# Prediction of Roadway Crashes Using Logistic Regression in SAS

Srinivasan Suresh<sup>#1</sup>

#1 Student, Doctorate in Computer Science, Aspen University, Denver, Colorado, USA

#### Abstract

Roadway crashes occur instantly with less time to respond. Predicting these crashes or identifying the major factors affecting these crashes can help to reduce these from occurring. As machine learning techniques help make these predictions and identify the impact factors, they can be applied to the roadway crash data set. The data set is obtained for the State of Virginia from the Department of Transportation. The logistic regression method was applied by grouping the dataset into fatal and non-fatal crashes. The model was built in SAS studio software and had an accuracy of 76%. The major factors were identified as Road, not lighted, Ramps, and Intersections on Divided roadways.

**Keywords** — Fatal roadway crashes, Machine Learning, Logistic Regression, State of Virginia

## I. INTRODUCTION

A road crash is one of the significant reasons for loss of life. A crash could occur due to various reasons and can be preventable in many cases. Safe and defensive driving is always advisable. As per the statistics obtained from the Department of Motor Vehicles in Virginia, almost 8.5 million vehicles are registered, and more than 120 thousand crashes occur every year. There is a very high volume of traffic movement in the state, and the number of fatal crashes varies from 600 to 1000 every year. There could be many factors involved in a fatal crash, and almost 35% is alcohol-related. Immense measures are taken to reduce these fatal crashes. In a few cases, the Road characteristics can also contribute to these types of crashes to some extent. Using a machine learning technique named logistic regression, the major causes of a fatal road crash can be analyzed. Also, a model can be built to identify possible fatal crashes. The logistic regression method is suitable for analyzing binary outcomes. SAS studio software comes up with different options to conduct machine learning studies. This study is performed in the SAS studio by consuming the roadway crash data from the Virginia Department of Transportation.

## **II. LITERATURE SURVEY**

Roadway crashes were predicted through various traditional and machine learning methods. Kononen, Flannagan & Wang (2011) have predicted serious injuries caused by motor vehicle crashes using the logistic regression method. The model was developed using variables available post-crash, and a different screening algorithm was developed to model injuries for each mode of the crash. Feng et al. (2014) have developed a multivariable linear regression model to reflect the various factors on traffic accidents and improve the accuracy of predicting accidents in China. Dong et al. (2019) have performed studies to understand the risk factors for road crashes and their impact on the crashes. A two-step method involving the support vector regression model and state-space models were used.

## **III. DATA PREPARATION**

## A. Data Source

The Department of Motor Vehicles gathers all reported crashes in the State of Virginia, and this data is available to the public from the Virginia Department of Transportation. The crash data required for the analysis is obtained from the State of Virginia (VDOT, 2018) for the period 2013 to 2017. The overall crash dataset is huge, and hence, the research is limited to only Fatal crashes. Fatal crashes are represented with KABCO code as K. The other levels are non-fatal crashes and classified based on the injury's intensity. The data is imported in the form of a rich text file from the website and later converted to a spreadsheet for analysis.

## B. Data preprocessing

The crash dataset contains many observations and details regarding the crash incident. A high-level cleanup is done to remove unwanted reported data. The details apart from the crash cause like Latitude, Longitude, and Route name are removed. These could add bias to the prediction. Most of the predictor variables required for the analysis are categorical ones, as they represent the road characteristics in which the crash has occurred. As the regression analysis is performed at the state level, any references to a specific county, jurisdiction was also removed from it. The attributes were renamed for easier understanding of data and to get a better display in the results. Fatal crashes are represented as 1, and non-fatal crashes are represented as 0 in the dataset. This would help to understand the reasons behind a fatal crash when compared to a nonfatal crash.

# **IV. REGRESSION ANALYSIS**

## A. Logistic Regression

Logistic regression performs classification in the provided data and groups them into a discrete set of classes. It can be a binary or multilinear function. This is used as a machine learning technique to classify the dataset or predict events based on probability. The cost function associated with logistic regression, also called a sigmoid function, is complex. The logistic regression method limits the cost function between 0 and 1.

$$0 \le h_{\theta}(x) \le 1$$

The formula for the sigmoid function is defined as

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

The hypothesis of the logistic regression is expected to give values between 0 and 1

$$Z = \beta_0 + \beta_1 X$$
$$h\Theta(x) = \text{sigmoid}(Z)$$
$$h\theta(X) = \frac{1}{1 + e^{-(\beta_0 + \beta_1 X)}}$$

#### **B.** Data description

SAS Studio is used to analyze the crash data. As the count of predictor variables is more, the regression analysis would be a complex process. Since SAS is a powerful analytical tool, it is chosen for this analysis. A similar analysis can be performed in Python or R as well. The spreadsheet data is loaded into a SAS work file using PROC IMPORT, as shown in figure 1.

