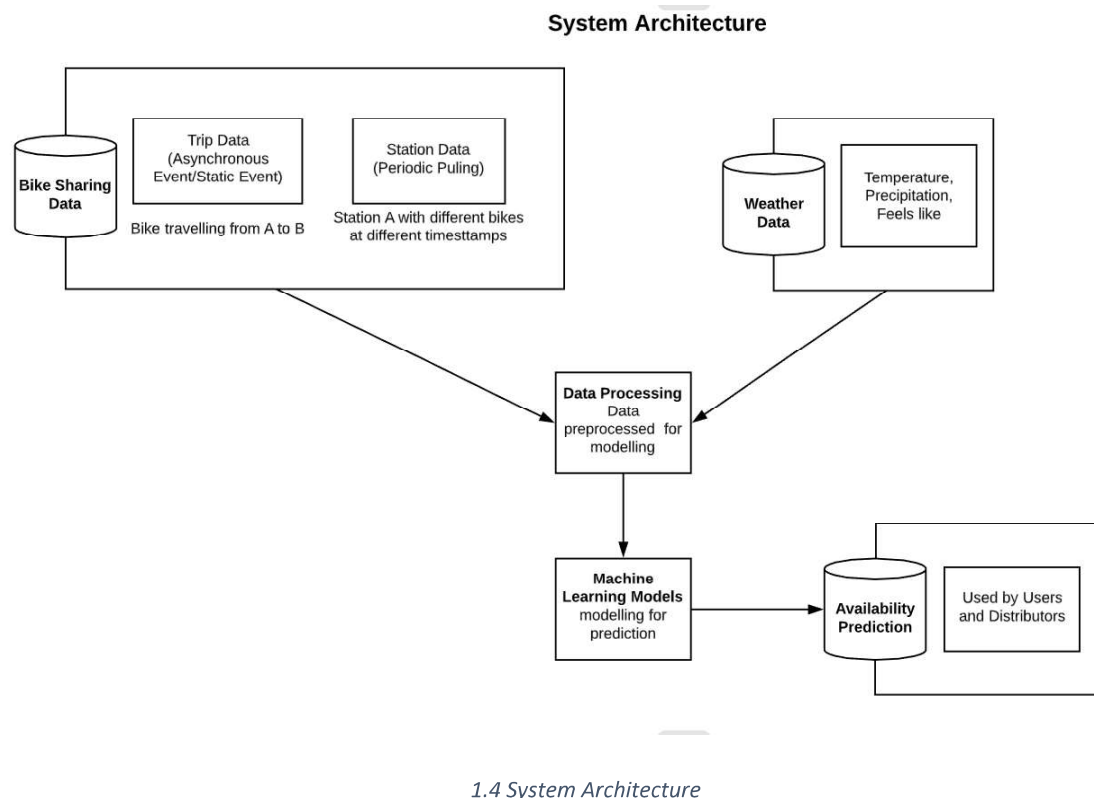


Figure 1.4 describes the system architecture of the data flow from data collection to data transformation and data usage. The bike data and weather data is collected through APIs which include trip data of bike from station A to station B and Station data with its respective number of bikes and the temperature, precipitation and feels like data respectively. The data is further processed and used in the machine learning models for prediction and evaluation of the model. The predicted data is stored and used by the commuters and the distributors in the end.

