

3-12-2018

# Automated Outlier Detection in Crime Data Using Programming

Aidan Connolly

*University of Nebraska-Lincoln*

Follow this and additional works at: <https://digitalcommons.unl.edu/honorsthesis>



Part of the [Journalism Studies Commons](#)

---

Connolly, Aidan, "Automated Outlier Detection in Crime Data Using Programming" (2018). *Honors Theses, University of Nebraska-Lincoln*. 38.

<https://digitalcommons.unl.edu/honorsthesis/38>

This Article is brought to you for free and open access by the Honors Program at DigitalCommons@University of Nebraska - Lincoln. It has been accepted for inclusion in Honors Theses, University of Nebraska-Lincoln by an authorized administrator of DigitalCommons@University of Nebraska - Lincoln.

AUTOMATED OUTLIER DETECTION IN CRIME DATA USING PROGRAMMING

An Undergraduate Honors Thesis  
Submitted in Partial Fulfillment of  
University Honors Program Requirements  
University of Nebraska-Lincoln

by  
Aidan Connolly, BJ  
Journalism  
College of Journalism and Mass Communications

March 12, 2018

Faculty Mentors:  
Matt Waite, BJ, Journalism  
Sue Bullard, MSA, Journalism

# Abstract

After the University of Nebraska-Lincoln Police department began publishing their Daily Crime and Fire Log online, journalists and other members of the public have been able to view updates almost instantly. They can see what incidents have been reported so far for that day, and they can view any day back to 2005. Using an advanced search, they can also filter the data by date range, location or crime type.

However, there is no way to analyze the data. There's no way to see how crime reports have evolved over time. Other people have developed programs to look at past trends and outliers to see how things have changed, but there was no way to know when new outliers were happening. The goal of this program is to fill that gap.

This program uses Python to calculate the average number of reports per month for each crime type. Then, as the reports come in each month, it checks to see if any crime type has an abnormally high number of crimes reported. At the end of the month, it checks to see if an unusually low number of crimes were reported for a crime type.

If an abnormality is found, a message is created and sent to a messaging platform common to newsrooms called Slack. This allows journalists to be notified of the abnormality. From there, they're able to look into the reports to determine if it is worth a story.

## Key words

journalism, crime, programming, statistics, automation

# Background

Automation has been a part of journalism for a while, and it has been used in more ways than just alerting journalists. In March 2014, Ken Schwenke developed a program for the Los Angeles Times that automatically wrote a story every time an earthquake is registered above a certain threshold. Using information from the U.S. Geological Survey, [“Quakebot”](#) fills in a pre-written template and puts together a story in seconds. All Schwenke has to do is quickly glance over the story before publication, allowing the Los Angeles Times to publish those stories much faster than other publications.

Automation can help with quantity, too. In July 2014, the Associated Press [began using automated journalism](#) for its quarterly earnings reports. Previously, staff would comb through earnings reports to show trends and numbers. In an effort to reevaluate its use of resources, AP designed a program to help with the task.

Now, an algorithm looks at the latest earnings report and uses information from previous reports to automatically generate stories. They went from 300 reports each quarter to 4,000. This allows the journalists to spend more time focusing on the unusual trends and exclusive stories that need a human to write them. The AP estimates the automation of earnings reports has [freed up 20 percent](#) of the time spent producing reports each quarter.

The Washington Post has also used automated journalism. With [Heliograf](#), the Post is able to automatically write stories about [local high school sports](#). The program gets its data from information about when and how each team scored, along with player statistics and weekly rankings.

For more examples and information, check out the Tow Center for Digital Journalism’s [Guide to Automated Journalism](#).

While the program I built does not yet automatically write stories, it’s in the works, and this program accomplishes the first steps to make that possible.

