

JKAU: Comp. IT. Sci., Vol. 11 No. 2, pp: 53 – 66 (1444 A.H. / 2022 A.D.) DOI: 10.4197/Comp.11-2.5 The journal of King Abdulaziz University (Computing and Information Technology Sciences) is licensed under a Creative Commons Attribution 4.0 International License (CC-BY 4.0).

Using data analytics and machine learning for road traffic accident prediction and causal factor analysis

Abdullah Hassan Asiri

Department of Information Systems, Faculty of Computing and Information Technology King Abdulaziz University, Jeddah, Saudi Arabia ahussainasiri@stu.kau.edu.sa

Abstract. Machine learning and data analysis together represent a powerful force in producing analysis and forecasting models which in turn support the growth of different fields. In terms of safety and people, road traffic accidents have a negative impact on individuals, governments, and companies. Therefore, predicting traffic accidents on the road has become very important in the decision-making process that will lead to the safety of others. It is also important to find the causes of these accidents. This study aims to predict road traffic accidents and their causative factors based on a data set from Kaggle [1] and to create a road vehicle accident prediction model using machine learning and decision tree algorithms. The results showed that one of the most important factors causing accidents on the road is the lack of distance between vehicles, and that weather and age play an important role. The model achieved an accuracy of 84.22% and a standard deviation of +/-1.08 %. The tools used in data analysis and machine learning and prediction were Python and RapidMiner an integrated software platform for data science that provides an integrated environment for machine learning, deep learning, and predictive analysis.

Keywords — data analysis, algorithm, vehicle accidents, machine learning, prediction.

I. INTRODUCTION

The world of cars is witnessing great developments, and although modern cars are equipped with modern technology and multiple safety features, the number of car accidents is increasing dramatically throughout the world. In a 2018 report, the World Health Organization reported that the number of traffic accident deaths was 1.35 million, of which a large proportion is disproportionately shared among pedestrians, motorcyclists, and cyclists [2].

To preserve human life, it is important to understand the factors causing traffic accidents so that the world will be able to avoid them in future. Among these factors are the following: weather, roads, age, alcohol, and driver behavior. Researchers have investigated the factors causing these accidents, such as: driver behavior, road design, and weather fluctuations [3]. Although these factors are considered major factors, this study did not delve deeply into some other factors, so some other reasons and factors are still not clear. Therefore, this paper aimed to create a model to predict vehicle accidents on the road and to analyze the factors causing these accidents. The data set was obtained from Kaggle [1], in an attempt to collect the largest possible number of factors which can cause accidents such as weather fluctuations, road conditions, driver age, time of day, etc.

In this paper, various techniques were used, starting with using the Python language to clean the data, through to using Rapid Miner for data analysis and using machine learning and the decision tree algorithm for making predictions. The aim was to discover effective predictions and detect the most important factors causing vehicle accidents.

II. RELATED WORKS

Data analytics is the science of analyzing raw data to make conclusions about that information [4]. Data analytics requires the automation of a set of techniques, processes, and some algorithms on raw data to produce information for human use.

Machine learning can be broadly defined as computational methods using experience to improve performance or to make accurate predictions [5]. Experience comes in the form of electronic data representing already existing information that can be collected, cleaned, and analyzed.

Data analytics is used to avoid some problems and discover solutions or improvements to others. When reliable and clean data are used, along with the various machine learning algorithms, it is possible to obtain reliable information to make predictions. Data analytics and machine learning are used in many fields, but they are used effectively in predicting the severity of road accidents and determining the factors that cause them [6]. Specifically in Bangladesh, a study worked on examining the most important factors causing accidents and the extent of the severity of these accidents. An analysis was conducted using several tools, including Decision Tree, K-Nearest Neighbors (KNN), Naïve Bayes, and AdaBoost, to classify accidents and their severity into fatal, moderate, or minor injuries, and the best tool used in this study was AdaBoost, which gave the highest performance.

What we are witnessing in our time of population density, with the requirements of sprawling urban life, is that everyone is using the car to move around and accomplish tasks, so car accidents on the road have increased significantly. There are many factors that cause these accidents. Various studies have been conducted in a variety of ways to search for the most important of factors and their main causes.

Table I. Some studies that used different methods to study and analyze the factors causing road traffic accidents.

Objectives	Method	Results	Limitations	Study
Establishes models to select a set of influential factors and to build up a model for classifying the severity of injuries.	Machine learning algorithms, AdaBoost, Logistic Regression (LR), Naive Bayes (NB), and Random Forests (RF)	RF algorithm has shown better performance with 75.5% accuracy than LR with 74.5%, NB with 73.1%, and AdaBoost with 74.5% Accuracy.	The study focused on only a certain group of factors that cause traffic accidents. We may need more than one group to study some other factors that may be more important than some of those factors selected. In other words, we need several models.	Comparison of Machine Learning Algorithms for Predicting Traffic Accident Severity [7].
To develop a method for novel Frequent Pattern tree (FP tree) based variable selection. The method works by identifying all the frequent patterns in the traffic accident dataset.	FP tree (Bayesian) and the random forest method.	Best model found was an FP tree-based Bayesian network model that can predict 61.11% of accidents while having a false alarm rate % of 38.16.	The method used in this study to predict road traffic accidents depends on the patterns of previous accidents, and the patterns of accidents cannot be counted, or there could be repeats in some patterns, so the model needed to be continuously updated	A novel variable selection method based on frequent pattern tree for real-time traffic accident risk prediction [8]

Table II. Some studies that used different methods to study and analyze the factors causing road traffic accidents.

