

Dear Editor and Reviewers,

We would like to thank you for your consideration and helpful comments on our work. Based on your comments, we investigated the data anomalies, and we addressed data coding issues related to the following variables used in the analysis.

1. Single Lane/Multiple lane road: Previously, if this variable was “unknown” for a record, it was re-classified it into “single lane”. This is not appropriate because the majority of records are “multiple lane”. Hence, this variable is now re-classified into “multiple lane” if it is unknown.
2. Critical event that made the crash imminent: Previously this variable was labeled incorrectly. This has now been fixed.
3. Other driver distracted: This variable records whether or not the drivers of motor vehicles other than the bus were distracted. It does not provide information on non-motorists. For non-motorists, this variable is usually coded as “NA” in the datasets. Previously, these “NA”s were assigned to “distracted”. These cases have been re-assigned to “not distracted”. Consequently, the bar plots of this variable reflect information on other motor vehicle drivers of multi-vehicle crashes only.
4. Pickup truck, SUV or van involved/ Light or heavy truck involved: Some additional body type codes were added in the later dataset.

Additionally, we reviewed the data coding more thoroughly to make sure that the analysis remains valid. Note that, after rerunning the clustering, the previous order is not preserved. Below is a point by point response to each reviewer’s comments. We hope that our proposed changes will sufficiently address all your concerns.

## Reviewer 1

Reviewer comment: The authors successfully applied SOM algorithm though the “kohonen” package in R and provided adequate visual and tabular summaries following the cluster analysis. Although authors did not extend the methodology to accommodate variable other than binary factor variable as in Prato and Kaplan (2013), the broader impact aspect of this exercise justifies the publication. Below are some minor suggestions/comments that the authors may wish to consider:

1. Figure 3 may be replaced by a table or mosaic/stacked side-by-side bar plot. Similar to pie chart, radar plots are not ideal for visualization.  
[Response from authors: Figure 3 has been replaced with side-by-side bar plots as suggested.](#)

2. Although the study focuses on bus crashes, multiple plots indicate that “single-lane road” shows significant different distribution for the two time periods. Authors could use some (more) discussion on the this difference as other related factor, such as “changing lane” was mentioned. Perhaps also look into any contributory link from this factor.

Response from authors: This difference no longer exists after the corrections made to this variable as described above.

3. In Section 2.1 “data extraction”, authors mentioned the use of a “procedure that involves some subjective assessment”. Record linkage is very important in analyzing this data set, given its not particularly large sample size. It would be great if the authors can assess the impact on the variance due to these assumptions.

Response from authors: We have added a description of the subjectivity involved in Section 2.1. In particular, in order to bring all the data to the accident level, for multi-vehicle crashes, raw variables at the vehicle-level were summarized as described in Section 2.1. The alternative was to create more variables which would have led to more complicated clustering. We followed the analysis described in Prato and Kaplan (2013).

4. The division at 2009 seems rather unnecessary as suggested by most of the plots. Different data standardization is a legitimate concern.

Response from authors: We wanted to compare how (if at all) the nature of bus crashes have changed over time. This meant a before-after type of analysis.

## Reviewer 2

Reviewer comment: My review is from the perspective of a highway safety engineer. I found some of the results comparing the 2005-2009 dataset with the 2010-2015 very interesting. When comparing Tables 2 and 3, there seems to be very significant differences in many of the characteristics in the two data sets that creates questions/concerns in my mind. Most of those are discussed on page 8. I will go over my questions in order of importance.

1. “Was the other driver/non-motorist distracted?” - you note that from 2005-2009, 100% of the non-motorists were distracted. I find it nearly impossible that all of those crash reports noted that the pedestrian was distracted, and suggest you check the data again. You also note that only 5.5% of the non-motorists were distracted from the 2010-2015 dataset. I also find that extremely unlikely. Anecdotal evidence and some of the discussions about the recent national increases in pedestrian fatalities point towards “distracted walking”, where pedestrians are listening to music though headphones, cars such as hybrids have become much quieter if not silent, and the fact that people walk around with their cell phones in front of them. I suggest you confirm there wasn’t some change in coding.

