

Proposal

Stat 385 Final Project Proposal

Team Member:

- Jiayi Chen (jchen246)
- Vetrie Senthilkumar (vetries2)
- Tae Kyung Lee (tkl2)

Part 1: Introduction

For our project, we chose to analyze crime in Chicago to identify patterns and trends that could potentially be used to reduce crime in one of America's most troubled cities. We obtained our dataset directly from the City of Chicago's publicly available data. We chose this topic because we considered it to be a very relevant issue that needs to be addressed. Chicago's reputation as a violent city has continued growing in recent years, to the point where many of us have been desensitized to acts of brutality such as assault, robbery, and homicide since we see it every day on the news. In fact, Chicago had 762 homicides in 2016, its highest homicide rate in 19 years, according to CNN. During this time span, it seems that the city has been unable to find any feasible, permanent solutions to its crime problems. Thus, due to the severity of the issue and the fact that many of our fellow students come from Chicago or neighborhoods near it, we felt invested in this topic and thought that it could be a rewarding challenge to take on.

We plan to approach this problem by dividing the city into locations based on the severity and prevalence of crime in that specific region. We will first divide the locations by district, then by street address, and finally by latitude and longitude. This will allow us to use Google Maps to isolate specific areas with high rates of crime. This information would be useful in directing the focus of the Chicago Police Department to the most problematic, crime-riddled areas, helping prevent crime before it occurs. Furthermore, we would like to analyze which types of crime (narcotics, assault, robbery, etc.) occur most frequently in each location since different types of crime could require different strategies by the police department. The dataset also provides us with the setting (residence, street, hotel, etc.) in which each crime occurred. We plan on scrutinizing this information to see if any correlations emerge between certain types of crimes and certain settings. This information might not be immediately useful to us but could prove valuable to forensics teams when they examine the setting of the crime. Finally, we also plan on tracking crime throughout the past couple years to gauge whether the crime situation is improving or worsening in certain locations and in the city of Chicago overall. To accomplish our goals for this project, we must use statistical programming methods to analyze, extract, and visualize data. Since the main purpose of our project is pinpointing locations within Chicago that have the highest crime rates, we need to be able to subset out data based on variables like district, street address, and latitude and longitude. Furthermore, our data has a tabular structure so knowledge about data frames and vectorized operations is essential to analyze and extract data by rows and columns within our data frame. In addition, our dataset contains a lot of text-based information which requires proficiency in regex to process and analyze. Also, we plan on creating graphs using ggplot to show how the crime rates for different types of crime and crime in general have changed over the years. These are just a few examples of how our project will utilize the statistical programming methods taught in this class.

Part2: Related Work

Link to project:

`related_project_page`

Link to dataset: [related_dataset_page](#)

I've included a link to a project we used as inspiration above. This project is a good example of the types of analysis and types of visualizations we would like to include in our project. One aspect of this project we found impressive was the time series analysis and the user interface that allowed users to interact with the time series graphs. These graphs made it easy to observe how crime rates varied over time. We also noticed how this project separated the different types of crimes into categories and provided visualizations on their frequencies. Thus, it was simple to understand which types of crime were more prevalent and problematic and which types were less of an issue. Another feature that caught our attention was the map of the Chicago region that showed what type of crime occurred where and the heatmap that accompanied it. We feel that including a similar geographical overview makes the project feel more applicable and realistic and gives a different perspective on crime's impact on Chicago.

Although we plan to use some ideas from existing projects, we also want our project to be original with its own unique components. We wanted to use the Google Maps API to create a more detailed map of Chicago that would allow users to mouse over regions to obtain the district the region belongs too and the crime info of that region. Furthermore, when doing our analysis, we want to narrow down problematic locations even further. The project above shows which districts have the highest crime rates but we want to be able identify locations on a smaller scope, such as certain streets or maybe even a range of longitude and latitude coordinates. In addition, we also wanted to do further analysis on whether certain types of crime are associated with a particular setting and how the crime rate varies based on the season.

