# Road and Traffic Violation Data Analytics Using Random Forest

**Cindy Agno-Balabat Aurora, Maria Visitacion Gumabay**

*Abstract--- This paper presents an in-depth analysis of road and traffic violations pattern using Data Analytics methods, aimed at improving road and traffic management, government planning and decision making. The study identified the road and traffic current management practice as basis of the design development and implementation of the road and traffic management system. The application managed all the road and traffic violation that will produce recorded set for analysis, which carried out from over of five years. Through data cleansing a total of twenty thousand six hundred forty record set was derived. It is important to find use of this record set, build analysis models, and use interactive tools to produce predictive data, understand the relevance, trends, and driving behaviors from the road and traffic violations data in terms of the following predictors: gender of the violator, vehicle owner address, location of violation, month and time the violation was committed and traffic enforcer who issued the citation. The study was able to establish a data analysis model by using a powerful classification and regression tool - random forest which was executed using an open source application named Orange. Finally, the developed application was evaluated by system users and IT experts using the ISO 25010 criteria.*

*Keywords: Road and Traffic, Violation, Data Analytics, Random Forest*

## I. INTRODUCTION

Due to increasing population, the demand of motorization also increases which further results to traffic congestion and violations. Congestion in traffic is a state wherein the road networks are characterized by slower speeds, longer trips, and an increase of vehicular queuing, which will affects both accessibility and mobility of people and goods. Moreover, the increasing economy boosts the need for mobility and thus increased the need for vehicle ownership.

Traffic volume studies have been conducted to determine the number, movements, and classifications of roadway vehicles at a given location. These data helped identify critical flow time periods, determine the influence of large vehicles or pedestrians on vehicular traffic flow, or document traffic volume trends.

The Road and Traffic Administration Department (RTA) of Cagayan de Oro City has increase it man power to address the pressing traffic problems. But no studies have ever been conducted on how to address road and traffic violation issued by Traffic Enforcer and Highway Patrol Police which usually amounts to five thousand violation ticket per month such data was never been process for analysis.

It is important to find use of this information, build analysis models, and use interactive tools to understand the

relevance, trends, and driving behaviors from the traffic violations data. This paper aimed to establish a traffic violations data analysis model by using the distributed random forest using traffic violation data.

Random Forest (RF) is a powerful classification and regression tool. When given a set of data, RF generates a forest of classification (or regression) trees, rather than a single classification (or regression) tree. It is one of the most accurate learning algorithms available. For many data sets, it produces a highly accurate classifier. [3]

Therefore, the researcher aimed to process, inspect, cleanse, and transform to develop, implement and evaluate a Road and Traffic Violation Application that will provide the data set to perform data analysis by using random forest.

*Conceptual Framework*

The conceptual framework of this study was adopted on the concept of the applied conceptual architecture of data analytics [5].
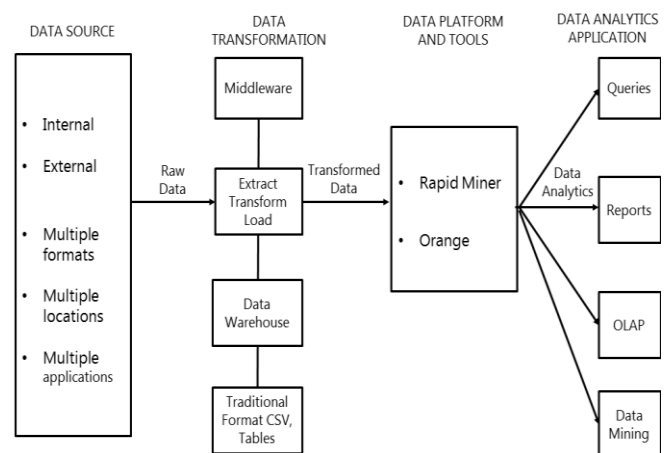


**Figure 1. An applied conceptual architecture of data analytics.**

Figure 1 indicates, a primary component is the data itself. The data can be from internal and external sources, often in multiple formats, residing at multiple locations in numerous legacy and other applications. All this data has to be pooled together for analytics purposes. The data is still in a raw state and needs to be transformed. Here, several options are available. A service-oriented architectural approach combined with web services (middleware) is one possibility. The data continues to be in the same state, and services are used to call, retrieve, and process the data. On the other hand, data warehousing is another approach wherein all the

data from the different sources are aggregated and made ready for processing. the data from However, the data is unavailable in real time. Via the steps of extract, transform, and load (ETL), diverse sources is cleansed and made ready. Depending on whether the data is structured or unstructured, several data formats can be input to the platform [4] [5].