# Process

## Downloading the data

---

To start, I needed to download as much UNLPD data I could get. Using the advanced search function of the Daily Crime and Fire Log, I was able to download data one year at a time, all the way back to 2005. Then, using `csvkit`, I was able to trim and stack the files.

```
1 | in2csv 2005.xls > 2005.csv --skip-lines 8
2 | in2csv 2006.xls > 2006.csv --skip-lines 8
3 | ...
4 | in2csv 2017.xls > 2017.csv --skip-lines 8
5 |
6 | csvstack 2005.csv 2006.csv > 0506.csv
7 | csvstack 0506.csv 2007.csv > 0507.csv
8 | ...
9 | csvstack 0516.csv 2017.csv > all_years.csv
```

Bash

`in2csv` converts the downloaded Excel files into csv files, and `--skip-lines 8` removes the unnecessary header at each of the files. `csvstack` joins two files by putting one file's rows at the end of another file.

## Importing the historical data

---

The next step was to read the final csv into the data analysis library `pandas`. I chose `pandas` for this project because of its ability to easily manipulate data and do simple statistics.

```
1 | all_years = pd.read_csv('original_data/all_years.csv')
```

Python

Side note: `pandas` is often imported into Python scripts using `import pandas as pd`, which renames the library when it's imported. I have followed this trend.

## Checking the data

---

Once the data is loaded, I wanted to check the validity of it. All of the records should have a case number, a reported date, a location, an amount stolen and an amount damaged. To ensure each record has these values, we can count how many values are present in the data and compare it to the number of rows.

```
1 | assert df['Case #'].count() == len(df) != 0
2 | assert df['Reported'].count() == len(df) != 0
3 | assert df['Location'].count() == len(df) != 0
4 | assert df['Stolen'].count() == len(df) != 0
5 | assert df['Damaged'].count() == len(df) != 0
```

Python

These `assert` statements check the count of values against the number of rows, found using the `len()` function. Checking to make sure the values do not equal 0 helps to double check the data downloaded and imported correctly.

## Cleaning the data

---

This program doesn't work with the stolen and damaged amounts specifically, but cleaned them anyway. To be able to treat them as numbers, I removed the extra punctuation and converted the data type of the column from string to float. I also converted the Reported column to datetime objects for time-based analysis.

```
1 | df['Stolen'] = df['Stolen'].str.replace(',','')
2 | df['Stolen'] = df['Stolen'].str.replace('$','')
3 | df['Stolen'] = df['Stolen'].astype(float)
4 | df['Damaged'] = df['Damaged'].str.replace(',','')
5 | df['Damaged'] = df['Damaged'].str.replace('$','')
6 | df['Damaged'] = df['Damaged'].astype(float)
7 |
8 | df['Reported'] = pd.to_datetime(df['Reported'])
```

Python

The `replace()` method removes the characters by replacing them with nothing. Then, the `astype()` method converts the Stolen and Damaged columns to float values, and the `to_datetime()` function converts the Reported column to datetime objects.

## Creating a month column

To filter the data by month, I needed to create a column with the month each crime was reported. I filtered by month to allow trends to appear that wouldn't be obvious to the standard journalist.

```
1 df['Month'] = df['Reported'].dt.to_period('M')
2 df2 = df.set_index(['Reported'])
```

Python

The `to_period()` method returns the time period a datetime is in, and in this case, it's finding the month. That info is saved to a new column named 'Month.' Then, the index of the dataframe is switched to the Reported column to allow slicing by month later.

## Counting the crimes

Next, I needed to count how many crimes were reported for each crime type. I separated the data by crime type and separated it into subsets by month. I counted how many rows were in those subsets and saved it to a dictionary. Then, all those dictionaries were added to a list, which was converted into a `pandas` dataframe.

```
1 #This holds the dictionary for each crime
2 months_count = []
3 #For each crime present in the dataframe
4 for crime in df['Incident Code'].unique():
5     crime_dict = {}
6     crime_dict['Crime'] = crime
7     #For each month in the dataframe
8     for month in df['Month'].unique():
9         #Slice the dataframe for one month's data
10        month_subset = df[str(month)]
11        #Filter the subset for instances of the crime
12        crime_subset = month_subset[month_subset['Incident Code'] == crime]
13        #Save the count with the month
14        crime_dict[str(month)] = len(crime_subset)
15        #Append the dictionary to the months_count list
16        months_count.append(crime_dict)
17 #Convert the list into another dataframe
18 months_count_df = pd.DataFrame(months_count)
```

Python

This is a sample of the result.

