

The Correlation Between Weather and Crime

Introduction

I have used the National Oceanic and Atmospheric Association(NOAA) weather dataset and Chicago Crime Dataset from the City of Chicago. I chose these because they both have daily entries, and I was able to merge the datasets by day. These will allow me to find a day by day relationship between weather and crime, as I felt that monthly was too far apart to find meaningful relationships. I also liked the crime dataset because it labels the type of each crime. This way I am able to investigate the relationship between weather and different types of crimes.

Getting a Weather Dataset

I used the Weather Dataset from the National Oceanic and Atmospheric Association(NOAA). I was able to find a dataset for daily weather reports in Chicago from 1949 to 2009, which I filtered down to only the 2008 reports. I only kept values such as precipitation, snow, and temperatur max and mins, along with the key values year, month, and day.

```
weather <- read.dly("ftp://ftp.ncdc.noaa.gov/pub/data/ghcn/daily/all/USW00094846.dly")

# keep only 2008 entries
weather_2008 <- weather %>%
  filter(YEAR == 2008)

# Select variables of interest
weather_2008 <- weather_2008 %>%
  select(YEAR, MONTH, DAY, PRCP.VALUE, SNOW.VALUE,SNWD.VALUE, TMAX.VALUE, TMIN.VALUE)

# check out the dataset
save(weather_2008, file='weather.RData')
rm(weather)
```

Gathering the Chicago Crime Dataset

I used the official dataset from the City of Chicago, recording all reported crimes in 2008. I then split the date into year, month, and day to prepare for merging with the weather dataset. Finally, I only kept the variables year, month, day, as well as the type of crime committed.

```
# Download one year of crime data from the open data portal of city of Chicago
# NOTE: This may take a while depening on the strength of your internet connection
# First I ran read_csv() to find the default col_types() then I updated them to this:
type=cols( `CASE#` = col_character(),
  `DATE OF OCCURENCE` = col_datetime(format="%m/%d/%Y %I:%M:%S %p"),
  BLOCK = col_factor(),
  IUCR = col_factor(),
  `PRIMARY DESCRIPTION` = col_factor(),
  `SECONDARY DESCRIPTION` = col_factor(),
  `LOCATION DESCRIPTION` = col_factor(),
  ARREST = col_factor(),
  DOMESTIC = col_factor(),
  BEAT = col_factor(),
  WARD = col_factor(),
  `FBI CD` = col_factor(),
  `X COORDINATE` = col_double(),
  `Y COORDINATE` = col_double(),
  LATITUDE = col_double(),
  LONGITUDE = col_double(),
  LOCATION = col_character()
)

# Read in data
crime_raw <- read_csv('Crimes_-_2008.csv', na='', col_types = type)

# Fix column names
names(crime_raw)<-str_to_lower(names(crime_raw)) %>%
  str_replace_all(" ","_") %>%
  str_replace_all("-","_") %>%
  str_replace_all("#","_num")

crime_2008 <- crime_raw %>%
  separate(date, c('MONTH', 'DAY', 'YEAR'), sep = c('/'))

crime_2008 <- crime_2008 %>%
  separate(YEAR, c('YEAR', 'TIME'), sep = c(' '))

crime_2008$YEAR <- as.numeric(crime_2008$YEAR)
crime_2008$MONTH <- as.numeric(crime_2008$MONTH)
crime_2008$DAY <- as.numeric(crime_2008$DAY)

crime_2008 <- crime_2008 %>% select(YEAR, MONTH, DAY, primary_type)

crime_2008 <- crime_2008 %>% arrange(YEAR)

save(crime_2008, file='crime.RData')

rm(crime_raw)
```

Merging the Crime and Weather Datasets

Next, I merged the two datasets by key ID's Year, Month, Day. First, I inspected the two datasets. The crime dataset had 427142 observations, while weather_2008 had 366 (one for each day as 2008 was a leap year). Performing an antijoin on the two resulted in a table with 0 observations, meaning that there were no mismatches between the datasets!. I performed a left join on the crime dataset with the weather dataset, resulting in a table with the same number of observations as the crime datasets since no rows had to be dropped. Finally, I reverted to a singular date column to make plotting a bit easier.

```
## inspect weather dataset

load('weather.RData')
load('crime.RData')
head(weather_2008)%>%
  kbl() %>%
  kable_styling(bootstrap_options = c("striped", "hover", "condensed"))
```

kable_styling(bootstrap_options = c("striped", "hover", "condensed"))			
YEAR	MONTH	DAY	primary_type
2008	1	1	CRIM SEXUAL ASSAULT
2008	10	24	DECEPTIVE PRACTICE
2008	7	24	SEX OFFENSE
2008	1	1	OFFENSE INVOLVING CHILDREN
2008	2	26	MOTOR VEHICLE THEFT
2008	6	28	THEFT