In this next stage in the conceptual framework, several decisions are made regarding the data input approach, distributed design, tool selection, and analytics models. Finally, to the far right the four typical applications of data analytics are shown. These include queries, reports, online analytic processing (OLAP), and data mining. Visualization is an overarching theme across the four applications.

A wide variety of techniques and technologies have been developed and adapted to aggregate, manipulate, analyze, and visualize data. These techniques and technologies draw from several fields, including statistics, computer science, applied mathematics, and economics [1].

*Statement of the Problem*

The study sought to develop, implement and evaluate Web-Based Road and Traffic Citation Management System with Mobile Application that provide the data set to perform data analysis by using Random Forest.

Specifically, it sought to answer to the following:

1. What is the assessment of the participants on the current practices on the Road and Traffic Management in terms of:

1.1 traffic violation citation issuance;

1.2 management of traffic violation citation ticket; and

1.3 monitoring of traffic violation status?

2. What proposed application can be developed to enhance the current system of road and traffic administration?

3. What is the result of the data analytics in terms of the following predictors of the Road and Traffic Violations:

3.1 gender of the violator;

3.2 vehicle owner address;

3.3 location of violation;

3.4 month and time the violation was committed and

3.5 traffic enforcer who issued the citation?

4. What is the extent of compliance of the developed application in terms of the following ISO 25010 criteria?

4.1 Functional Sustainability;

4.2 Performance Efficiency;

4.3 Usability;

4.4 Reliability;

4.5 Security;

4.6 Maintainability; and

4.7 Portability.

## II. METHODOLOGY

*Research Design*

The study made use of descriptive survey and research system development to develop and evaluate an Web-Based Road and Traffic Citation Management System with Mobile Application that provide the data set to perform data analysis by using Random Forest.

*Participants of the study*

The participant in this study was chosen using purposive sampling. A purposive sample is a non-probability sample that is selected based on characteristics of a population and the objective of the study. Purposive sampling is also known as judgmental, selective, or subjective sampling. This type of sampling can be very useful in situations when you need to reach a targeted sample quickly, and where sampling for proportionality is not the main concern. A kind of purposive sampling called the expert sampling is a form of purposive sampling used when research requires one to capture knowledge rooted in a particular form of expertise. It is common to use this form of purposive sampling technique in the early stages of a research process, when the researcher is seeking to become better informed about the topic at hand before embarking on a study. Doing this kind of early-stage expert-based research can shape research questions and research design in important ways.

*Instrumentation*

The proposed system was tested by the selected IT experts and the system users. The system evaluation tools included items to describe the system in terms of functional suitability, performance efficiency, usability, reliability, security, maintainability, and portability. The main research instrument that was used in this study is the questionnaire that was developed based on the ISO 25010 quality standard. The evaluation was administered using the online version of the software.

*Data Gathering Procedure*

Before gathering and collecting data, the researcher sought the proper authorization and permission to conduct research from the Road and Traffic Administration Department. The request for proper communication protocol was observed in the dissemination of questionnaires to system evaluators and end users.

*Data Analysis*

The data collected were tabulated, analyzed, interpreted, and summarized using both descriptive and inferential statistics. The data were analyzed using the Statistical Package for Social Science for Windows (SPSS for Windows).

Frequency Count and Percentage Distribution were used to describe the demographic profile of the participants, the variety of rice used, and field practices of the farmer participants.

Mean was used to analyze the average rating of the IT experts with respect to the compliance of the application that was developed in this study in relation to the ISO standard.

## III. RESULTS AND DISCUSSIONS

The ultimate aim of this study was to develop, implement and evaluate Web-Based Road and Traffic Citation Management System with Mobile Application that provide the data set to perform data analysis by using Random Forest.