Part3: Method

3.1 General Overview of Code Structure

1. We will use `read.csv` to load our Chicago crime .csv file into R and store it into a data frame.
2. We will initially use methods like `head()`, `tail()`, `summary()`, etc. to gain some insight about our data set before we proceed with our analysis.
3. Next, we will utilize different subset methods to extract the information of certain columns within the data frame. For categorical variables that contain textual info, we will need to use regex and stringr methods like `str_detect()` and `str_extract()`
4. Then, we must pipe in the subsetted data to ggplot2 functions like `geom_hist()` to graph the distribution of numeric variables and `geom_bar()` to graph the distribution of categorical variables
5. Next, we will plot different time series graphs that shows fluctuations in Chicago's crime rate using ggplot's `geom_line()`. To do this, we need to extract dates and different types of crimes from our data set using dplyr methods such as `filter()` and `select()`. We can also use `facet_wrap` to plot similar time series graphs for individual districts.
6. Afterward, we will identify the locations with the highest crime rates by subsetting our data using dplyr based on variables such as street name and latitude and longitude. We can then use ggplot, particularly `geom_bar()`, to display the crime composition for these locations.
7. Now that we've identified problematic areas within Chicago, we will write code that accesses the Google Maps API and pinpoints these locations, creating a heatmap of crime across Chicago. We also want to create another map where individual acts of crime are plotted in different colors, with each color representing a different type of crime.
8. Finally, we will use Shiny to transform our project into a web app with a user interface. We would like to allow users to interact with the maps created in the step above. For example, we plan on using methods such as `selectInput()` and `numericInput()` to allow users to choose the type of crime and year they are interested in analyzing. The maps would then change according to their input.

9. We would also like our web app to include some of the times series graphs we created in step 3. We expect to use the `renderPlot()` function for this part. A feature we are interested in adding to our time series graphs is a scroll bar so users can observe in detail how crime patterns have varied over time by scrolling through our plots. We also want to add buttons to our plots to allow users to select different periods of time (month, year, all time). To achieve this, we plan on using the `sliderInput()` and `actionButton()` functions.

3.2 Packages used in implementation of the project.

3.2.1 To load data

RMySQL - used to read in data from a database

XLConnect, xlsx - used to read and write Microsoft Excel files from R. You can also just export your spreadsheets from Excel as .csv's.

3.2.2 To manipulate/translate data

dplyr - Essential shortcuts for subsetting, summarizing, rearranging, and joining together data sets. dplyr is our go to package for fast data manipulation.

tidyr - Tools for changing the layout of data sets.

stringr - Tools for regular expressions and character strings.

lubridate - Tools that make working with dates and times easier.

dpqr - Summarizing data

3.2.3 To visualize data

ggplot2 - R's package for making graphics

ggvis - Interactive, web based graphics built with the grammar of graphics. maps - Easy to use map polygons for plots.

ggmap - Download street maps straight from Google maps and use them as a background in ggplots. gganimate - Show An Animation Of A Ggplot2 Object

3.2.4 To report results

shiny - A perfect way to explore data and share findings with non-programmers.

3.5 Future Use of Project

This information would be useful in directing the focus of the Chicago Police Department to the most problematic, crime-riddled areas, helping prevent crime before it occurs. Furthermore, we would like to analyze which types of crime (narcotics, assault, robbery, etc.) occur most frequently in each location since different types of crime could require different strategies by the police department. The dataset also provides us with the setting (residence, street, hotel, etc.) in which each crime occurred. We plan on scrutinizing this information to see if any correlations emerge between certain types of crimes and certain settings. This information might not be immediately useful to us but could prove valuable to forensics teams when they examine the setting of the crime. Finally, we also plan on tracking crime throughout the past couple years to gauge whether the crime situation is improving or worsening in certain locations and in the city of Chicago overall.

Part 4: Feasibility

One of the challenges of the project was gathering a dataset that satisfied all of the requirements and captured our interest. Using the time series analysis and other types of visualizations, we are planning to create a project that illustrates how Chicago has been impacted by crime. Our project will use aspects from a previous project, but will implement our own unique ideas and vision. Specifically, our user interface will allow the user to interact with the details and info of criminal activity in Chicago. We feel the most difficult part of this project will be creating the user interface and using new concepts like shiny. We anticipate that this will take the most amount of time, but we are confident that the core idea behind our project can be completed by the end of the semester. After completing the project proposal, the group will spend the next two weeks working on the project demo video. We are planning to meet at least once or twice a week to ensure fluent communication across all members and combining our individual work to form a cohesive project. Thus, we have decided to meet after each class. Each member will have specific roles to complete. One member will focus on data visualizations and initial data analysis. This member will create graphs such as the box-and-whisker plot, line plot, and matrix plot. The other two members will primarily devote their time to developing the user interface using Shiny and the Google Maps API. They will be responsible for implementing features like the heatmap and for adding tools like scroll bars, sliders, and buttons to allow users to interactively explore our data. Since Shiny is a new concept and the Google Maps API is something that hasn't been covered in class, we anticipate this part of the project to be time consuming and potentially frustrating. Therefore, we believe that this is a fair distribution of work across group members that will ensure that our project gets completed by the deadline.

