



Data Science Seminar - MSAI 339

Checkpoint 5

December 2nd, 2021

Professor

Jennie Rogers

Students

Anery Patel

Chukwueloka Obi

Ana Cheyre

"We want to study the cause and effect of an allegation being declared unsustained vs sustained. For this, we want to find what are the most common contents among unsustained allegations, and how they differ from the content of a declared sustained as this will be the entire method for our analysis." By finding what are the most common contents among the unsustained and finding how they differ from the content of a declared sustained, we can see if the degree of the allegation is the defining factor between a sustained and unsustained allegation in our NLP Model.

First, we did a Wordcloud for sustained and unsustained allegations, to have a first visual approach about the difference between both contents. We can identify the main differences that attract attention, for sustained allegations we see clearly person names: 'james', 'joseph', 'john'; instead for unsustained allegations we see words with negative connotations, like: 'falsely arrested', 'without justification', 'failed'.



Second, to have a better perspective we identified the top 20 more frequent words for sustained and unsustained allegations, which can be seen in Figures 1 and 2 respectively. For sustained, we can see again the names mentioned before 'james', 'halper', 'joseph', 'webb', so it can be inferred that this type of allegations are mostly related to certain officers/detectives that declare the allegations as sustained. For unsustained, we see 'failed', 'without', 'falsely', 'refused' which gives us an idea of a scenario with unjustified actions.