Crime	2005-01	2005-02	2005-03	...
LOST OR STOLEN ITEM	10	5	8	...
FRAUD - CREDIT CARDS/ATM/BANK CARD	2	4	2	...
ACCIDENTS - P.D. REPORTABLE	4	7	4	...
...	...	...	...	...

## Calculating the statistics

With the crimes counted, I could calculate the average number of crimes reported per month. I also calculated the standard deviation to help create an upper and lower threshold for alerts.

```

1 #Creates a dataframe with the unique crimes
2 std_df = df.filter(['Crime'])
3 #Adds a column with the mean count for each crime
4 std_df['mean'] = df.mean(axis=1)
5 #Adds a column with the standard deviation for each crime
6 std_df['std'] = df.std(axis=1)
7 #Adds a column with a lower threshold
8 std_df['lower'] = std_df['mean'] - std_df['std']
9 #Adds a column with an upper threshold
10 std_df['upper'] = std_df['mean'] + std_df['std']

```

Python

The `filter()` method created a new dataframe with only the crime types column. The `mean()` method calculates the mean for each row, and `axis=1` ensures it's calculating the mean for each crime type, rather than the mean for each month. That information is saved to a new column. Then, the standard deviation is calculated for each crime type using the `std()` method, and it's saved to a new column.

The lower bound is calculated by subtracting one standard deviation from the mean, and the upper bound is calculated by adding one standard deviation to the mean. I chose to create the thresholds at one standard deviation away, as that seemed to be a good balance between missing outliers and having too many false positives.

## Getting this month's data

Because UNLPD's Daily Crime and Fire Log is a `.asp` page, standard web scraping processes couldn't be used, as the website's URL never changes. Instead, I needed to use `selenium`, which was developed to automate web browser interaction.

```

1 #This is needed to set up selenium
2 #os.path.expanduser allows the use of a '~'
3 path_to_chromedriver = os.path.expanduser('~Downloads/chromedriver')
4 browser = webdriver.Chrome(executable_path=path_to_chromedriver)

```

Python

`selenium` uses a file called `chromedriver` to manipulate an instance of Google Chrome.

`os.path.expanduser` allows other users to simply place the file in their Downloads folder without having to change the file path provided.

```

1 #The URL to the Daily Crime and Fire Log
2 url = "https://scsapps.unl.edu/policereports/MainPage.aspx"
3 #Go to the URL
4 browser.get(url)
5 #Find the advanced search button and click it
6 browser.find_element_by_id('ctl00_ContentPlaceholder1_AdvancedSearchButton').click()
7 #Find the first date field, hit tab and hit '01'.
8 #This sets the date to the first day of the month
9 date_box = browser.find_element_by_id('ctl00_ContentPlaceholder1_DateRange_MonthText1')
10 date_box.send_keys('\t01')
11 #Find the search button and click it
12 browser.find_element_by_id('ctl00_ContentPlaceholder1_SearchButton').click()
13 #Switch to the iframe on the page
14 browser.switch_to.frame(browser.find_element_by_id('ctl00_ContentPlaceholder1_ViewPort'))
15 #Find the export button once the iframe loads and click it
16 export_button = WebDriverWait(browser, 10).until(
17     EC.presence_of_element_located((By.ID, 'ExportButton')))
18 )
19 export_button.click()

```

Python

`selenium` then goes through the pages, filling out the form to get this month's data. The form has automatically moves your cursor when you enter the date. So, to set the beginning date to the first of the

month, it finds the field for the first month, hits tab to move to the day field, and enters '01.'

Then, when it loads the data, it loads it in an `iframe`, which is basically a webpage inside of a webpage. `selenium` has to switch to the `iframe` and then wait for it to load. Once it does, it can click the export button.

Because the data always downloads with the same filename, it's important the file is deleted before the program is run again. Otherwise, the new file will have '(1)' after it, causing the program to continue using the old data.