Response from authors: The variable “Other driver distracted” does not reflect information on non-motorists. The numbers seen previously were a result to mis-classifying “NA”s corresponding to non-motorists as “distracted”. This has been fixed as can be seen from the table.

2. “Critical event that made the crash eminent”- Related to Concern 1, you note “other vehicle encroaching into lane” for Cluster 2 more than doubled from 2005-2009 to 2010-2015 to 30.7%. Since Cluster 2 is single-vehicle crashes involving non-motorists, I’m assuming the “other vehicle” would have to be non-motorists (i.e. pedestrians). You note that non-motorist distraction was only 5.5% for period 2. So unless these pedestrians intentionally encroached into the lane, how do you explain this?

Response from authors: The variable “Critical event that made the crash eminent” was mislabeled previously. This has been corrected. The doubling in non-motorist cluster is seen for category “Non-motorist at fault”. As mentioned previously, the variable “other driver distracted” does not reflect information on non-motorists.

3. “Was there a pick-up truck/van/SUV involved?” - For Cluster 2 on the 2005-2009 data set, you show that 24.8% involved a pick-up truck/van/SUV. Above, it shows that this is only single vehicle involvement. So is this indicating that 24.8% of the buses were actually vans/SUVs? Again, this seems unlikely and could indicate an error in coding.

Response from authors: After data corrections, these proportions are now negligible.

4. “Bus Driver Distracted” - Comparing Cluster 2 for 2005-2009 to 2010-2015, it shows that bus driver distraction dropped in half, from 33.9% to 14.3%. Clusters 1,3, and 4 also show significant drops. Again, with the proliferation of cell phones, I find this result counterintuitive. I’m sure more companies and agencies have policies restricting their use, along with the FMCSA regulation that prohibits commercial vehicle operators have had some effect, but I still would have expected more crashes to note distraction, as its become something that law enforcement agencies are reporting more frequently as contributing to crashes.

Response from authors: This is still the case after we reran our analysis. In 2011, the National Transportation Safety Board called for a “No call, no text, no update behind the wheel” nationwide ban and also started imposing hefty fines on employers failing to adhere to policy. This could have a significant impact on reducing driver distraction levels.

5. “Single lane vs. multiple lanes” - While I would expect more of the bus-involved crashes on multilane roads to be predominant, the percentage on single lane roads being under 4% for all four clusters seems extremely low.

Response from authors: This is still the case after we reran our analysis. Classifying the records with the variable populated as “unknown” might have made a significant difference in proportion in the pre-2009 dataset.

However, we chose to group these records as “multi-lane” because it was the more frequent class. Post 2009, a new variable was introduced by GES. In this new variable, the distribution of “unknown” was significantly different than those in pre-2009 period.

6. You note on page 8, line 32, that the proportion of crashes occurring on in the absence of traffic control is roughly 50%. Considering that the majority of mileage of the transportation system is segments, it doesn’t surprise me that no traffic control was present.

### Reviewer 3

Reviewer comment: This paper is excellent, but there are a few problems in one specific area. I assume these problems can be corrected in a relatively brief timeframe.

1. The GES variables are coded incorrectly several times. Such as Figure 2, showing single lane roads have dropped greatly. This isn’t logical, and should have caused the author to look into the category changes in GES pre and post 2010. There was obviously not a massive drop in the percent of roads that are single lane roads in the US. This is a coding error, since probably the computer code wasn’t changed to represent the new FARS/GES categories for roadway type.

Response from authors: This has been corrected.

2. Distraction coding, and surface conditions, also changed. Distraction was due to GES changes, as it was merged with FARS. Please be aware of all the new FARS/GES files that started in 2010. Before that, it was just accident, vehicle and person. After that many other files were started. That is great for organization, but annoying for a programmer who is forced to merge many files (new in 2010) together to locate all the desired variables. Examples include roadway, precrash....

Response from authors: We double checked all the variables and their derivations and reran the analysis.