*A. Project Description*

The project was designed to operate on a web-based and

mobile application that operate on desktop computer, laptop, tablet and android smart phones. The web-based application is designed using HTML, CSS, JQuery, Javascript and PHP was used for the front-end and MySql for the back-end. The notepad++ as the IDE and Xampp for local development were used. Also, Adobe Photoshop CS3 for web/graphic design. The mobile application was developed using Android Software Development Kit (SDK) and Android Developers Tool (ADT). It uses a MySQL database for its information storage needs. The general function of the application is the management and monitoring of issued citation ticket of road and traffic violation.

*B. Web-Based Application Interface*



**Figure 2. Administrators Dashboard**



**Figure 2. Administrators Report**



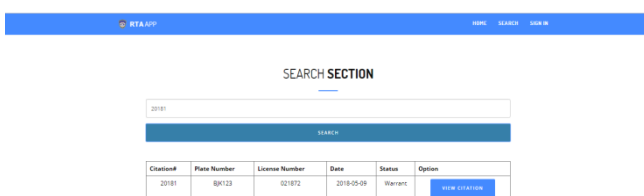**Figure 3. Traffic Enforcer - Citation Issuance**


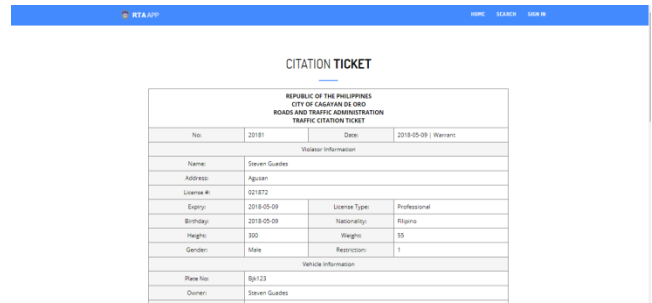
**Figure 5. Violator - Search Section**



**Figure 6. Violator - Citation Ticket**

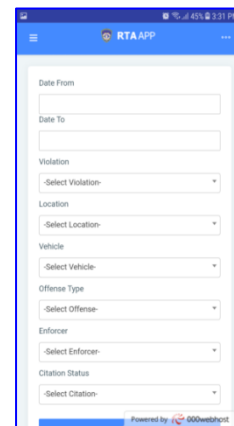*C. Mobile Application Interface*
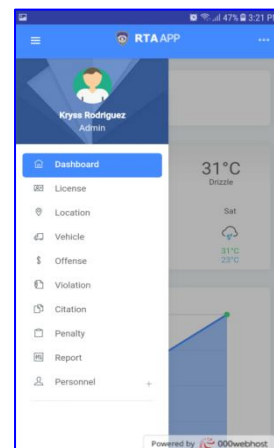


**Figure 7. Administrator Dashboard**
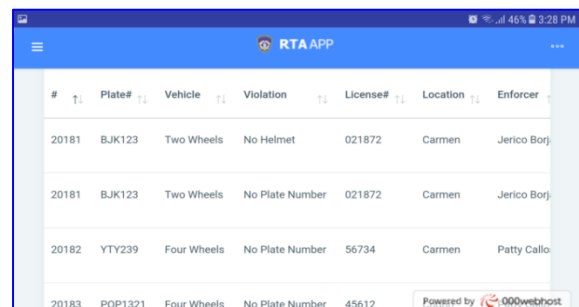


**Figure 8. Citation Issuance**



**Figure 9. Citation Report**

*D. Result of the data analytics in terms of the following*

*predictors of the Road and Traffic Violations:*

The researcher has initially gathered not less than fifty-thousand actual citation ticket. This data went through the data cleaning or data cleansing process where the researcher identifies incomplete, incorrect or irrelevant parts of the data by replacing, modifying, or removing records from the record set. After the data cleaning twenty-thousand two hundred forty-five record set remains.