Objectives	Method	Results	Limitations	Study
To study the pattern of road traffic accidents, and antecedent factors influencing the road traffic accidents.	Cross- sectional study (Various parameters like age and sex distribution, time of occurrence, alcohol consumptio n, etc.)	There was a marked male preponderanc e (88.77%), Most of the accidents had taken place in the evening hours (6 pm - 12 midnight).	A cross-sectional study was used, and as we know that this method is used to collect data for a specific period and in a specific place, and therefore the data and numbers vary from month to month, from year to year, and from place to place. This study has been included to cover technical and non- technical studies and their use.	An epidemiologica l study of road traffic accident cases at a tertiary care hospital in rural Haryana [9].
This study aims to use a consensual neural network to predict with great accuracy the occurrence of traffic accidents.	Convolutio n neural network.	The prediction of (TAP-CNN) model accuracy comes to 78.5%, which is 7.7% higher than the model (TAP-BP).	The study focused on different specific sets of accident cases to train the neural network, but as also reported in this paper, it did not explicitly address the causative factors of accidents, including the road environment on which cars are traveling and so on. This method may succeed in predicting traffic accidents, but it will not tell us the causative factors to avoid them in the future	A model of traffic accident prediction based on convolutional neural network [10]

Objectives	Method	Results	Limitations	Study
The aim of this study was the analysis of human behavior or the conditions which can lead to an accident or influenced the cause of an accident.	Qualitative analysis and Quantitative analysis (Pearson s chi- squared test).	The most common contributing factor (or accident cause) was inattention (regardless of the age or gender). The Pearson chi-squared test indicated statistically significant differences between the men and women and between age groups (young drivers under 25, middle age, and seniors up to 65).	This study focused on one factor, which is human behavior, and it was supported by age and gender. It was based on qualitative and quantitative analysis. Therefore, we need to analyze some other factors using better technical methods.	Human factors contributing to the occurrence of road traffic accidents [11].

Table III. Some studies that used different methods to study and analyze the factors causing road traffic accidents.

A study references a major and very important factor, which is the age group, specifically the adolescent stage. This study stated that people between 18-24 years represent 23% of the number of deaths caused by traffic accidents [12]. Research needs to address the other age groups, especially the elderly and people with disabilities, so we need to detect, analyze, and study these other factors.

The main objective [13] is to analyze traffic accidents with the aim of developing models for accurate prediction of collision frequency in Anambra, Nigeria, using auto regressive integrated moving average (ARIMA) and auto regressive integrated moving average with explanatory variables (ARIMAX). The results showed that the ARIMAX model outperformed the ARIMA model. This study's findings reveal that incorporating human, vehicle, and environmental related factors in a time series analysis of the crash dataset produces a more robust predictive model than using the aggregated crash count. Some of the classification models used in Mohanta's study [14], specifically Logistic Regression, Artificial Neural Network, Decision Tree, K-Nearest Neighbors, and Random Forest have been implemented to predict the severity of accidents. These models have been verified, and the first results prove that these classification models have attained considerable accuracy.

III. METHOD

In this study, a variety of techniques for analysis and prediction were used, starting with the Python language to clean the data because it is easy to use, accurate and highly flexible. Then Rapid Miner was used to analyze the data and create a prediction model using the decision tree. Rapid Miner was chosen because of its ease of use and fast creation of models and analysis. In this section we will see some more details.



Fig 1. Workflow

A. Data source

The data for car accidents in Addis Ababa for the years 2017-2020 was collected by Shahane [15], obtained from Kaggle [1]. The data consisted of 32 columns and 12,316 rows. The dataset contains many factors that cause traffic accidents on the road, such as the age of the driver, the type of road, the weather condition, the driver's experience, the driver's gender, etc.

• Data description

The descriptive table shows us the count, uniqueness, top entry, and frequency for each column. For example, the time column contains 12,316 rows, the number of unique entries is 1074, the top entry is 15:30 with a frequency of 120, which means that 15:30 is repeated 120 times in this column. The other columns work in the same way. Another example is the driver gender column containing 12,316 rows, the number of unique entries is 3 (males, females and unknown), the top entry is male with a frequency of 11,437.