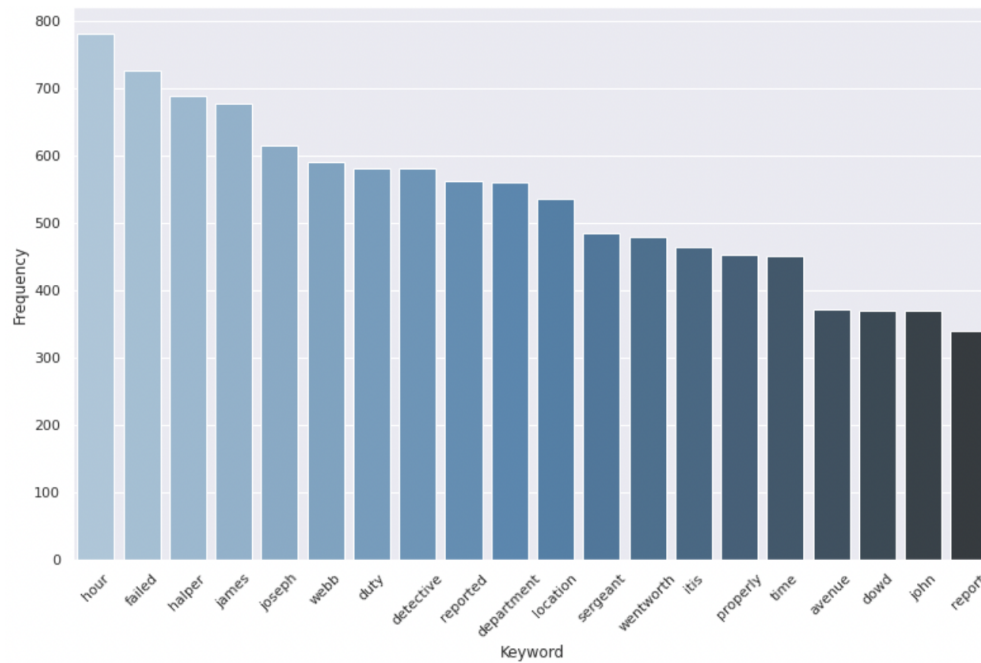


Figure 2: Word frequency for sustained allegations

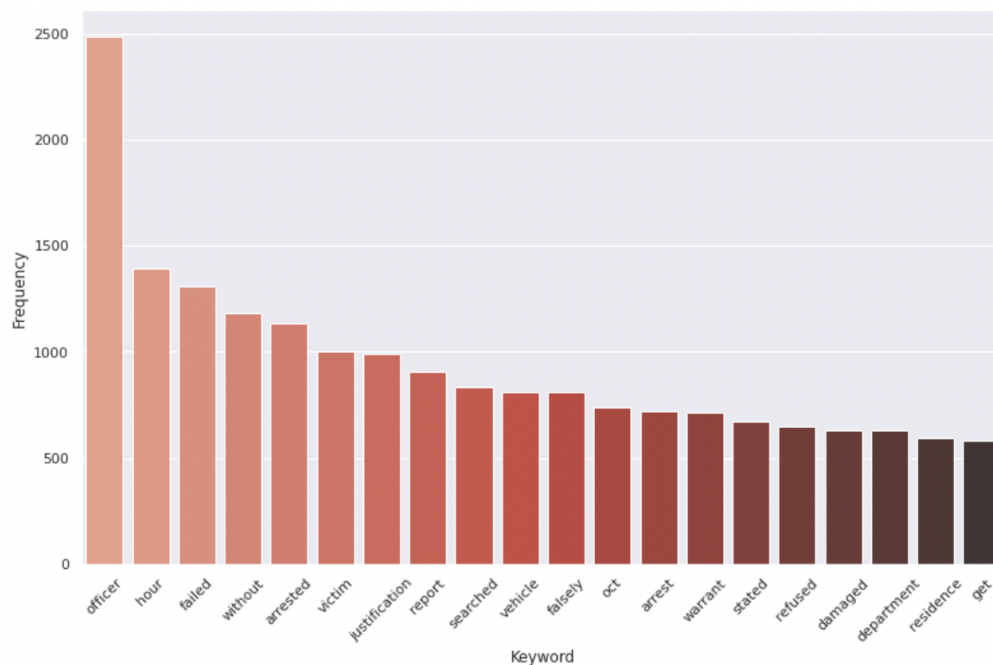


Figure 3: Word frequency for unsustained allegations

Then, to have more clarity of the content, we counted the frequency of bi-grams, referring to words that appear together in an allegation, which can be seen in Figures 4 and 5. For sustained allegations, we can notice that 'detective james halper' is a center key of this type of allegations, also 'sergeant john dowd' and 'joseph webb'. Based on information from the database, we know that only 6.5% of all allegations are sustained, realizing that it corresponds to a very limited group of officers and detectives who declare it this way.

For unsustained allegations, which corresponds to 93.5% of the total allegations, has content more focused on arrests of victims without clear justification by the officers. We can see clearly that the most frequent bi-grams are 'without justification' and 'falsely arrested'.