Part 5: References

5.1 list (5+) of papers or items read to write this proposal:

- Azadeh Ansari and Rosa Flores, "Chicago's 762 homicides in 2016 is highest in 19 years", 2017, <https://www.cnn.com/2017/01/01/us/chicago-murders-2016/index.html>
- Garrett Golemund, "Quick list of useful R packages", 2018, <https://support.rstudio.com/hc/en-us/articles/201057987-Quick-list-of-useful-R-packages>
- Vivek Mangipudi, "ANALYSIS OF CRIMES IN CHICAGO 2001 - 2017", 2017, https://rstudio-pubs-static.s3.amazonaws.com/294927_b602318d06b74e4cb2e6be336522e94e.html
- Andrew V. Papachristos, "48 YEARS OF CRIME IN CHICAGO: A Descriptive Analysis of Serious Crime Trends from 1965 to 2013", 2013
- Sai Krishna Vithal Lolla, "Crime Occurrence Analysis in Chicago City", 2013, https://www.jmp.com/about/events/summit2013/resources/Poster25_Lolla_Liu.pdf
- Udeh Tochukwu, "CRIME MINING AND INVESTIGATION USING R", 2015, https://www.researchgate.net/publication/277020625_CRIME_MINING_AND_INVESTIGATION_USING_R