## Importing the monthly data

---

Since we were already manipulating data in Python, I didn't want the user to have to quit the program to load in the monthly data. So, I used the `subprocess` library to run the `in2csv` command.

```

1 #Runs in2csv on the downloaded file and converts it to UTF-8
2 csv_data = subprocess.check_output([
3     "in2csv",
4     os.path.expanduser("~/Downloads/DailyCrimeLogSummary.xls"),
5 ]stderr=subprocess.DEVNULL,).decode("utf-8")
6 #Creates a file instance for pandas to use on the next line
7 csv_file_instance = StringIO(csv_data)
8 #Reads in the csv to a dataframe, skipping the first eight rows
9 month_df = pd.read_csv(csv_file_instance, skiprows=8)

```

Python

`subprocess.checkoutput()` runs the `in2csv` command and captures the output, which in this case is a csv version of the file. Then using the `decode()` method, I ensured the text was using the UTF-8 character encoding.

`StringIO` converts a string of text to a text stream, which `pandas` needs to read it into a data frame.

## Cleaning and counting the data

---

To clean the monthly data, I used the same process as I used for the historical data. For counting the crimes, I used the same process except for one change.

```
1 | #Save the count with the word 'Month'  
2 | crime_dict['Month'] = len(crime_subset)
```

When saving the crime count to the dictionary, I just save it with the key 'Month.' This allows me to get the count later by just using the key 'Month,' rather than trying to figure out which month it currently is. This is a sample of the result.

Crime	Month
ALCOHOL - MINOR IN POSSESSION	2
LARCENY - STOLEN BIKE	6
WEAPONS - CONCEALED	1
...	...

## Filtering newsworthy crimes

---

Now that I had the historical averages and this month's data, the only thing left was to see if this month's data exceeded the thresholds.

However, there were some crimes that weren't newsworthy that met the threshold because they were so rare. One standard deviation above the mean for False Security Alarms - Mechanical is 0.126. If one crime was reported under that category, it would trigger an alert, but it's not worth a story.

At the same time, there are crimes that would trigger an alert if they happened once, and they would be newsworthy. For example, one standard deviation above the mean for suicides is also 0.126. But, we figured those reports would be covered anyway, and it would be more obvious if there was a strange spike in those crimes. So, this program focuses on the crimes that happen more than a couple times per month that aren't normally covered by the media. For us, those crimes include:

Python

```

1 | flagged_crimes = [
2 |     "LOST OR STOLEN ITEM",
3 |     "FRAUD - CREDIT CARDS/ATM/BANK CARD",
4 |     "LARCENY - FROM MOTOR VEHICLE",
5 |     "NARCOTICS - POSSESSION",
6 |     "BURGLARY",
7 |     "LARCENY - FROM BUILDING",
8 |     "ALCOHOL - DWI",
9 |     "ALCOHOL - DRUNK",
10 |    "ALCOHOL - MINOR IN POSSESSION",
11 |    "VANDALISM - OTHER",
12 |    "LARCENY - STOLEN BIKE",
13 |    "VANDALISM - BY GRAFFITI",
14 |    "NARCOTICS - OTHER",
15 |    "NARCOTICS - SALE/DELIVER",
16 | ]

```

If someone wanted to remove or add a crime, they would just need to remove it from or add it to this list.

## Checking the upper threshold

If crime reports crossed the upper threshold at any point, it would be newsworthy. So, the upper threshold is checked every time the program runs. First, I needed to combine the historical data with the current month's data.

```

1 | merged = pd.merge(all_years_stats, month_count, on='Crime', how='inner')

```

Python

The `merge()` function combines two dataframes on a common column. In this case, I wanted to join the dataframes on the Crime column. This way, I could have the thresholds and the current counts in the same rows.

Since I was checking the upper threshold, I only wanted crimes that had occurred this month. So, I used an 'inner' join, which only includes data that is present in both sets.

Next, I needed to check if the current month's count had exceeded the upper threshold.

Python

```

1 plural_msg = "This month, there have been {month_total} {crime} incidents reported.
2 sing_msg = "This month, there has been {month_total} {crime} incident reported. There are
3
4 for index, row in merged.iterrows():
5     if row['upper'] < row['Month'] and row['Crime'] in flagged_crimes:
6         #If it has happened more than once, use plural words
7         if row['Month'] != 1:
8             plural_msg.format(
9                 crime=row['Crime'],
10                bound=round(row['upper'], 2),
11                month_total=row['Month'],
12                mean=round(row['mean'], 2),
13                direction='above',
14            )
15        #Otherwise, use singular words
16        else:
17            message = sing_msg.format(
18                crime=row['Crime'],
19                bound=round(row['upper'], 2),
20                month_total=row['Month'],
21                mean=round(row['mean'], 2),
22                direction='above',
23            )

```

I created two template messages. It iterates through each crime type in the set. If the count is over the upper bound, and if the crime is in `flagged_crimes`, it will create a message. It checks how many reports there are and uses the appropriate template, substituting in values using the `format()` method.