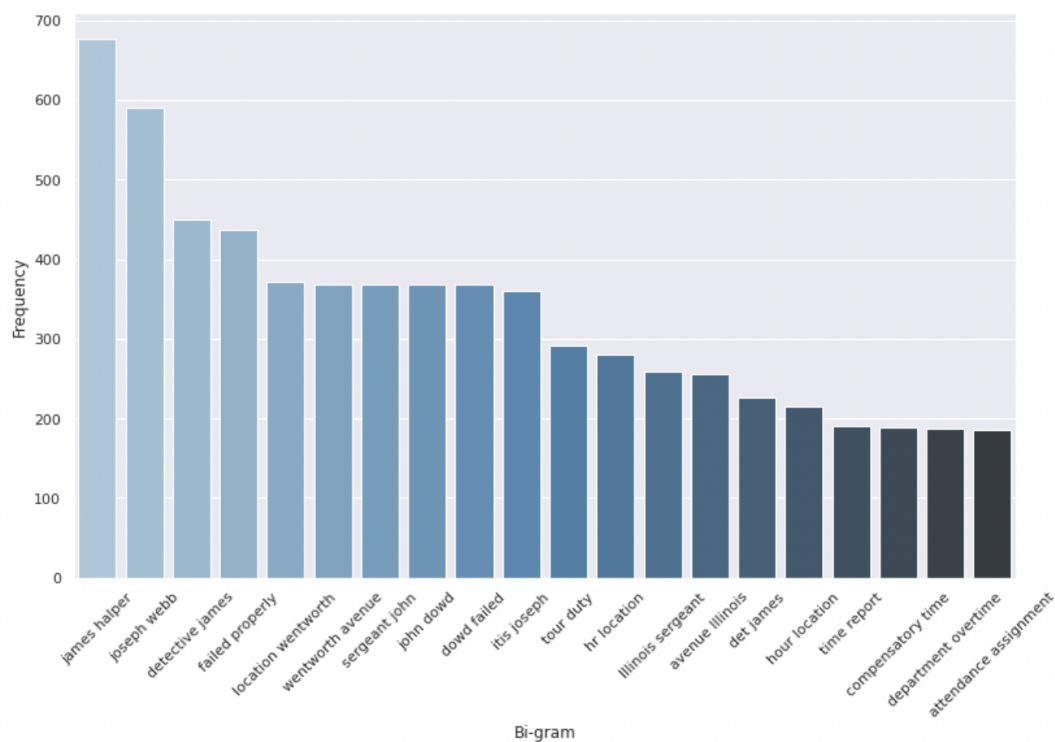


Figure 4: Bi-gram of words frequency for sustained allegations

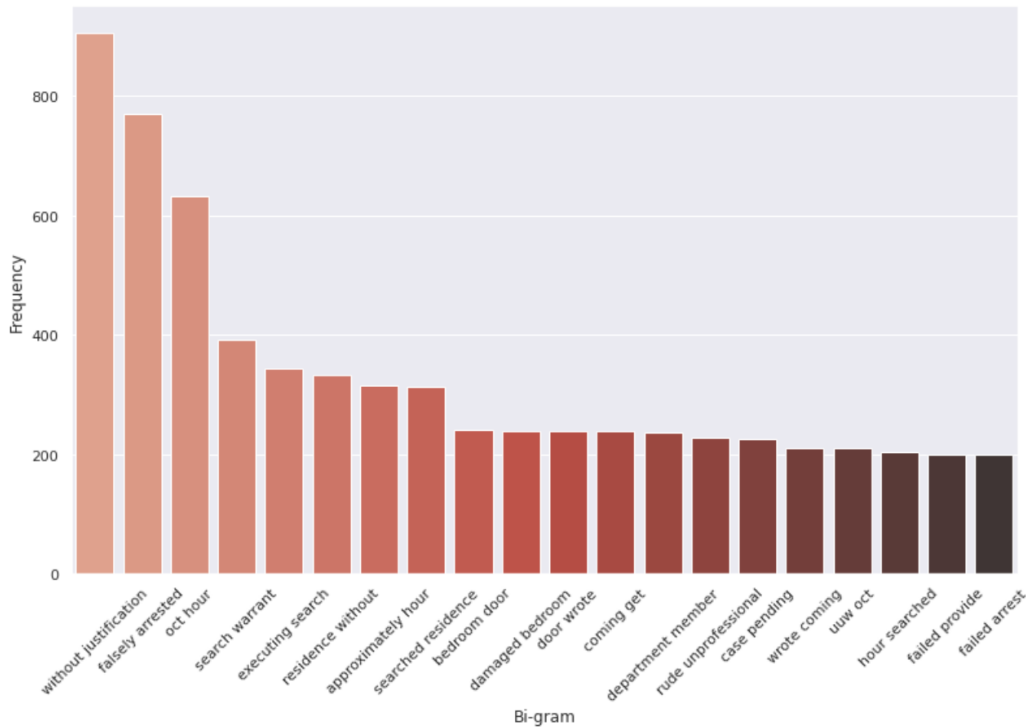


Figure 5: Bi-gram of words frequency for unsustained allegations

Another technique used was TF-IDF, to remove ubiquitous words and identify the most distinctive ones from each type of allegation, applied for all the narratives. The bar charts with the results can be seen in Figure 6 and 7 below.