## 1.7 Research Methodology

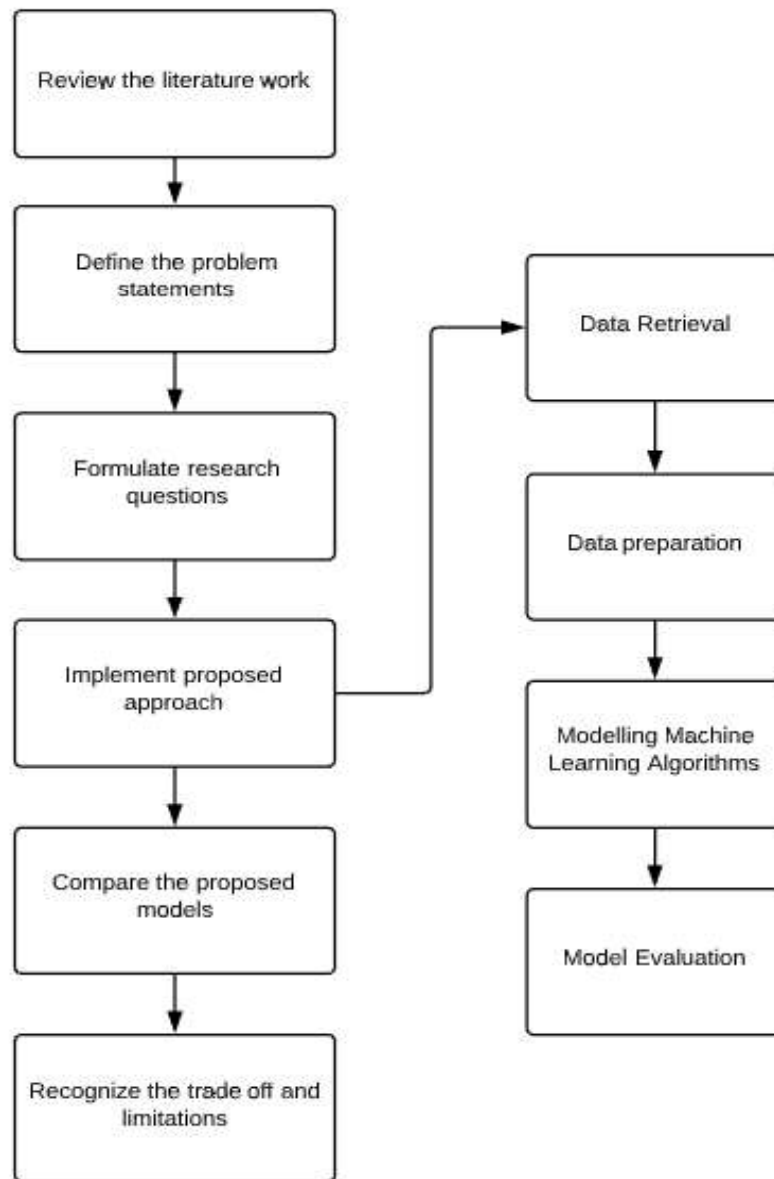


Figure 1.5 Methodology Flowchart



The methodology followed in this thesis is shown in the above fig 1.6 and is explained in the following steps:

- Initially, review of literature work related to Bike sharing systems around the world and their challenges were studied.
- The problems that exists which were gathered from different researches were formulated.
- Research questions were formulated that is planned to study and answer through the model developed in the thesis.
- Two machine learning timeseries models (LSTM and ARIMA) have been developed to solve the defined problem of Bike sharing which follows CRISP-DM steps. The detailed pipeline is explained in Chapter 4.
- After each model's evaluation, its performance is checked through comparing respective root mean square errors.
- Limitations of the proposed models is recognised, and future work is proposed.

## 1.8 Contributions

The contributions of this work are manifold:

- The main challenge and contribution were to collect the historical and real time bike data which resulted in extracting the real time Dublin bike data from JCDecaux API managed by Dublin City Council. A python script was run for a month to retrieve data. Similarly, weather data has been collected using pwyapi library in python.
- The data was collected, cleaned and processed for modelling and has been made publicly available as open source for future use.
- Machine learning timeseries models have been developed based on multivariate attributes to analyse and predict the forecast. Previous researches include study and implementation on single feature but this thesis deals with multiple features.
- According to the previous researches and study, it is observed that there is no work done on prediction of bikes considering multiple factors like weather, event, day of the week for the city Dublin. This thesis contributes by working on these factors by modelling multivariate timeseries model and comparing their performances.

- The various factors affecting the availability of the bikes such as temperature, wind speed, precipitation, events, day of the week is considered and is visualized based on time, day and station in Tableau.

## 1.9 Thesis Outline

The rest of the thesis is organised as follows:

- *Chapter 2: State of Art* – This chapter elaborates about the timeseries related work and reviews the previous study done on the bike sharing problems. It discusses about the Regression, Long Short-Term Memory and Auto Regressive Integrated Moving Average model work done on bike sharing systems and their limitations with the takeaway points which has been implemented in the thesis, discussed later in chapter 4.
- *Chapter 3: Dataset Description* – This chapter includes brief description about the method of dataset collected, explains its attributes and the use as features for future feature engineering procedures explained in Chapter 4.
- *Chapter 4: Method and Implementation* – This chapter explains the in-depth description about the methodology of the thesis which was mentioned in section 1.6. It provides the complete knowledge about the CRISP-DM pipeline which includes elaborative explanation of all the steps from data retrieval to model evaluation along with comprehensive study of LSTM and ARIMA model.
- *Chapter 5: Results and discussions* - This chapter discusses the results of the implementation of the model, comparison of their performances and discussions about the better model for the bike series dataset.
- *Chapter 6: Conclusion and future work* – This chapter concludes the thesis with listing the advantages and limitations of the work and discusses the potential impact along with future work.
- *Chapter 7: References* – This includes all the referenced sources cited in Harvard style.