		Time	Day_of_week	Age_band_of_driver	Sex_of_driver	Educational_level	Vehicle_driver_relation	Driving_experience	Type_of_vehicle	Owner_of_vehic	le Service_y	ear_of_vehicle	
	count	12316	12316	12316	12316	11575	11737	11487	11366	1183	14	8388	
	unique	1074	7	5	3	7	4	7	17		4	6	
	top	15:30:00	Friday	18-30	Male	Junier high school	Employee	5-10yr	Automobile	Own	er	Unknown	
	freq	120	2041	4271	11437	7619	9627	3363	3205	1045	19	2883	
	4 rows x	30 colur	ns										
: 01	id_data . Vehici	30 colum describ e_moverne	NS =(include=[* nt Casualty_o	object']) lass Sex_of_casualty	Age_band_of_	casualty Casualty	severity Work_of_casua	ity Fitness_of_casu	ility Pedestrian_r	novement Cause,	of_accident	Accident_seve	ri
: 0	id_data . Wehici	30 colum describ e_movem 120	NS e(include=[* nt Casualty_o 08 1	abject']) lass Sex_of_casualty 2316 12316	Age_band_of_	casualty Casualty	sevenity Work_of_casua 12316 9	ity Fitness_of_casu	iity Pedestrian_r 681	novement Cause, 12316	of_accident 12316	Accident_seve	ri 3
	id_data . Vehici	30 colum describ e_movem 120	ns =(include=[" nt Casualty_c 08 1 13	zbject")) Hass Sex_of_casualty 1316 12316 4 3	Age_band_of_	casualty Casualty 12316 6	sevenity Work_of_casua 12316 9 4	ity Fitness_of_casu	iity Pedestrian_1 681 5	novement Cause, 12316 9	.of_accident 12316 20	Accident_seve	ni 3
- 01 	id_data . Wehici	30 colum describ e_moverne 120 Going straig	ns +(include=(* nt Casualty_ 08 1 13 ht Driver or	object')) lass Sex_of_casualty 1316 12316 4 3 ider Male	Age_band_of_	casualty Casualty 12316 6 na	severity Work_of_cessue 12316 9 4 3 Dr	ity Fitness_of_casu 118 S 7 Ver No	i lity Pedestrian_r 681 5 mal Notai	novement Cause, 12316 9 Pedestrian I	of_accident 12316 20 No distancing	Accident_seve 12 Slight In	rit 31

Fig 2. Descriptive table

B. Data cleaning

Data cleaning is a necessary part of the data analysis procedure. As we saw in the data description section, there are many missing data values in several columns and the total number of missing data values in the complete data set is 20,057, and there are also completely missing rows and useless columns. This empty data needs to be removed so that the missing data can become complete and reliable data. In this section, we will discuss in detail the data cleaning process.

Python is an interpreted, object-oriented, highlevel programming language with dynamic semantics, a multi-functional, maximally interpreted programming language with several advantages that are often used to streamline massive and complex data sets.

Jupyter notebooks provide an easy-to-use private data science environment that can be used in presentations or educational settings. Jupyter is often used with Python because of its flexibility and efficiency. Python and the Jupyter tool were used to understand and clean the data.

This code allows us to call the Panda library, which is one of the Python libraries that deals with data, and then read the data file via the file path.

	11210	Insert	Cel K	erne wogets	Help				Trusted	Python 3 (pykeme)
+ ×	0.6	+ +	► Bun	C 🗰 Code	× =	1				
In (1)	impor	t pandas	s as pd							
In (2)	old_d	ata= pd.	.read_csv('	/Users/abdullat	ha/Desktop/	adisababa_data	/olddata.csv')			
In (3)	old d	ata								
Out[3]:		Time	Day_of_week	Age_band_of_driver	Sex_of_driver	Educational_level	Vehicle_driver_relation	Driving_experience	Type_of_vehicle	Owner_of_vehicle
	0	17:02:00	Monday	18-30	Maie	Above high school	Employee	1-2yr	Automobile	Owner
	1	17:02:00	Monday	31-50	Male	Junior high school	Employee	Above 10yr	Public (> 45 seats)	Owner
	2	17:02:00	Monday	18-30	Male	Junior high school	Employee	1-2yr	Lony (417100Q)	Owner
	3	1:06:00	Sunday	18-30	Male	Junior high school	Employee	5-10yr	Public (> 45 seats)	Governmental
	4	1:06:00	Sunday	18-30	Male	Junior high school	Employee	2-5yr	NaN	Owner
	-		-	-		2	-			-
	12311	16:15:00	Wednesday	31-50	Male	NaN	Employee	2-5yr	Lony (11740Q)	Owner
	12312	18:00:00	Sunday	Unknown	Male	Bementary school	Employee	5-10yr	Automobile	Owner
	12313	13:55:00	Sunday	Over 51	Male	Junior high school	Employee	5-10yr	Bajaj	Owner
	12314	13:55:00	Sunday	18-30	Female	Junior high school	Employee	Above 10yr	Lony (417100Q)	Owner
	12315	13:55:00	Sunday	18-30	Male	Junior high school	Employee	5-10yr	Other	Owner
Г	12316	rows × 32	columns							

Fig 3. Data read

Missing data

The missing data may be incomprehensible text or numeric data, and it may be corrupt data or empty fields that distort the data set.

This code shows us each column in the data with the number of missing values for each column.



Fig 4. Number of missing data per column

	. 18	View	Inset	Cell	Kernel	V	Vidgets	Help					T	usted	- 1	Python 3	3 (ipyka	ornei)
s +	k d	6	+ +	► Bun			Code	~	83									
	1	Tehicle	driver	relat	ion		579											
	-	riving	experi	ence			829											
		ype_of	_vehicl	e			950											
	1.1	wner_c	of_vehic	le			482											
	-	lervice	year_o	f_vehi	cle		3928											
		erect.	of_vehi	cle	31		4427											
	- 12	trea_ac	cident_	occure	3		239											
		anes c	1 increase	ns			303											
		bines d	f Junch	ion			897											
		load a	irface t	vne			172											
		toad st	irface c	onditio	ons		0											
		ight o	conditio	ns			0											
	1	leather	_condit	ions			0											
		ype_of	collis	ion			155											
	1	umber_	of_vehi	cles_i	svolve	\$	0											
	1	lumber_	of_casu	alties			0											
		Tehicle	noveme	nt			308											
	- 13	asualt	y_class				0											
		ex_or	casualt	y			0											
	10	agual4	W BAUAT	itu			0											
		lork of	casual	inu			3198											
		itness	of cas	uality			2635											
		edest	ian nov	ement			0											
	- 04	ause_c	f_accid	ent			0											
		locides	it_sever	ity			0											
		itype:	int64				1.1.1	_										

Fig 5. Missing values in all rows and columns.