	ODE	LO	3	RES	ULTS										
×	Ð-	88	ß	Ð	2	5	e 1	- 1 <sub>k</sub>	8	Lina #	0	×н	20	먨	8
	FILE	NAME	REFFI	LE .	/fol	ders	/myfo]	ders	/Cras	h/Summa	iry_Tr	end_d	ata-	к -	Limited Fields - Cleanup_v1.xlsx';
	PROC	IMPO	RT DA	TAFI	LE=R	EFFI	LE DEM	S=XL!	SX OU	T=WORK.	IMPOR	T:			
	4	GETNA	MES=)	ES:											
	PUN-														
	5														
	5 PRO	CONT	ENTS	DATA	-	W TH	DORT.								

Fig 1: Loading the excel spreadsheet into SAS

PROC FREQ is used to analyze the dataset. From the results, the various categorical values present in the dataset can be observed in figure 2.

	Alcohol Notalcohol								
Alcohol_Notalcoho	Frequenc	y Percer	Cumulativ Trequence	e Cumulativ y Percer					
ALCOHOL	661	3 17.8	1 661	3 17.8					
Not_ALCOHOL	3051	7 82.1	9 3713	0 100.0					
	Date at the band								
Belted_Unbelted	Frequency	Percent	Cumulative Frequency	Cumulative Percent					
BELTED	29047	78.23	29047	78.23					
UNBELTED	8083	21.77	37130	100.00					
Bike_Nonbike	Bil	ke Nonbike Percent	Cumulative Frequency	Cumulative Percent					
Bike_Nonbike BIKE	Bil Frequency 803	e Nonbike Percent 2.16	Cumulative Frequency 803	Cumulative Percent 2.18					
Bike_Nonbike BIKE Not_BIKE	Bik Frequency 803 36327	Percent 2.16 97.84	Cumulative Frequency 803 37130	Cumulative Percent 2.16 100.00					
Bike_Nonbike BIKE Not_BIKE	Bil Frequency 803 36327 CR	Percent 2.18 97.84	Cumulative Frequency 803 37130	Cumulative Percent 2.16 100.00					
Bike_Nonbike BIKE Not_BIKE CRASHYEAR	Bik Frequency 803 36327 CR Frequency	Percent 2.18 97.84 ASHYEAR Percent	Cumulative Frequency 803 37130 Cumulative Frequency	Cumulative Percent 2.18 100.00 Cumulative Percent					
Bike_Nonbike BIKE Not_BIKE CRA SHYEAR 2013	Bil Frequency 803 36327 CF Frequency 7652	Percent 2.18 97.84 A SHYEAR Percent 20.81	Cumulative Frequency 803 37130 Cumulative Frequency 7852	Cumulative Percent 2.16 100.00 Cumulative Percent 20.81					
Bike_Nonbike BIKE Not_BIKE CRASHYEAR 2013 2014	Bii Frequency 803 36327 CF Frequency 7652 6811	Percent 2.16 97.84 Percent 20.61 18.34	Cumulative Frequency 803 37130 Cumulative Frequency 7852 14483	Cumulative Percent 2.18 100.00 Cumulative Percent 20.81 38.95					
Bike_Nonbike BIKE Not_BIKE CRASHYEAR 2013 2014 2015	Bit   Frequency   803   36327   CF   Frequency   7652   6811   7237	Percent 2.16 97.84 ASHYEAR Percent 20.61 18.34 19.49	Cumulative Frequency 803 37130 Cumulative Frequency 7682 14483 21700	Cumulative Percent 2.16 100.00 Cumulative Percent 20.61 38.05 58.44					
Bike_Nonbike BiKE Not_BIKE CRASHYEAR 2013 2014 2015 2016	Bill   Frequency   803   36327   CR   Frequency   7652   6811   7237   7323	Re Nonbike Percent 2.18 97.84 A SHYEAR Percent 20.81 18.34 19.49 19.72	Cumulative Frequency 803 37130 Cumulative Frequency 7852 14483 21700 29023	Cumulative Percent 2.16 100.00 Cumulative Percent 20.61 38.95 58.44 78.17					

Fig 2: A part of PROC FREQ output

The crashes are observed over the year using the SGPLOT, as shown in figure 3. The percentage of fatal crashes is at a peak in 2017.