Using the data analytic tool called Orange. The record set are processed using the model random forest and the result are as follow:
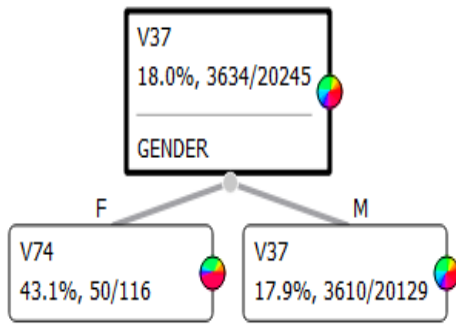


**Figure10. Result of data analytics for violation - gender.**

As shown in Figure10, out from 20245 road and traffic violation and violators, male drivers got 20129 violation with the highest violation of 17.9% or 3610 which is V37 (No helmet when driving or riding a motorcycle). While female driver out of the total 20245 violation got 116 violation 43.1% or 50 of which is V74 (Violation of Procedures Involving Traffic Accidents).
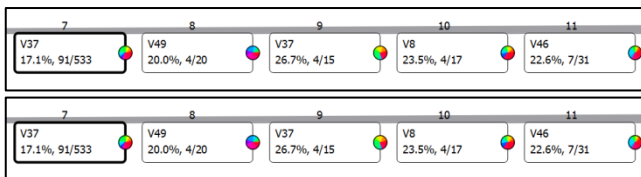


**Figure 11. Result of data analytics for violation – violators address.**

As shown in Figure 11, out from 20245 road and traffic violation and violators, 5783 of the violators are non-resident of the city, 15.7% or 906 out of 5783 of the violation is V37 (No helmet when driving or riding a motorcycle). While total 533 of the violators are residing from 7 (Barangay Carmen) 17.1% or 91 of which is V37 (No helmet when driving or riding a motorcycle).
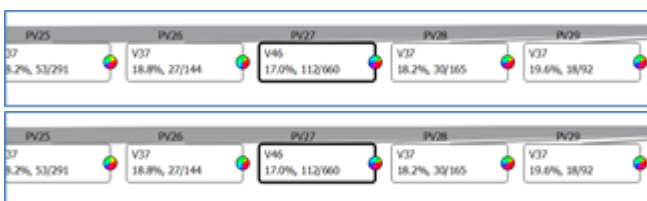


**Figure 12. Result of data analytics for violation – location of violation.**

As shown in Figure 12, out from 20245 place of violation, 819 of the violations happened in PV61 (Rizal Street,

Barangay 11), 18% or 148 out of 819 of the violation is V37 (No helmet when driving or riding a motorcycle). While 660 of the violations happened in PV27 (Osmeña Street, Barangay 31), 17% or 112 out of 660 of the violation is V46 (No Sticker). On the other hand, PV5 (Balongis, Balulang) only has a total of 13 violations 46.2% or 13 of which is V74 (Violation of Procedures Involving Traffic Accidents). PV3 (Agustin-Velez Street, Barangay 1) only has a total of 11 violations 27.3% or 3 of which is V37 V37 (No helmet when driving or riding a motorcycle).
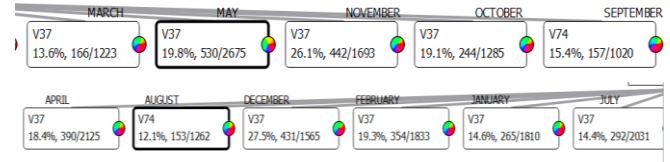


**Figure 13. Result of data analytics for violation –month.**

As shown in Figure 13, out from 20245 violation, 2675 of the violations happened in the Month of May 19.8% or 530/2675 of the violation is V37 (No helmet when driving or riding a motorcycle). While 1262 of the violations happened in Month August with 12.1% or 153 of the violation is V74 (Violation of Procedures Involving Traffic Accidents).
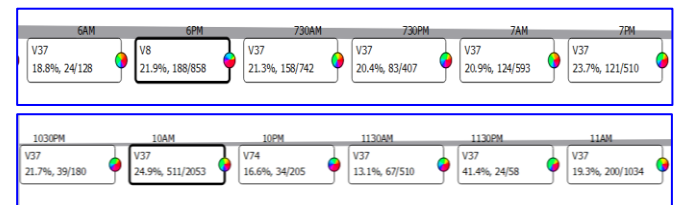