For sustained allegations we can see that the words 'property', 'bag' and 'department' refer to the evidence of the case kept in a bag, which could be money, driver's license, drugs, etc. We would need to go deeper to find which are the most common things kept inside these property bags.

For unsustained allegations we can see that the words 'tax', 'seven year' and 'social security'. 'Tax' could be the lemmatization of 'taxi' and 'taxes', where the taxis broke traffic laws and the second are related to tax fraud. 'Seven year' may refer to 7-year-old children or to people who have stayed in a place for 7 years, both types of allegation do not seem to be as relevant when looking for them in the databases. Regarding 'social security', just refers to the social security number that is suddenly requested to the victims.

The TF-IDF analysis does not seem to provide such fundamental information when it comes to understanding the content of sustained and unsustained allegations, since it provides certain content exclusive of each type of allegation, but that does not indicate that it is something highly frequent in a transversal way to each of them.

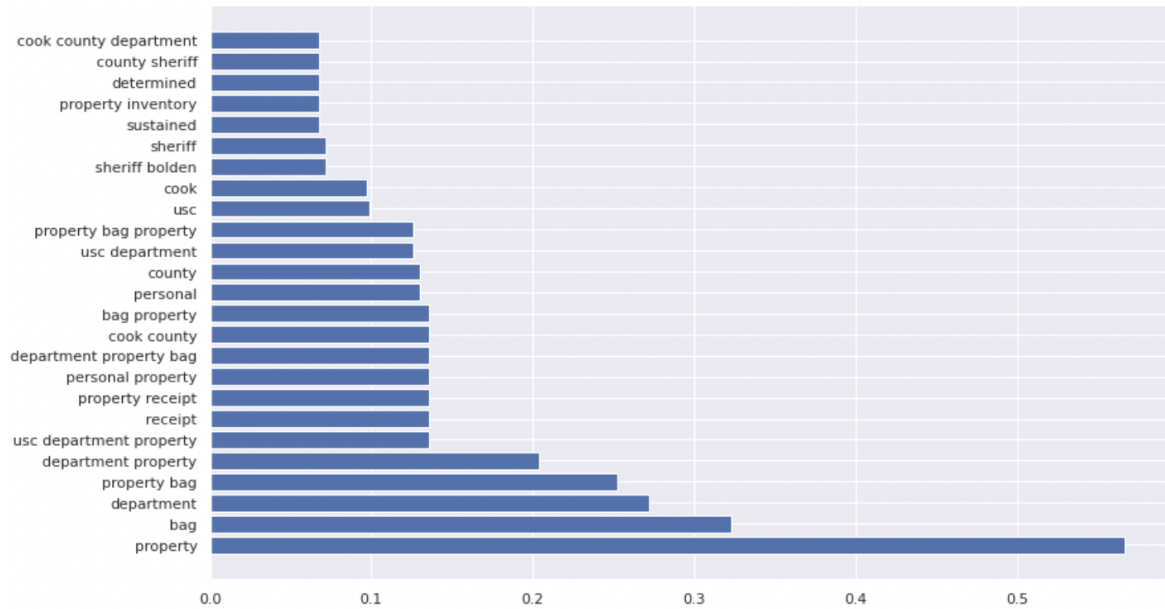


Figure 6: Bar chart of TF-IDF score for sustained allegations

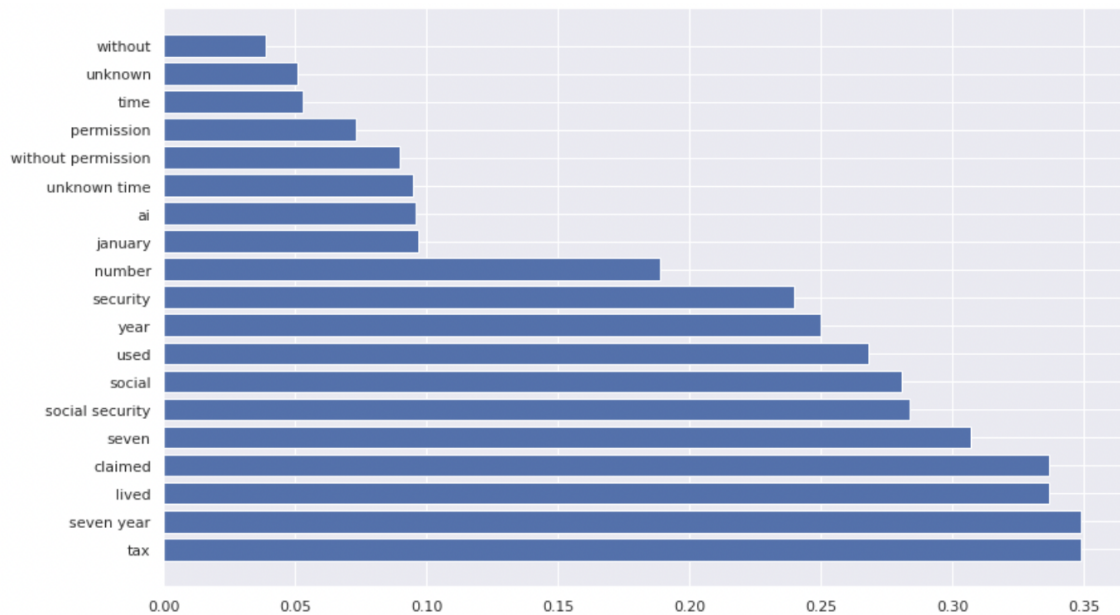


Figure 7: Bar chart of TF-IDF score for unsustained allegations

Finally, to have a deeper analysis of the content we used BERT, a deep neural network model to get the main keywords, which can be seen in the tables below. First, analyzing the sustained keywords, we see that it makes main reference to real 'evidence', such as 'marijuana', 'heroin', 'narcotics'. They are substances that an officer could not hide and are strictly prohibited in their behavior. This could be why officers may be compelled to declare these actions as sustained.