Fig 3: Percentage of Fatal crashes by year

The weather condition during fatal crashes is observed over the years, as shown in figure 4.



Fig 4: Fatal crashes by year sliced by weather conditions

It seems most of the crashes have occurred during No adverse condition.

## C. Regression model

Usually, principal component analysis is performed to identify the component with the most variance in Numeric variables. Since most of them involved in this analysis are Categorical, the Multiple Correspondence analysis is performed here (NCBI, 2009). The SAS program, PROC CORRESP is used to identify the variables with most variance based on the chi-square value and cumulative percent as in figures 5, 6, and 7



Fig 5: Command to perform the Multiple Correspondence Analysis



Fig 6: Displaying the best contributors to the inertia



Fig 7: Showing the cumulative percent of the Chi-Square value

The input data is split into the Train and Validation data set using the PROC STRATA is shown in figure 8.



NOTE: There were 37130 observations read from the data set WORK.DEVELOP\_SAMPLE. NOTE: The data set WORK.VALID has 12475 observations and 30 variables. NOTE: DATA statement used (Total process time): real time 0.03 seconds cpu time 0.04 seconds

Fig 8: Command to split the dataset into Train and Validate and log statements

As the regression model contains categorical variables, the logistic regression model was chosen. Using PROC LOGISTIC, the model was initially built with the identified Categorical variables from the Multiple Correspondence Analysis. The NULL Hypothesis is these variables do not have an impact on the occurrence of Fatal Crashes.

When calculating the logistic model, the quasi separation of data is observed, as shown in figure 10.



Fig 9: Logs showing quasi complete separation issue

To overcome it, initially, the clustering of data is performed based on business knowledge. But, a similar issue was faced. Hence, the smooth weight of the evidence approach was performed.

CO	DE LOG RESULTS OUTPUT DATA
* (	9- 🔒 😡 🕼 🗈 🕒 🥙 🛩 🐜 🏨 Line# 🥑 🕆 註 🗯 請 👯
66	/* SWOE for Train - FUN */
67	%GLOBAL RH01;
68	PROC SQL NOPRINT;
69	SELECT MEAN(CRASH) INTO :RH01
70	FROM WORK.TRAIN;
71	QUIT;
73	PROC MEANS DATA=WORK.TRAIN SOM NWAY NOPRINT;
74	CLASS FONCTIONAL_CLASS;
75	VAR CRASH; OUTDIT OUTSHOP COUNTS SUM_EVENTS.
70	DIAL
78	non,
79	FILENAME WOE "/folders/myfolders/Crash/swoe branch1.sas":
80	· · · · · · · · · · · · · · · · · · ·
81	DATA NULL ;
82	FILE WOE;
83	SET WORK.COUNTS END=LAST;
84	LOGIT=LOG((EVENTS + &RH01*24)/(_FREQ EVENTS + (1-&RH01)*24));
85	<pre>IF _N_=1 THEN PUT "select (Functional_Class);";</pre>
86	<pre>PUT " when ('" FUNCTIONAL_CLASS + (-1) "') fun_swoe = " LOGIT ";" ;</pre>
87	IF LAST THEN DO;
88	LOGIT = LOG(&RHO1/(1- &RHO1));
89	<pre>PUI = otherwise tun_swoe = " LOGIT -;" / "end;";</pre>
90	END;
91	RUN;

Fig 10: Showing the commands to perform Smooth Weight of Evidence

With this smoothed data, the logistic model was rebuilt, and it does remove the quasi separation issue (Paul D. Allison, 2008).

The PROC LOGISTIC command used to build the model is shown below (Fig 13). It includes the categorical variables qualifying to 95% of the variance from the Multiple Correspondence Analysis. The standardized estimates are calculated along with the model. The forward selection method is used here as there are more categorical variables, and using the stepwise method might result in overfitting.



Fig 11: Command to build a Logistic Regression model

The final predictors and interactions determined by the logistic regression model are shown below in figure 12.



Fig 12: Summary of the Forward Selection approach

Based on the standardized estimates output value, the top 5 categorical values are shown in Table 1.