Seven unwanted columns were removed until we were left with 25 out of 32 columns. The unwanted columns do not add any value to the analysis process and do not help in the search for the causes of accidents. For example, the Fitness_of_casuality column. This column talks about the fitness of casualty, but this does not help the process of analysis and detection of the causes of accidents. Another example, Type_of_vehicle, is a column that talks about the type of vehicle this also does not the help analysis process, etc.

ie Edt	View	Insert	Cell K	ernel Widgets	Help				Trusted	Python 3 (ipykerne
+ >	0	* •	▶ Run	C 🗰 Code	۷					
In [8]: 1	lata=o.	ld_data.	drop(('Typ	e_of_vehicle',	'Service_y	ear_of_vehicle	','Fitness_of_cas	uality','Defec	t_of_vehicle'	'Casualty_se
In [9]:	new_d	ata								
Out[9]:		Time	Day_of_week	Age_band_of_driver	Sex_of_driver	Educational_level	Vehicle_driver_relation	Driving_experience	Owner_of_vehicle	Area_accident_
	0	17:02:00	Monday	18-30	Male	Above high school	Employee	1-2yr	Owner	Resident
	1	17:02:00	Monday	31-50	Male	Junior high school	Employee	Above 10yr	Owner	05
	2	17:02:00	Monday	18-30	Male	Junior high school	Employee	1-2yr	Owner	Recreation
	3	1:06:00	Sunday	18-30	Male	Junior high school	Employee	5-10yr	Governmental	05
	4	1:06:00	Sunday	18-30	Male	Junior high school	Employee	2-5yr	Owner	industr
	-	-		-	-		-	-	-	
	12311	16:15:00	Wednesday	31-50	Male	NaN	Employee	2-5yr	Owner	Outside ru
	12312	18:00:00	Sunday	Unknown	Male	Elementary school	Employee	5-10yr	Owner	Outside ru
	12313	13:55:00	Sunday	Over 51	Male	Junior high school	Employee	5-10yr	Owner	Outside ru
	12314	13:55:00	Sunday	18-30	Female	Junior high school	Employee	Above 10yr	Owner	05
_	12315	13:55:00	Sunday	18-30	Male	Junior high school	Employee	5-10yr	Owner	Outside ru
	12316	rows × 25	columns							

Fig 6. Deleted columns

This code deletes any row that contains at least one missing field.

File Edit	Ver K	Insert	Cel	Kernel W	idgets Help				Trusted	Python 3 (pykernel)
3 7 6	0.0	T V	P nut		C008 *	8				
In (24):	cleaned	_data=c	ev_data	() () () ()				_		
In [25]	cleaned	_data.t	o_csv{	/Users/abd	ullaha/Desktop/	adisababa_d	lata/ new_data	(.csv*)		
Out (25) :	ı	innamed: 0	Time	Day_of_week	Age_band_of_driver	Sex_of_driver	Educational_level	Vehicle_driver_relation	Driving_experience	Owner_of_vehicle Are
	3	3	1:06:00	Sunday	18-30	Male	Junior high school	Employee	5-10yr	Governmental
	8	8	17:20:00	Friday	18-30	Maie	Junior high school	Employee	Above 10yr	Owner
	9	9	17:20:00	Friday	18-30	Male	Junior high school	Employee	1-2yr	Owner
	11	11	14:40:00	Saturday	31-50	Male	Above high school	Employee	No Libence	Owner
	14	14	17:45:00	Thursday	31-50	Male	Junior high school	Employee	Above 10yr	Owner
	141	-		-	1.14	- 14-		-		
	12304	12304	7:10:00	Friday	18-30	Male	Junior high school	Employee	2-6yr	Owner
	12305	12305	7:10:00	Friday	18-30	Male	Junior high school	Employee	2-5yr	Owner
	12309	12309	9.05:00	Friday	31-50	Female	Elementary school	Employee	5-10yr	Owner
	12313	12313	13:55:00	Sunday	Over 51	Male	Junior high school	Employee	5-10yr	Owner
	12315	12315	13:55:00	Sunday	18-30	Male	Junior high school	Employee	5-10yr	Owner
	4658 row	s x 26 co	lumns							





After the data cleaning process, we see in the descriptive table that the number of columns has shrunk to 25 out of 32 columns, and we see that the number of rows has become 4,658 rows instead of 12,316 rows. Also, the number of rows in the data after cleaning is equal in all columns, which is the opposite of what we had before cleaning. There was a discrepancy in the number of rows per column.