**Figure 14. Result of data analytics for violation – time.**

As shown in Figure 14, out from 20245 violation, 2053 of the violations happened during 10 in the morning with 24.9% or 511/2053 of the violation is V37 (No helmet when driving or riding a motorcycle), 858 of the violations happened during 6 in the evening with 21.9% or 188 of the violations is V8 (Disregarding traffic sign).
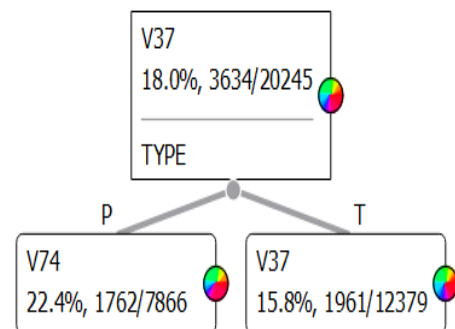


**Figure 15. Result of data analytics for violation – officer type.**

As shown in Figure 15, out from 20245 road and traffic violation 12379 was issued by traffic enforcer 15.8% or 1961 of which is V37 (No helmet when driving or riding a motorcycle), 7866 violation was issued by police 22.4% or 1762 of which is V74 (Violation of Procedures Involving Traffic Accidents).
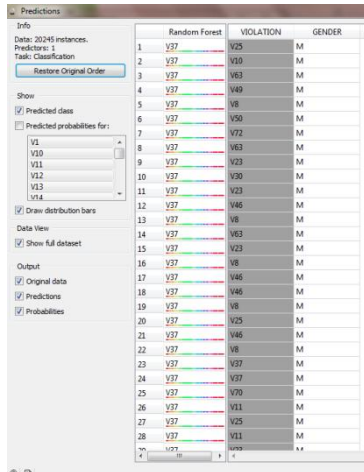
*E. Prediction Result using Random Forest*



**Figure 16. Prediction result for violation - gender.**

As shown in Figure 16, Prediction report for violation-gender using random forest. It is evident in the prediction report that male drivers are more likely to commit violation V37 (No helmet when driving or riding a motorcycle) and V74 (Violation of Procedures Involving Traffic Accidents) for female drivers.
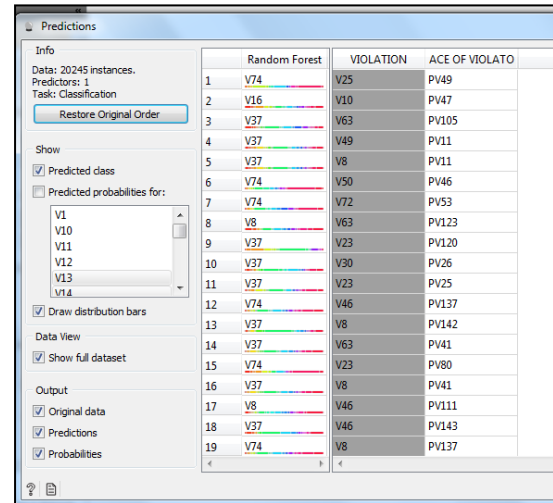


**Figure 17. Prediction result for violation – violators address.**

As shown in Figure 17, Prediction report for violation-violators address using random forest. It is evident in the prediction report that the violation V37 (No helmet when driving or riding a motorcycle) would be more likely to be the violation of this violators.