**TABLE 1: Analysis of Maximum Likelihood Estimates** 

			Wald	Pr	
			Chi-	>	
Para	D		Squar	Chi	Std
meter	F	Estimate	e	Sq	Est
Intercept	1	-0.1162	0	0.99	
Darkness					
-					
Road					
Not					
Lighted	1	2.2489	0.002	0.98	0.515
On/Off					
Ramp	1	6.9227	0.005	0.94	0.370
Grade -					
Straight	1	1.7833	0.003	0.95	0.324
Intersecti					
on -					
Divided				<.0	
Roadway	1	0.5458	55.26	001	0.278

The top factors for fatal crashes are listed below:

- a) Darkness Road Not Lighted
- b) On/Off Ramp on a Divided Roadway
- c) Grade Straight
- d) The intersection at a Divided Roadway

The smooth weight of evidence is done on the validation dataset, too, and the model is scored with it.



Fig 13: Smooth weight of evidence performed on the VALID dataset

The logistic regression model built with the TRAIN dataset is scored with the VALID dataset.



Fig 14: Logistic regression model with evaluating the VALID dataset

A ROC curve is built for the same (Fig 18), and the model turned out to be 76% accurate.



Fig 15: ROC Curve for the VALID dataset

## D. Conclusion

From the regression analysis, it could be inferred that the major root causes for a fatal crash are the darkness in road condition and ramps/intersections on a divided roadway by improving the light conditions on the routes and improvising the ramps/intersections on a divided roadway could reduce the fatal crashes. For further research, one option is to perform the analysis in detail for each county, and improvisation can be made. Another option is to utilize the traffic data, identify the busy routes, and reduce crashes. This could help in saving a lot of travel time.

## REFERENCES

- VDOT. (2018, June 6). Crash Analysis Tool. https://public.tableau.com/profile/publish/Crashtools8\_2/Main #!/publish-confirm
- [2] National Center for Biotechnology Information (NCBI). (2009, Nov 06). Correspondence analysis is a useful tool to uncover the relationships among categorical variables. https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3718710/
- [3] Annapoorani Anantharaman, "A Study of Logistic Regression And Its Optimization Techniques Using Octave" SSRG International Journal of Computer Science and Engineering 6.10 (2019): 23-28.
- [4] Virginia Roads (2018, March 15). Crash Data (Full Details). http://www.virginiaroads.org/datasets/crash-data-full-details
- [5] Analytics Vidhya (2015, July 28). Beginners Guide to Learn Dimension Reduction Techniques
- [6] https://www.analyticsvidhya.com/blog/2015/07/dimensionreduction-methods/
- Paul D. Allison (2008). Convergence failures in Logistic Regression. https://pdfs.semanticscholar.org/4f17/1322108dff719da6aa0d3 54d5f73c9c474de.pdf
- [8] Virginia Department of Motor Vehicles. (, 2019). https://www.dmv.virginia.gov/safety/#crash\_data/index.asp
- [9] Kononen, D. W., Flannagan C. A. & Wang S. C. (2011 Jan). Identification and validation of a logistic regression model for predicting serious injuries associated with motor vehicle crashes. Accid Anal Prev.;43(1):112-22. doi: 10.1016/j.aap.2010.07.018.
- [10] Feng, Z. X., Lu, S. S., Zhang, W. H., & Zhang, N. N. (2014). Combined prediction model of the death toll for road traffic accidents based on independent and dependent variables. *Computational intelligence and neuroscience*, 2014, 103196. https://doi.org/10.1155/2014/103196
- [11] Dong C, Xie K, Sun X, Lyu M & Yue H. (2019) Roadway traffic crash prediction using a state-space model-based support vector regression approach. PLOS ONE 14(4): e0214866. https://doi.org/10.1371/journal.pone.0214866
- [12] T.Maris Murugan, Dr.M.Kandasamy and G.Revathy,(2017). "Accident Prevention System through Intelligent Transport System". SSRG International Journal of Industrial Engineering 4(1), 22-25.
- [13] Annapoorani Anantharaman,(2019). "A Study of Logistic Regression And Its Optimization Techniques Using Octave". SSRG International Journal of Computer Science and Engineering 6(10), 23-28.
- [14] Sangeethu Sharma and Santini (2017). "Accident Avoidance and Safety System for Vehicular Communication". SSRG International Journal of Industrial Engineering 4(2), 1-4.