Clean	ed_data.	describe(inc	lude=['object'])								
	Time	Day_of_week	Age_band_of_driver	Sex_of_driver	Educational_level	Vehicle_driver_relation	Driving_experience	Owner_of_vehicle	Area_accident_occ	cured Lanes_or_Median	s
count	4658	4658	4658	4658	4658	4658	4658	4658	. 4	4658 465	3
unique	927	7	4	2	6	3	7	4		13	i
top	18:30:00	Friday	18-30	Male	Junior high school	Employee	5-10yr	Owner		Two-way (divide Other with broken line road marking	i i
freq	42	742	1908	4622	3068	3868	1344	4081		1509 172	1
Cle	aned_d Road_s	ata.descr urface_condi	ibe(include=[fions Light_condit	'object']; ions Weathe) r_conditions Ty	pe_of_collision Vehi	cle_movement C	isualty_class Se	x_of_casualty Ag	ge_band_of_casualty	
3			1070								
			4006	4658	4658	4658	4658	4658	4658	4658	
3			4006 4	4	4658 8	4658 9	4658 12	4658 3	4658 2	4658	
; ; ;)			4008 4 4 Dry Day	4 4 fight	4658 8 Normal	4658 9 Vehicle with vehicle collision	4658 12 Going straight	4658 3 Driver or rider	4658 2 Male	4658 5 18-30	
	count unique top freq 4 rows > : Cle 3	Time count 4658 unique 927 top 18:30:00 freq 42 4 rows × 23 colum : Cleaned_d	Time Day_of_week count 4558 4558 wnique 927 7 top 153000 Friday free 42 742 drams x 23 columns Cleaned_data.descr L Road, surface_condition	Time Day,d_yeek Ap,3and_d_dhire count 4553 4555 4653 unique 227 7 4 top 183000 Friday 18-30 free 42 742 1068 stass x 23 columns Cleaned_data.describe [include=[i = Road surface_conditions Light_conditions	Time Day, af yeek Age, band, af, driver Sex, af, driver count 4553 4558 4558 4558 unique 327 7 4 2 top 1930:00 Friday 18-30 Male free 42 742 1908 4522 class x 23 columes Cleaned_data.describe (include=['object'] 4 i Road, surface_conditions Light conditions Weather	Time Day, adjustek Age, Land, adjustive Sex, adjustive Educational juent count 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4559 4558 4558 4558 4558 4558 4558 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 4559 <	Time Day,af-yeek App.land,af.driver Soc.af.,driver Educational,level Vehicle,driver_relation count 4558 4658 4658 4658 4658 wrigue 827 7 4 2 6 3 top 183000 Riday 18-30 Male Joint high school Employee fine 42 742 1958 4522 3568 3868 cross x 22 columes State 3568 3868 file datadescribe (include=['object ']) 3688 3688	Time Day_of_week App. band, of_chive Filter Educational_invel Velocity Velocity	Time Doy_of_week App.band_of_chrive: Sex_of_chrive: Educational_invel Vehicle_driver_priation: Driving_superiore: Owner_pl_which count: 4558 4558 4558 4558 4558 4558 winpue 507 7 4 2 6 3 7 4 top 18/30:00 Friday 19-30 Maie Jour high school Employee 5-10y Owner fine 4/2 742 19:08 4222 50:68 3584 40:01 fines x22 columes 5-10y Owner Cleased_data_describe (includes['object']) s Road surface_conditions Light conditions Weather_conditions Type of collision Weikele_movement Casually_class So	Time Day_of_yeek App.band.pl_chiner Soc.pl_finiter Educational_level Vehicla_driver_pristics Driving_reperience Dense_pl_vehicle Area_accident_accident_coc count 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558 4558	Time Day, d/yeek App, band, d/ drive Sex, ed, drive Educational, Jueel Welde, driver, pratrice Driver, ed., which is Ana, societart, sociand Lanes, yell Median count 4558 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658 4658

Fig 9. Descriptive table after cleaning

C. Data analysis

The data collector collected data from manual records of road accidents for the years 2017-20, and all sensitive information was excluded. The data

collector also expressed the belief that these are the most important causes of traffic accidents based on what he found in the manual records.

This chart contains the causes of accidents, and the biggest factor causing accidents is "No distancing," as the number of accidents due to this factor reached 861, which is followed by "Changing the lane to the right" with 677 accidents. The least relevant factor causing accidents was "Improper parking" with 8 accidents, then "Drunk driving" with 10 accidents.



Fig 10. The causes of accidents

Here is the driver's gender column. This column shows us the gender of the drivers who caused accidents with any cause. The number of males who were one of the parties to an accident was 4,622, but there were only 36 females.



Fig 11. Driver's gender

This chart shows us the most vulnerable group to accidents, and they are the youth group in this case. The ages of the drivers who caused the accident were divided into four groups. 1,908 accidents involved drivers aged between 18-30 years, followed by drivers aged 31-50 years, where the number of accidents reached 1,731, then drivers over 51 years, where the number of accidents reached 662, and finally those under the age of 18 years, where number of accidents reached 357.





Weather conditions are an important reason that may cause traffic accidents, depending on the condition: rain, fog, dust, etc. In the weather column, the number of accidents is shown by different weather conditions. As "normal weather" has the highest proportion with 3,882 accidents, then the "rainy" state with 516 accidents, then followed by "other weather" with 108 accidents, then the "cloudy" case with 57 accidents, then "winds" with 45 accidents, "snow" with 25 accidents, and finally "rain and wind" with 20 accidents.