## Checking the lower threshold

To check the lower threshold at the end of the month, I needed a way to check if it was the last day of the month.

Python

```

1 today = datetime.today()
2 #monthrange() returns weekday of first of the month and number of days in month.
3 if today.day == monthrange(today.year, today.month)[1]:
4     return True
5 else:
6     return False

```

This function uses the `datetime` library to get today's date. Then, it uses the `monthrange()` function from the `calendar` library to get the last day of the month. If they're the same, it's the last day of the month.

## Comparing the values

Comparing the counts to the lower threshold is similar to upper threshold process, but with a few changes.

```
1 | merged = pd.merge(all_years_stats, month_count, on='Crime', how='outer')
```

Python

When joining the two data frames, an outer join is performed. Outer joins keep all information, even if it's not present in one of the data frames. This allows us to check crime types that have not occurred in a month.

```
1 | for index, row in merged.iterrows():
2 |     if row['lower'] > row['Month'] and row['Crime'] in flagged_crimes:
3 |         #If it has happened more than once, use plural words
4 |         if row['Month'] != 1:
5 |             message = plural_msg.format(
6 |                 crime=row['Crime'],
7 |                 bound=round(row['lower'], 2),
8 |                 month_total=row['Month'],
9 |                 mean=round(row['mean'], 2),
10 |                direction='below',
11 |            )
12 |         #Otherwise, use singular words
13 |     else:
14 |         message = sing_msg.format(
15 |             crime=row['Crime'],
16 |             bound=round(row['lower'], 2),
17 |             month_total=row['Month'],
18 |             mean=round(row['mean'], 2),
19 |             direction='below',
20 |         )
```

Python

Then, I compared the count to the lower threshold, instead of the upper threshold. If the count is lower than the threshold, it will create a message, using the same templates from before.

# Posting to Slack

---

Having the messages in the program is great, but it requires someone to constantly check to see if something has popped up. Instead, we can send the messages to Slack, alerting people when a new message is created. Slack is a common program for newsrooms to use, making this an easy integration for them.

## Setting up Slack

This feature requires a Slack workspace with the ability to add an [Incoming Webhook](#). Once you add a configuration and choose a channel for it to post to, you're good to go. Just save the Webhook URL as an environment variable `SLACK_URL`, or replace `os.environ.get('SLACK_URL')` with your URL.

```
1 | slack_url = os.environ.get('SLACK_URL')
```

Python

## Sending the message

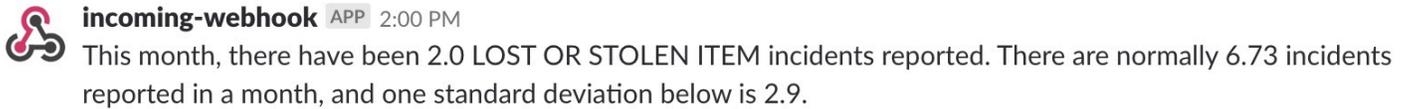
If the URL is setup, the only thing left is to send the message.

```
1 | #Put the message in a dictionary
2 | slack_data = {'text': message}
3 | #Send the message
4 | response = requests.post(
5 |     #Convert the dictionary to a JSON object
6 |     webhook_url, data=json.dumps(slack_data),
7 |     #These headers help Slack interpret the message
8 |     headers={'Content-Type': 'application/json'})
9 | )
10 | if response.status_code != 200:
11 |     raise ValueError(
12 |         'Request to slack returned an error {code}, the response is:\n{text}'.format(
13 |             code=response.status_code,
14 |             text=response.text,
15 |         )
16 |     )
```

Python

This uses the `requests` library to send the message, which is formatted using the `json` library. A status code of '200' means it worked, so if it doesn't equal '200,' an error is raised.