**Figure 18. Prediction result for violation – location of violation.**

As shown in Figure 18, Prediction report for violation-location of violation using random forest. It is evident in the prediction report that majority of the violation V37 (No helmet when driving or riding a motorcycle) is predicted in some location but V74 (Tampered/marked plate or stickers) are also predicted to be committed in location such as PV49 (J.R. Borja-Osmena Street, Barangay 37), PV46 (J.R. Borja-Corrales Street, Barangay 32), PV56 (Junction, Kauswagan) and PV80 (Patag), V16 (Driving with Student Permitt) in PV47 (J.R. Borja-Daumar Street, Baranga. 38), V8 (Disregarding traffic sign) in PV123 (Velez-Yacapin Street, Barangay. 29) and PV111 (Tiano-Neri Street).
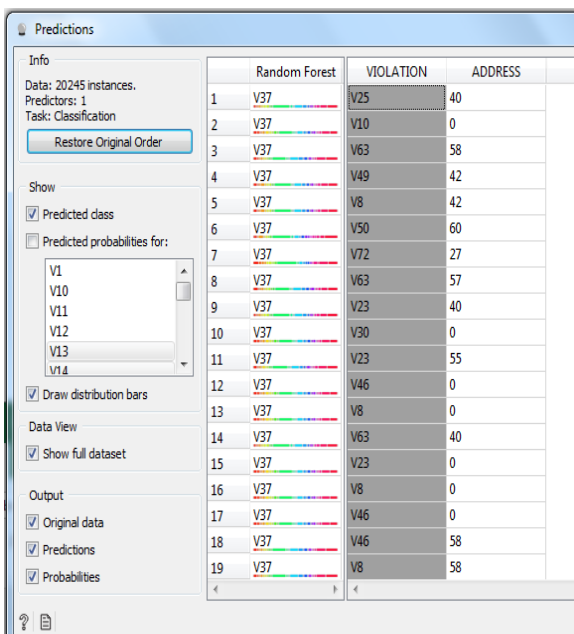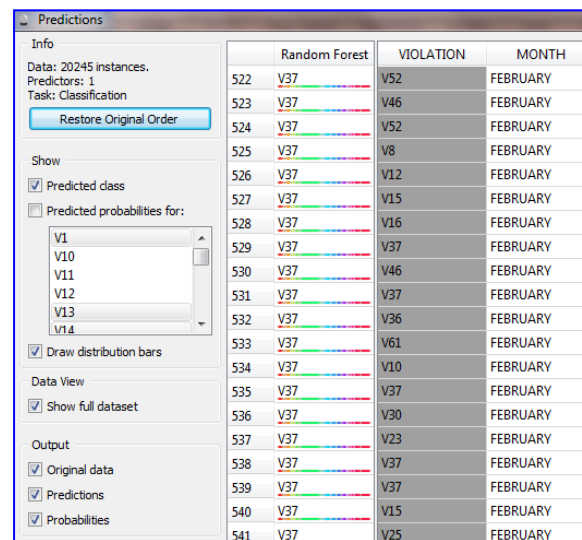


**Figure 19. Prediction result for violation – month.**

As shown in Figure 19, Prediction report for violation-month of violation using random forest. It is evident in the prediction report that majority of the violation V37 (No helmet when driving or riding a motorcycle) is predicted in majority of the months but V74 (Tampered/marked plate or stickers) are also predicted to be committed.
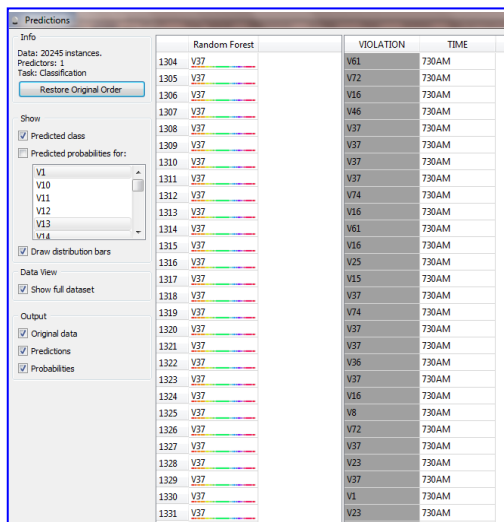
**Figure 20. Prediction result for violation – time.**

As shown in Figure 20, Prediction report for violation-time using random forest. Starting from 3 AM until 11:30 AM the most likely violation to be committed would be V37 (No helmet when driving or riding a motorcycle).
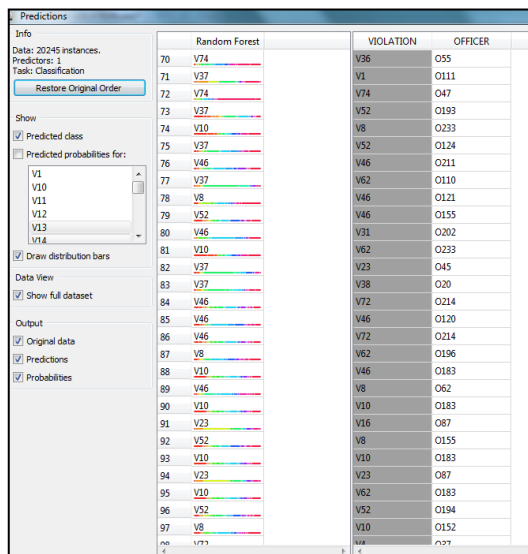