Fig 13. Weather conditions

The type of road the vehicle is traveling on is also an important factor that may cause accidents. There are several different circumstances such as dry, wet, snowy, and more. "Dry roads" got the highest rank with 3,482 accidents, then "wet roads" with 1,152 accidents, then "snowy roads" with 23 accidents, and "floods" with only one accident.



Fig 14. Road conditions

The lights on the road and on the vehicle are other important traffic safety factors for drivers on the road, and they reduce risk while driving. This chart shows us the accidents that occurred in different lighting conditions: the number of accidents in "daylight" reached 3,347, and the number of accidents in "darkness with lights lit" was 1,232, followed by "darkness without lighting" with 72 accidents, and finally "darkness with unlit lights" with 7 accidents.



Fig 15. Lighting on the road

D. Machine learning

To create machine learning, the cleaned data was input, and a duplication filter was applied to remove any duplicate data. This operator removes duplicate examples from an ExampleSet by comparing all examples with each other based on the specified attributes. Two examples are considered duplicate if the selected attributes have the same values.



Fig 16. Cleaned data and a duplication filter

The search is a classification problem that we need for our predictions. A decision tree was chosen, and the decision tree algorithm was applied to the data

that had been cleaned. This Operator generates a decision tree model, which can be used for classification and regression.



Fig 17. The decision tree algorithm

The huge decision tree could be used to support decision making manually, but this will be tedious, tiring and take a lot of time, so we will keep applying the model in the next steps until the decision support becomes automated after a short time.



Fig 18. Huge decision tree

To apply the model, the data was divided into two parts. The first section contains 70% of the entries used for training the model. The training data is the initial dataset used to train the machine learning algorithms. The second part contains 30% of the entries used for testing the model. The test data is unseen data to test your model and evaluate its performance. The Split Data operator is used to partition data into subsets according to the specified relative sizes.



Fig 19. Applying the model

The weather column was chosen to predict accidents according to the weather conditions. There were several reasons behind choosing the weather column as the basis for prediction, including that the weather is one of the most important factors causing accidents. Also, the weather column is a comprehensive list of all different weather conditions. As we see in Figure 20, two columns are colored in green: the first column is the original data, and the second column is the prediction data.

sult History		ExampleSet (Apply Mod	e0 ×							Repository X	
-	Ones in	Turbo Pres	Model			Filter (1	197 / 1.197 exat	te the	,	🗘 Import Data	Ξ,
Data										+ Training Resources Icon	neradi
	Row No.	Weather_conditions	prediction/Weather_conditions)	confidencei	confidence(confidence(confidence(_	confidence(confidenc	🕴 🚞 Samples	
	605	Cloudy	Normal	0.863	0.088	0.002	0.020	0.009	0.009	Community Samples to	mediad
Σ	606	Raining	Normal	0.863	0.068	0.002	0.020	0.009	0.029	Local Repository Croit My Period Local Repository	
Statistics	607	Normal	Normal	0.863	0.088	0.002	0.020	0.009	0.009	+ Corrections	-
	608	Normal	Normal	0.863	0.088	0.002	0.020	0.009	0.009	🕶 🛅 Data	
	609	Normal	Normal	0.863	0.068	500.0	0.020	0.009	0.009	Cleaned_data 💷	(2)22.9
	610	Normal	Normal	0.863	0.088	0.002	0.020	0.009	0.009	* Processes	
	611	Normal	Normal	0.863	0.088	0.002	0.020	0.009	0.009	d 02- huld deckin	n tree
1	612	Raining	Raining	0.125	0.875	0	0	0	0	a 03- apply the mo	del 11
1	613	Raining	Kaining	0.125	0.875	0	0	0	0	d 04- testing the m	odel (
nnotations	614	Normal	Normal	0.863	0.068	0.002	0.020	0.009	0.009	@ 05- cross validat	on I III
	615	Raining	Normal	0.863	0.058	0.002	0.020	0.009	0.009	• B DE Lepert	
	616	Normal	Normal	0.863	0.088	0.002	0.020	0.009	0.009		
	617	Raining	Normal	0.863	0.068	0.002	0.020	0.009	0.009		
	618	Normal	Normal	0.973	0.002	0.008	0.012	0.002	0.004		
	619	Normal	Normal	0.863	0.088	500.0	0.020	0.009	0.079		
	620	Normal	Normal	0.863	0.088	0.002	0.020	0.009	0.029		
	621	Normal	Normal	0.863	0.058	0.002	0.020	0.009	0.009		
	622	Normal	Normal	0.863	0.088	0.002	0.020	0.009	0.009		
	623	Normal	Normal	0.863	0.088	0.002	0.020	0.009	0.009		
	624	Normal	Normal	0.863	0.058	0.002	0.020	0.009	0.009		
	625	Normal	Normal	0.863	0.055	0.002	0.020	0.009	0.009		
	626	Normal	Normal	0.863	0.088	0.002	0.020	0.009	0.079 V		

Fig 20. Original data and prediction data

In this step we need an operator to measure the accuracy of the model's performance for the training data, so we need a performance operator. This operator is used for performance evaluation. It delivers a list of performance criteria values.



Fig 21. Performance

The accuracy of the model that we obtained was 84.75%. This is the accuracy of the model based on the training data, which totaled 70% of the entire data. Therefore, this should represent the accuracy of the model training.