5.2 list all R packages or software referenced:

```
citation(package="XLConnect")
```

```
##
## To cite package 'XLConnect' in publications use:
##
##   Mirai Solutions GmbH (2018). XLConnect: Excel Connector for R. R
##   package version 0.2-15. http://www.mirai-solutions.com
##   https://github.com/miraisolutions/xlconnect
##
## A BibTeX entry for LaTeX users is
##
```

```
## @Manual{,
##   title = {XLConnect: Excel Connector for R},
##   author = {{Mirai Solutions GmbH}},
##   year = {2018},
##   note = {R package version 0.2-15},
##   url = {http://www.mirai-solutions.com
## https://github.com/miraisolutions/xlconnect},
## }
```

```
citation(package="dplyr")
```

```
##
## To cite package 'dplyr' in publications use:
##
##   Hadley Wickham, Romain Francois, Lionel Henry and Kirill Müller
##   (2017). dplyr: A Grammar of Data Manipulation. R package version
##   0.7.4. https://CRAN.R-project.org/package=dplyr
##
## A BibTeX entry for LaTeX users is
##
## @Manual{,
##   title = {dplyr: A Grammar of Data Manipulation},
##   author = {Hadley Wickham and Romain Francois and Lionel Henry and Kirill Müller},
##   year = {2017},
##   note = {R package version 0.7.4},
##   url = {https://CRAN.R-project.org/package=dplyr},
## }
```

```
citation(package="tidyr")
```

```
##
## To cite package 'tidyr' in publications use:
##
##   Hadley Wickham and Lionel Henry (2018). tidyr: Easily Tidy Data
##   with 'spread()' and 'gather()' Functions. R package version
##   0.8.0. https://CRAN.R-project.org/package=tidyr
##
## A BibTeX entry for LaTeX users is
##
## @Manual{,
##   title = {tidyr: Easily Tidy Data with 'spread()' and 'gather()' Functions},
##   author = {Hadley Wickham and Lionel Henry},
##   year = {2018},
##   note = {R package version 0.8.0},
##   url = {https://CRAN.R-project.org/package=tidyr},
## }
```

```
citation(package="lubridate")
```

```
##
## To cite lubridate in publications use:
##
##   Garrett Golemund, Hadley Wickham (2011). Dates and Times Made
##   Easy with lubridate. Journal of Statistical Software, 40(3),
##   1-25. URL http://www.jstatsoft.org/v40/i03/.
##
```

```
## A BibTeX entry for LaTeX users is
```

```
##
```

```
## @Article{,  
##   title = {Dates and Times Made Easy with {lubridate}},  
##   author = {Garrett Grolmund and Hadley Wickham},  
##   journal = {Journal of Statistical Software},  
##   year = {2011},  
##   volume = {40},  
##   number = {3},  
##   pages = {1--25},  
##   url = {http://www.jstatsoft.org/v40/i03/},  
## }
```

```
citation(package="dplyr")
```

```
##
```

```
## To cite package 'dplyr' in publications use:
```

```
##
```

```
## Hadley Wickham, Romain Francois, Lionel Henry and Kirill Müller  
## (2017). dplyr: A Grammar of Data Manipulation. R package version  
## 0.7.4. https://CRAN.R-project.org/package=dplyr
```

```
##
```

```
## A BibTeX entry for LaTeX users is
```

```
##
```

```
## @Manual{,  
##   title = {dplyr: A Grammar of Data Manipulation},  
##   author = {Hadley Wickham and Romain Francois and Lionel Henry and Kirill Müller},  
##   year = {2017},  
##   note = {R package version 0.7.4},  
##   url = {https://CRAN.R-project.org/package=dplyr},  
## }
```

```
citation(package="ggplot2")
```

```
##
```

```
## To cite ggplot2 in publications, please use:
```

```
##
```

```
## H. Wickham. ggplot2: Elegant Graphics for Data Analysis.  
## Springer-Verlag New York, 2009.
```

```
##
```

```
## A BibTeX entry for LaTeX users is
```

```
##
```

```
## @Book{,  
##   author = {Hadley Wickham},  
##   title = {ggplot2: Elegant Graphics for Data Analysis},  
##   publisher = {Springer-Verlag New York},  
##   year = {2009},  
##   isbn = {978-0-387-98140-6},  
##   url = {http://ggplot2.org},  
## }
```

```
citation(package="maps")
```

```
##
```

```
## To cite package 'maps' in publications use:
```

```
##
```

```
## Original S code by Richard A. Becker, Allan R. Wilks. R version
## by Ray Brownrigg. Enhancements by Thomas P Minka and Alex
## Deckmyn. (2018). maps: Draw Geographical Maps. R package version
## 3.3.0.
##
## A BibTeX entry for LaTeX users is
##
## @Manual{,
##   title = {maps: Draw Geographical Maps},
##   author = {Original S code by Richard A. Becker and Allan R. Wilks. R version by Ray Brownrigg. E
##   year = {2018},
##   note = {R package version 3.3.0},
## }
##
## ATTENTION: This citation information has been auto-generated from
## the package DESCRIPTION file and may need manual editing, see
## 'help("citation")'.
```

```
citation(package="ggmap")
```

```
##
## To cite ggmap in publications, please use:
##
## D. Kahle and H. Wickham. ggmap: Spatial Visualization with
## ggplot2. The R Journal, 5(1), 144-161. URL
## http://journal.r-project.org/archive/2013-1/kahle-wickham.pdf
##
## A BibTeX entry for LaTeX users is
##
## @Article{,
##   author = {David Kahle and Hadley Wickham},
##   title = {ggmap: Spatial Visualization with ggplot2},
##   journal = {The R Journal},
##   year = {2013},
##   volume = {5},
##   number = {1},
##   pages = {144--161},
##   url = {http://journal.r-project.org/archive/2013-1/kahle-wickham.pdf},
## }
```

```
citation(package="shiny")
```

```
##
## To cite package 'shiny' in publications use:
##
## Winston Chang, Joe Cheng, JJ Allaire, Yihui Xie and Jonathan
## McPherson (2017). shiny: Web Application Framework for R. R
## package version 1.0.5. https://CRAN.R-project.org/package=shiny
##
## A BibTeX entry for LaTeX users is
##
## @Manual{,
##   title = {shiny: Web Application Framework for R},
##   author = {Winston Chang and Joe Cheng and JJ Allaire and Yihui Xie and Jonathan McPherson},
##   year = {2017},
```

```
##      note = {R package version 1.0.5},  
##      url  = {https://CRAN.R-project.org/package=shiny},  
##    }
```