**Figure 21. Prediction result for violation – officer.**

As shown in Figure 21, Prediction report for violation-officer using random forest. It has shown that officer has various road and traffic violation issuance

*F. Evaluation Result*

**TABLE I SUMMARY ON THE EVALUATION OF THE SYSTEM USER AND EXPERT ON THE DEVELOPED SYSTEM IN TERMS OF THE ISO 25010 CRITERIA**

| Characteristics | Mean (expert) | Mean (user) | Descriptive Rating |
|---|---|---|---|
| Functional Suitability | 4.00 | 4.00 | Very Great Extent |
| Performance Efficiency | 4.00 | 4.00 | Very Great Extent |
| Usability | 4.00 | 4.00 | Very Great Extent |
| Reliability | 3.83 | 4.00 | Very Great Extent |
| Security | 4.00 | 4.00 | Very Great Extent |
| Maintainability | 4.00 | 4.00 | Very Great Extent |
| Portability | 4.00 | 4.00 | Very Great Extent |
| **Overall Mean** | **4.00** | **4.00** | **Very Great Extent** |

Table I shows the result of the evaluation of the web-based and mobile application by the system user and experts with regards to ISO 25010 criterion of functional suitability, performance efficiency, usability, reliability, security, maintainability and portability. The respondent could give a score of 1 to 4, with 1 as the lowest and 4 as the highest. The respondents gave mean score of 4.00; the score indicate a very great extent.

## IV. CONCLUSION

The web-based and mobile application provides accessibility anywhere and anytime with an Internet connection. This further extends the RTA personnel's ability to provide or retrieve data in a way that suits them. In this way, the up to date information is always at the fingertips of the people who need it.

The use of data analysis especially to agencies such as the local is not new in IT industry as cited in various systems and studies. The importance of data analysis and interpreting data is proven to be helpful in the road and traffic management, especially on reviewing their road plans and strategies or decision and policy making that can further enhance our road and traffic flow.

The findings of the study also confirmed after conducting the ISO 25010 evaluation on the proposed system with the IT experts and RTA personnel's that the proposed system's degree of capabilities to be of very great extent using the ISO 25010 criteria namely functional suitability, performance efficiency, usability, reliability, security, maintainability and portability.

## V. RECOMMENDATIONS

The following recommendations are offered based on the conclusions of the study:

1. All RTA personnel that can issue road and traffic citation ticket can be able to have access to computers especially in the RTA Department or an android mobile phone capable to run the RTA mobile app with an Internet connection of at least 1MBps in speed.

2. The office that handle the road and traffic citation ticket can be issued with more appropriate computer hardware capable to handle the proposed system.

3. All RTA personnel in charge should have user training prior to using the system or be familiar using the system through the user manual given.

4. RTA department can make use of the data analytics result as a tool to further enhance road and traffic system of the city.

5. Encode in the system the previous road and traffic citation ticket.

## REFERENCES

1. Kudyba,S. (2014). Big Data Mining and Analytics: Components of Strategic Decision Making. CRC Press.
2. Olorunfemi, S.O. (2013). Examination of On-Street Parking and Traffic Congestion Problems in Lokoja. Department of Transport Management Technology Federal University Akure Nigeria.
3. Prinzie, Anita; Poel, Dirk (2007). "Random Multiclass Classification: Generalizing Random Forests to Random MNL and Random NB". Database and Expert Systems Applications. Lecture Notes in Computer Science. 4653. pp.349. doi:10.1007/978-3-540-74469-6_35. ISBN 978-3-540-74467-2.
4. Rapnipathi, V. et.al.(2014). Big Data Analytics Architectures, Frameworks, and Tools. Retrived : http://www.ittoday.info/ITPerformanceImprovement/Articles/2014-07Raghupathi.html
5. Sathi, A. (2012). Big Data Analytics. MC Press Online LLC, 2012.
6. Zikopoulos, P.C., et. al (2013). Harness the Power of Big Data—The IBM Big Data Platform. New York: McGrawHill.