

CE C265 & PBHLTH 285
Spring 2025
Assignment 3: Crash Severity Modeling

Submitted by: Chun-hin (Anson) Sit, Jesus Hinojosa
Mar 11, 2025

In this assignment you will develop a crash severity model using an ordered logit regression approach implemented in Python. Your objective is to identify and quantify the factors that influence the severity of crashes by analyzing a variety of explanatory variables such as road characteristics, driver behavior, environmental conditions, and vehicle types. You will use data extracted from TIMS (<https://tims.berkeley.edu>). Details on the data can be found at (<https://tims.berkeley.edu/help/SWITRS.php>).

Option 2

As an alternative to building your own model from scratch, a simple model code will be made available to you (on bcourses). With this code, you can calibrate the model across various counties or cities and then compare and discuss the differences in the estimated coefficients across these geographic regions. This approach will allow you to explore how local factors might influence crash severity and to interpret the variations in the results. Regardless of the option you choose, you are expected to provide a thorough interpretation of the model outcomes, similar to option 1.

Introduction

We are doing Option 2.

We have selected 3 counties to look at and compare the estimated coefficients and odd ratios across them. The 3 chosen counties are:

1. Alameda (Where Berkeley is located in)
2. San Francisco (Where a comprehensive public transit systems is developed)
3. Santa Clara (One of the wealthiest counties in the Bay Area)

For all counties, the range of data is from Jan 1 2022 to Dec 31 2023.

Result Interpretation

Summarizing the independent variables, their coefficients and level of statistical significance into the table below.

	Alameda		San Francisco		Santa Clara	
	Coef	P> z	Coef	P> z	Coef	P> z
Weekend	0.1156	0.002	0.1176	0.035	0.2531	0
Time of Day - Evening	0.2907	0	-0.0462	0.46	0.3037	0
Time of Day - Morning	-0.0997	0.02	-0.1186	0.073	-0.0368	0.415
Time of Day - Night	0.5999	0	0.2751	0.001	0.7751	0
Threshold Parameter 1/2	0.3999	0	0.4881	0	0.4439	0
Threshold Parameter 2/3	0.7549	0	0.625	0	0.767	0
Threshold Parameter 3/4	0.6244	0	0.7387	0	0.5195	0

Note: the coefficients are compared against a baseline scenario, which is the weekday afternoon.

Coefficient

Weekend Effect

- Consistent positive effect across all counties
- Magnitude: Strongest in Santa Clara (0.2531), similar between Alameda and San Francisco (~0.12)
- All statistical significant

Weekend crashes are consistently more severe than weekday crashes, with the effect being about twice as strong in Santa Clara compared to the other counties.

Evening Effect

Notable variation across counties:

- Positive and significant in Alameda (0.2907, $p < 0.001$)
- Positive and significant in Santa Clara (0.3037, $p < 0.001$)
- Negative but not significant in SF (-0.0462, $p = 0.460$)
- Magnitude: Strongest in Santa Clara (0.2531), similar between Alameda and San Francisco (~ 0.12)

Evening crashes are significantly more severe in Alameda and Santa Clara but show no difference from afternoon crashes in San Francisco.

Morning Effect

- Consistent negative direction across all counties:
- Magnitude: Strongest in SF (-0.1186), less pronounced in Alameda (-0.0997), weakest in Santa Clara (-0.0368)
- Significance: Significant in Alameda ($p = 0.020$), marginally significant in San Francisco ($p = 0.073$), not significant in Santa Clara ($p = 0.415$)

Morning crashes tend to be less severe than afternoon crashes, with the effect being more reliable in Alameda and somewhat in San Francisco.

Night Effect

- Strongest positive effect among all variables across all counties
- Magnitude: Highest in Santa Clara (0.7751), followed by Alameda (0.5999), lowest in SF (0.2751)
- Significance: Highly significant in all counties ($p < 0.001$ for Alameda and Santa Clara, $p = 0.001$ for SF)

Night crashes are consistently the most severe across all three counties, with the effect being nearly three times stronger in Santa Clara than in San Francisco.

Overall Pattern

- Night driving is associated with higher crash severity across all counties, with the effect being strongest in Santa Clara.
- San Francisco shows less time-of-day variation in crash severity.
- Santa Clara shows the strongest effects for both weekend and night variables.
- Evening effect is notably absent in San Francisco but strong in the other counties.
- Morning and Evening effects show more geographic variation in significance
- Time of day appears to have a stronger influence than day of week in all counties, with the night period being the most dangerous time for crash severity.