On the other hand, in the unsustained keywords, we see the word 'homicide', a felony that can be considered very severe, but in the eyes of an officer could be excused as 'self defense'. We

also see multiple references to 'rape', which is something difficult to prove and they may decide not to believe the victim. These reasons may be sufficient for the allegation to be declared as unsustained.

Sustained Keywords
marijuana
heroin
narcotics
pornographic
arrestee
mislaide
robbery
burglary
criminal
felony

Unsustained Keywords
homicide
rapist
rape
raped
urinating
burglar
burglary
burglarized
burglarizing
robbed

Conclusion

After the analysis made based on unstructured data of allegations descriptions using natural language processing, we can conclude the following about the content of the sustained and unsustained allegations:

1. **Sustained allegations:** a very characteristic feature is the presence of certain detectives and officers frequently mentioned in these allegations, having a marked bias to declare the allegations as sustained. Also, these types of allegations are mainly related to drugs and other verifiable things, since they have clear evidence and inadmissible behavior.
2. **Unsustained allegations:** a cross-cutting theme is that the accusations described in the allegations were "without justification", cataloging the allegation as unsustained. This makes sense with the main content topics referring to 'homicide', considered a justified action to do by an officer, and 'rape' something very difficult to prove. In both cases, it is understood as an action not justified to accuse an officer. This became a general practice applied to almost 93.5% of all allegations.

To summarize, we can conclude that for an allegation to be declared sustained, the presence of certain officers or detectives is necessary, and also to have evidence that proves the bad behavior of an officer, otherwise the allegation will be declared most of the time as unsustained.