Here's an example message.



## Outcome/applications

To make this program run on a regular basis, I recommend using a task scheduler like `cron`. This will allow the program to run as often as you'd like, but it will only send messages when the thresholds are reached. Now, if a crime does reach the upper threshold at some point during the month, it will continue sending that message until the end of the month. To prevent it from doing so, you can remove it from the `flagged_crimes` list until the next month.

Also, if you want the program to run faster, you can just use the saved version of the statistics, rather than calculating them again each time.

```
1 | all_years_stats = pd.read_csv('std.csv')
```

Python

This line uses the most recent version of the data that was saved after calculating the thresholds. This allows you to comment out the commands which count the historical crimes and calculate the thresholds until you update the historical data.

While this program does a lot on its own, it does require some maintenance. As mentioned earlier, the monthly data that's downloaded each time must be deleted before the program runs again. Otherwise, the program will not use the new data. Also, this program could be used by almost any person who had an interest in UNLPD data. While that audience may be fairly small, this program could also be used for other data sources, but it would take some tweaking.

At the same time, the historical data is not set to update automatically. If you want data past December 31, 2017, to be included in the calculations of the thresholds, you'll need to download it from the Daily Crime and

Fire Log and stack it with the `all_years.csv` file, as mentioned at the beginning of the Process section.

## Conclusion

In a world where data is created every second, journalists need to be able to use tools to help sort through that data in an efficient manner. While not a complete product, this program is a proof of concept for a program that could fill this role with UNLPD crime data.

This program is not a definitive solution for finding newsworthy trends in UNLPD crime. It should be used as a tool, in conjunction with standard journalism research and intuition. It may miss some outliers, and it may produce some false positives. It is therefore imperative that users take the time to investigate any messages this program produces.

For example, during the development of this product, an alert was created for four rape reports during the month of February. The average number of rape reports in a month is 0.35. After closer inspection, none of the incidents occurred in February; they were just reported in the same month. That may still warrant a story, but it may not be the story a journalist would expect at first glance.

There were also 27 narcotics possession reports in February, and the average number of reports in a month is 9.5. The most incidents reported in a single day during the month was four, so the rise in reports might not have been completely obvious to a journalist scanning the data. Using this program, he or she would have received a message as soon as the report count crossed the upper threshold, alerting them to the high number of reports.

Without a live newsroom to deploy this program, it's currently uncertain just how useful this program is in its current state. But, with testing over time and continued adjustments, this could be a viable product to help journalists find stories in data.

# Appendix I: The full program

Python

```

1  import os #To get Slack API key/expand user directory path
2  import pandas as pd #For most data manipulations
3  import json #To prepare message for Slack
4  import requests #To send message to Slack
5  import subprocess #To run in2csv
6  import time #To pause after downloading the file
7  from io import StringIO #To convert a string to a file
8  from datetime import datetime #To check today's date
9  from calendar import monthrange #To find last day of month
10 from selenium import webdriver #To scrape UNLPD's data
11 from selenium.webdriver.common.by import By
12 from selenium.webdriver.support.ui import WebDriverWait
13 from selenium.webdriver.support import expected_conditions as EC
14
15 def clean_data(df):
16     # Make sure there are no null values for the Case Number, Reported time, Location,
17     # Stolen amount and Damaged amount
18     assert df['Case #'].count() == len(df) != 0
19     assert df['Reported'].count() == len(df) != 0
20     assert df['Location'].count() == len(df) != 0
21     assert df['Stolen'].count() == len(df) != 0
22     assert df['Damaged'].count() == len(df) != 0
23
24     #Replace non-numerical characters and cast data type to float
25     df['Stolen'] = df['Stolen'].str.replace(',','')
26     df['Stolen'] = df['Stolen'].str.replace('$','')
27     df['Stolen'] = df['Stolen'].astype(float)
28     df['Damaged'] = df['Damaged'].str.replace(',','')
29     df['Damaged'] = df['Damaged'].str.replace('$','')
30     df['Damaged'] = df['Damaged'].astype(float)
31
32     #Cast data type to datetime
33     df['Reported'] = pd.to_datetime(df['Reported'])
34
35     #Double-check data types
36     print(df.dtypes)
37
38     #Create a new column with just the year and month from the Reported column
39     df['Month'] = df['Reported'].dt.to_period('M')
40