Threshold Parameters

These thresholds represent cut points on an underlying latent continuous variable (y^*) that determines the observed crash severity:

- 1/2: Threshold between severity levels 1 and 2
- 2/3: Threshold between severity levels 2 and 3
- 3/4: Threshold between severity levels 3 and 4

Contrasting the thresholds across the 3 sets of data:

Threshold	Alameda	San Francisco	Santa Clara
1/2	0.4	0.488	0.444
2/3	0.755	0.625	0.767
3/4	0.624	0.739	0.52

The thresholds define the "difficulty" of reaching each severity level:

- Wider gaps between thresholds mean it's harder to move from one severity level to the next
- The specific values determine the baseline probability distribution across severity levels

As observed in the table, the thresholds don't monotonically increase for Alameda and Santa Clara (3/4 is smaller than 2/3). The larger gap between thresholds 1/2 and 2/3 in Alameda and Santa Clara suggests it takes more factors to push a crash from severity 2 to severity 3 than from 1 to 2.

Odd Ratios

The Odd ratios help us understand the different factors like the day of the week and time of day have on the severity of collisions. The ratios found for Alameda, San Francisco, and Santa Clara when the ratio is greater than 1 tells us that the odds of severity increased, and when it is lower than 1 the severity decreased.

As we can see from the Data for all counties looking at the variable "DAY_OF_WEEK_WEEKEND" Collisions are 1.123-1.125 (depending on the county) more likely to happen on the weekends compared to weekdays for each county. This could be due to people taking more trips or engage in risky behaviors on the

weekend like drinking, speeding, or distracted driving due to social events that occur on the weekends.

For “TIME_OF_DAY_E” in Alameda the time in the evening (18:00 - 23:59), collisions are 1.337320 times higher in the afternoon. Compared to San Francisco and Santa Clara counties where their ratio is 0.955 which means it decreases. This could be due to a more trip occurring in Alameda do most jobs are in San Francisco and Santa Clara, and traffic is moving away from these cities and moving toward Alameda.

For “TIME_OF_DAY_M” all counties have lower ratios than 1 of collisions occurring in our of the morning (06:00-11:59) which means they are decreasing. While the morning commute does increase the traffic flow in the morning leading to the belief that more collisions can occur this could indicate that traffic jams cause slower traffic lowering the severity of collisions. Additionally when traffic frees because people are working those on the roads are long-haul drivers or drivers that have strict driving restrictions.

For “TIME_OF_DAY_N” all counties have higher ratios of severe collisions at night (00:00- 5:59) this could be because of lower visibility at night, and a higher risk of ricker behavior due to the decrease in traffic. Alameda is the highest at 1.8 factors that could lead to this possible poor road design which is highlighted when it is darker and people engage in risky behaviors.

Now look at the threshold parameter which shows the odds ratio for the transition from the severity going from level 1 to 2, 2 to 3, and 3 to 4. The higher the ratio the more likely it the to continue to the next level. Level 1 is the least serve, Level 2 is moderately serve, 3 is serve and 4 is the most serve collisions. For levels 1 to 2, the range for all counties is 1.5 to 1.6 showing the high jump for all. For levels 2 to 3 the jump is much higher for all counties meaning more serve collisions are occurring. Levels 3 to 4 have the highest jump for them all meaning that more severe collisions are occurring.

Discussion on the Methodology

Model Fit Statistics

Statistic	Alameda	San Francisco	Santa Clara
Log-Likelihood	-13,058	-5,944	-12,038
AIC/n	1.87	1.82	1.88
BIC/n	1.87	1.83	1.88
Sample Size (n)	14,005	6,545	12,846
Function Value*	0.932	0.908	0.937

*Note: *Function Value = $(-\text{Log-Likelihood})/n$, smaller values indicate better fit*

The differences of the function values are small, suggesting similar predictive power across all three models.

Regarding the AIC/n and BIC/n values, they are also roughly in the same range, with San Francisco having the smallest values, confirming that San Francisco has slightly better model fit.

Threshold Parameters

Thresholds that differ across countries suggest that the same underlying crash factors might be categorized differently in different jurisdictions.

Prediction on Crash Severity

The current set of independent variables failed to discriminate between different severity levels, as shown in the confusion matrices. They predicted no high severity level crashes in the dataset of San Francisco and Santa Clara. While the chosen set of variables are statistically significant, the models fail to translate these relationships into meaningful severity classifications.

Possible Improvements

- Incorporate other variables, as the current model relies solely on temporal variables (time of day and weekday/weekend), and crash severity is influenced by many factors other than that.

Appendix (Jupyter Notebook)