	• •			Vexs: De	sign R	iesults T	urbo Prep	Auto Model	Deployments		Find da	ia, operators…etc	PALS	studio *
Result History	% Perform.	unceVector (Pe	formance)	×								Repository	×	
OZ.	criterion accuracy	Table View	O Plot View									🗘 Import I	lata	
40 Performance	kappa	accuracy: 84	75N									 Training Res Samples 	ounces (con	rected)
			true Normal	true Raining	true Rainin	true Other	true Windy	true Cloudy	true Snow	true Fog or	class preci	Community 1	amples in	mered
3		pred. Nor	1158	131	6	28	13	17	7	1	85.08%	 Local Reposit We Deplace to 	tory (col)	and the state
Description		pred. Rain	7	24	0	2	1	0	0	0	70.59%	+ Connection	ns	and some
		pred. Rain	0	0	0	0	0	0	0	0	0.00%	💌 🛅 Data		
3		pred. Other	0	0	0	z	0	0	0	0	100.00%	Cleane	d_data (1)	(2/22 9:48)
Annotations		pred. Windy	0	0	0	0	0	0	0	0	0.00K	0° 01-im	port cleane	ed data (1
		pred. Cloudy	0	0	0	0	0	0	0	0	0.00K	💣 02 - bi	Ad decisio	n tree 🔅
		pred. Snow	0	0	0	0	0	0	0	0	0.00%	() ⁴ 03- 20	ply the m stinn the m	odel (10)4 notel (10)
		pred. Fog	0	0	0	0	0	0	0	0	0.00%	ef 05- cr	oss validat	tion (10/4/2
		class recall	99.40X	15.48%	0.00%	6.25%	0.00%	0.00%	0.00%	0.02%		🕨 📕 DB (Legacy)		
												<		>

Fig 22. The accuracy of the model training

To obtain true model accuracy based on test data, a cross validation must be applied. The crossvalidation operator performs a cross validation to estimate the statistical performance of a learning model.

Cross-Validation is a statistical method of evaluating and comparing learning algorithms by dividing data into two segments: one used to learn or train a model and the other used to validate the model.



Fig 23. Cross-Validation

The model's accuracy is 84.22% and the standard deviation is +/- 1.08%. The lower the standard deviation, the more stable the model. The model's accuracy is not bad, but this needs to be proven. In the weather column, there are 8 different elements: "normal weather," "Rain," "snow," "wind," "rain and wind," "cloudy," "fog," and "others". In addition to that, there is a large amount of data. To be clearer, the driver's gender column contains "male" and "female" and the model's accuracy is

99.23% and the standard deviation is $\pm 0.11\%$, but the driver's gender column is weak and ineffective for predicting accidents. Thus, we conclude that the greater the number of different elements in the target column, the lower the accuracy of the model and vice versa.



Fig 24. The model's accuracy

IV. RESULTS AND DISCUSSION

The data analysis showed that there are multiple factors in the causes of accidents. In Figure 25, the results of all possible influential factors of accidents are presented. The most influential factors are "changing lane to right," "moving backward," "No distancing," "No priority to vehicle," "changing lane to left," and "driving carelessly". The order of the influential factors of accidents is shown in Figure 26. The order shows that "No distance between vehicles on the road" is the most common factor that causes traffic accidents on the road, with a rate of 861 accidents.



Fig 25. The results of all factors of accidents



Fig 26. The order of the most influential factors of accidents

There are other factors that can be influential factors of accidents, which are drivers' gender and drivers' age. In Figure 27, the drivers' gender percentage shows that "male" had a high percentage, which is equal to 99%.



Fig 27. Percentage of drivers' gender

In addition, the drivers' age plays a major role. Figure 28 presents the percentages of drivers' ages. The age group that caused the most accidents was "between 18-30" years, which is equal to 41%, followed by "between 31-50" years.



Fig 28. The percentages of drivers' ages

One of the unexpected results shown in Figure 29, is that dry roads are the most common road type that causes accidents. In addition, most accidents happen in normal weather and daylight as presented in Figure 30 and Figure 31. This gives us an indication that accidents do not depend on certain factors or special circumstances, as these accidents abound in normal conditions and factors.



Fig 29. Percentages of road surface conditions



Fig 30. Percentages of weather conditions



Fig 31. Curve of light conditions

The model used a decision tree algorithm based on the weather column, and weather is an important factor in predicting vehicle accidents on the road. This column contains 8 sub-factors which are the different weather conditions. This model achieved an accuracy rate of 84.22%, exceeding several previous studies by a difference of up to 10 degrees Additionally, this model is considered more stable based on the standard deviation achieved of +/-1.08%. The closer this number is to zero, the more stable the model.

V. CONCLUSION AND FUTURE WORK

It would be interesting to create more than one model using different machine learning algorithms and then compare the results of each algorithm with the other algorithms. It would also be worthwhile to create a dashboard to track data for a specific process, for example tracking the number of incidents on a quarterly basis.

Acknowledgements

I would like to express my gratitude to my supervisor, Dr. Osama Rabie, who guided me throughout this project and spared no effort to help me. I would also like to thank Professor Mohammad Ramadan, who guided me at the beginning of this project. And thanks to my family who have supported me throughout this journey.