```

```

41     #Set index to Reported column; allows for slicing by month
42     df2 = df.set_index(['Reported'])
43
44     return df2
45
46 def count_crimes(df, all_years=False):
47     #This holds the dictionary for each crime
48     months_count = []
49     #For each crime present in the dataframe
50     for crime in df['Incident Code'].unique():
51         print(crime)
52         crime_dict = {}
53         crime_dict['Crime'] = crime
54         #For each month in the dataframe
55         for month in df['Month'].unique():
56             #Slice the dataframe for one month's data
57             month_subset = df[str(month)]
58             #Filter the subset for instances of the crime
59             crime_subset = month_subset[month_subset['Incident Code'] == crime]
60             #If multiple months, save the count with the month
61             if all_years:
62                 crime_dict[str(month)] = len(crime_subset)
63             #Otherwise, just save it with "Month"
64             else:
65                 crime_dict['Month'] = len(crime_subset)
66             #Append the dictionary to the months_count list
67             months_count.append(crime_dict)
68     #Convert the list into another dataframe
69     months_count_df = pd.DataFrame(months_count)
70     #To help with speed, save it to a csv
71     if all_years:
72         months_count_df.to_csv('month_count.csv', index=False)
73     return months_count_df
74
75 def calculate_stats(df):
76     #Creates a dataframe with the unique crimes
77     std_df = df.filter(['Crime'])
78     #Adds a column with the mean count for each crime
79     std_df['mean'] = df.mean(axis=1)
80     #Adds a column with the standard deviation for each crime
81     std_df['std'] = df.std(axis=1)
82     #Adds a column with a lower threshold
83     std_df['lower'] = std_df['mean'] - std_df['std']
84     #Adds a column with an upper threshold
85     std_df['upper'] = std_df['mean'] + std_df['std']

```

```
86     #Save the data to a csv
87     std_df.to_csv('std.csv', index=False)
88     return std_df
89
90 def check_last_day():
91     #Get today's date
92     today = datetime.today()
93     #monthrange() returns weekday of first of the month and number of days in month.
94     if today.day == monthrange(today.year, today.month)[1]:
95         return True
96     else:
97         return False
98
99 def post_to_slack(message):
100     #Put the message in a dictionary
101     slack_data = {'text': message}
102     #Send the message
103     response = requests.post(
104         #Convert the dictionary to a JSON object
105         webhook_url, data=json.dumps(slack_data),
106         #These headers help Slack interpret the message
107         headers={'Content-Type': 'application/json'})
108
109     if response.status_code != 200:
110         raise ValueError(
111             'Request to slack returned an error {code}, the response is:\n{text}'.format(
112                 code=response.status_code,
113                 text=response.text,
114             )
115         )
116
117 def find_outliers(all_years_stats, month_count):
118     #This is the list of crimes we decided we were interested in
119     flagged_crimes = [
120         "LOST OR STOLEN ITEM",
121         "FRAUD - CREDIT CARDS/ATM/BANK CARD",
122         "LARCENY - FROM MOTOR VEHICLE",
123         "NARCOTICS - POSSESSION",
124         "BURGLARY",
125         "LARCENY - FROM BUILDING",
126         "ALCOHOL - DWI",
127         "ALCOHOL - DRUNK",
128         "ALCOHOL - MINOR IN POSSESSION",
129         "VANDALISM - OTHER",
130         "LARCENY - STOLEN BIKE",
```

```

131     "VANDALISM - BY GRAFFITI",
132     "NARCOTICS - OTHER",
133     "NARCOTICS - SALE/DELIVER",
134 ]
135 #These two templates are used for the messages.
136 plural_msg = "This month, there have been {month_total} {crime} incidents reported. T
137 sing_msg = "This month, there has been {month_total} {crime} incident reported. There
138
139 #If it's the last day, merge the data and keep everything
140 #Then, check the low thresholds
141 if check_last_day():
142     merged = pd.merge(all_years_stats, month_count, on='Crime', how='outer')
143     for index, row in merged.iterrows():
144         if row['lower'] > row['Month'] and row['Crime'] in flagged_crimes:
145             #If it has happened more than once, use plural words
146             if row['Month'] != 1:
147                 message = plural_msg.format(
148                     crime=row['Crime'],
149                     bound=round(row['lower'], 2),
150                     month_total=row['Month'],
151                     mean=round(row['mean'], 2),
152                     direction='below',
153                 )
154             #Otherwise, use singular words
155             else:
156                 message = sing_msg.format(
157                     crime=row['Crime'],
158                     bound=round(row['lower'], 2),
159                     month_total=row['Month'],
160                     mean=round(row['mean'], 2),
161                     direction='below',
162                 )
163             #Print the message here
164             print(message)
165             #Post the message to Slack
166             post_to_slack(message)
167 #Otherwise, only keep the data for crimes that have happened this month
168 else:
169     merged = pd.merge(all_years_stats, month_count, on='Crime', how='inner')
170 #For each row, check if the count has crossed the upper bound
171 for index, row in merged.iterrows():
172     if row['upper'] < row['Month'] and row['Crime'] in flagged_crimes:
173         #If it has happened more than once, use plural words
174         if row['Month'] != 1:
175             plural_msg.format(