References

- [1] Kaggle, "https://www.kaggle.com/datasets," 2010. [Online].
- [2] World Health Organization. Available: https://www.who.int/publications/i/item/9789 241565684, 2018. [Online].
- [3] Rolison, J. J., Regev. S., Moutari. S. & Feeney. A.
 "https://www.sciencedirect.com/science/articl e/pii/S0001457518300873," 2018.
- [4] Runkler, T. A., "https://link.springer.com/book/10.1007/978-3-658-14075-5," 2016.
- [5] Mohri, M., Rostamizadeh, A. & Talwalkar, A. "https://cs.nyu.edu/~mohri/mlbook/," 2018.
- [6] Labib, M. F., Rifat, A. S., Hossain, M. M., Das, A. K. & Nawrine, F. "https://ieeexplore.ieee.org/document/884364 0," 2019.
- [7] AlMamlook, R. E., Kwayu, K. M., Alkasisbeh, M. R., & Frefer, A. A. "https://ieeexplore.ieee.org/abstract/document /8717393," 2019.
- [8] Lin, L., Wang, Q. & Sadek, A. W. "https://www.sciencedirect.com/science/articl e/abs/pii/S0968090X15000947," 2015.

- [9] Singh, A., Bhardwaj, A., Pathak, R. & Ahluwalia, S.
 "https://www.semanticscholar.org/paper/AN-EPIDEMIOLOGICAL-STUDY-OF-ROAD-TRAFFIC-ACCIDENT-A-Singh-Bhardwaj/a0a2fdd8c7bc9b0209d11b7b1d4df 0b9d2f90adb," 2012.
- [10] Wenqi, L., Dongyu, L. & Menghua, Y. "https://ieeexplore.ieee.org/document/805690 8," 2017.
- Bucsuházya, K., Matuchová, E., Zůvalaa, R., Moravcová, P., Kostíková, M. & Mikuleca, R. "https://www.sciencedirect.com/science/articl e/pii/S2352146520302192," 2020.
- [12] Gicquel, L., Ordonneau, P., Blot, E., Toillon, C., Ingrand & P. Romo, L. "https://pubmed.ncbi.nlm.nih.gov/28620324/, " 2017.

- [13] Ihueze, C. & Onwurah, O., "https://www.sciencedirect.com/science/articl e/abs/pii/S0001457517304542," 2018.
- [14] Mohanta, B. K., Jena, D., Mohapatra, N., Ramasubbareddy, S., Rawal, B.S., "https://content.iospress.com/articles/journalof-intelligent-and-fuzzy-systems/ifs189743," 2022.
- [15] Shahane, S. "Kaggle," 2020. [Online]. Available: https://www.kaggle.com/datasets/saurabhsha hane/road-traffic-accidents.

استخدام تحليلات البيانات والتعلم الآلي للتنبؤ بحوادث الطرق وتحليل العوامل السببية

عبدالله حسن العسيري

قسم نظم المعلومات، كلية الحاسبات وتقنية المعلومات، جامعة الملك عبد العزيز، جدة، المملكة العربية السعودية

مستخلص. يمثل التعلم الآلي وتحليل البيانات معًا قوة قوية في إنتاج نماذج التحليل والتنبؤ والتي بدورها تدعم نمو المجالات المختلفة. من حيث السلامة والأشخاص، فإن حوادث المرور على الطرق لها تأثير سلبي على الأفراد والحكومات والشركات. لذلك، أصبح التنبؤ بالحوادث المرورية على الطريق أمرًا مهمًا للغاية في عملية اتخاذ القرار الذي سيؤدي إلى سلامة الأخرين. من المهم أيضًا معرفة أسباب هذه الحوادث. تهدف هذه الدراسة إلى التنبؤ بحوادث الطرق والعوامل المسببة لها بناءً على من المهم أيضًا معرفة أسباب هذه الحوادث. تهدف هذه الدراسة إلى التنبؤ بحوادث الطرق والعوامل المسببة لها بناءً على من المهم أيضًا معرفة أسباب هذه الحوادث. تهدف هذه الدراسة إلى التنبؤ بحوادث الطرق والعوامل المسببة لها بناءً على مجموعة بيانات من [1] Kaggle وإنشاء نموذج للتنبؤ بحوادث الطرق باستخدام التعلم الآلي وخوارزميات شجرة القرار . وأظهرت النتائج أن أحد أهم العوامل المسببة للحوادث على الطريق هو عدم وجود مسافة بين المركبات، وأن الطقس والعمر وأظهرت النتائج أن أحد أهم العوامل المسببة للحوادث على الطريق هو عدم وجود مسافة بين المركبات، وأن الطقس والعمر والعمر والعمر والعرب الزار . والتوار . وأنتاج من والعسبة للما والعس والعمر والعوامل المسببة لموادث على الطريق هو عدم وجود مسافة بين المركبات، وأن الطقس والعمر وأظهرت النتائج أن أحد أهم العوامل المسببة للحوادث على الطريق هو عدم وجود مسافة بين المركبات، وأن الطقس والعمر والتهان دورًا مهمًا. حقق النموذج دقة ٤٤,٢٨٪ وانحراف معياري +/– ١٠,٩٠٪. كانت الأدوات المستخدمة في تحليل البيانات والتعلي والتنبؤ هي Python وعامل المسببة وهي منصة برمجية متكاملة لعلوم البيانات توفر بيئة متكاملة للتعلم يوالتعلم الآلي والتنبؤ هي والتبئي.

الكلمات المفتاحية: تحليل البيانات، الخوارزمية، حوادث المركبات، التعلم الآلي، التنبؤ