```

```

176         crime=row['Crime'],
177         bound=round(row['upper'], 2),
178         month_total=row['Month'],
179         mean=round(row['mean'], 2),
180         direction='above',
181     )
182     #Otherwise, use singular words
183     else:
184         message = sing_msg.format(
185             crime=row['Crime'],
186             bound=round(row['upper'], 2),
187             month_total=row['Month'],
188             mean=round(row['mean'], 2),
189             direction='above',
190         )
191     #Print the message here
192     print(message)
193     #Post the message to Slack
194     post_to_slack(message)
195
196 #Read in the csv file
197 all_years = pd.read_csv('all_years.csv')
198 #Clean the data
199 all_years_clean = clean_data(all_years)
200 #Count the crime occurrences
201 all_years_count = count_crimes(all_years_clean, all_years=True)
202 #Calculate the thresholds
203 all_years_stats = calculate_stats(all_years_count)
204
205 #This is needed to set up selenium
206 #os.path.expanduser allows the use of a '~'
207 path_to_chromedriver = os.path.expanduser('~Downloads/chromedriver')
208 browser = webdriver.Chrome(executable_path=path_to_chromedriver)
209 #The URL to the Daily Crime and Fire Log
210 url = "https://scsapps.unl.edu/policerreports/MainPage.aspx"
211 #Go to the URL
212 browser.get(url)
213 #Find the advanced search button and click it
214 browser.find_element_by_id('ctl00_ContentPlaceholder1_AdvancedSearchButton').click()
215 #Find the first date field, hit tab and hit '01'.
216 #This sets the date to the first day of the month
217 date_box = browser.find_element_by_id('ctl00_ContentPlaceholder1_DateRange_MonthText1')
218 date_box.send_keys('\t01')
219 #Find the search button and click it
220 browser.find_element_by_id('ctl00_ContentPlaceholder1_SearchButton').click()

```

```
221 #Switch to the iframe on the page
222 browser.switch_to.frame(browser.find_element_by_id('ctl00_ContentPlaceHolder1_ViewPort'))
223 #Find the export button once the iframe loads and click it
224 export_button = WebDriverWait(browser, 10).until(
225     EC.presence_of_element_located((By.ID, 'ExportButton'))
226 )
227 export_button.click()
228
229 #Wait for the file to download
230 time.sleep(5)
231
232 #Runs in2csv on the downloaded file and converts it to UTF-8
233 csv_data = subprocess.check_output([
234     "in2csv",
235     os.path.expanduser("~/Downloads/DailyCrimeLogSummary.xls"),
236 ], stderr=subprocess.DEVNULL).decode("utf-8")
237 #Creates a file instance for pandas to use on the next line
238 csv_file_instance = StringIO(csv_data)
239 #Reads in the csv to a dataframe, skipping the first eight rows
240 month_df = pd.read_csv(csv_file_instance, skiprows=8)
241
242 #Clean the data
243 month_clean = clean_data(month_df)
244 #Count the crime occurrences
245 month_count = count_crimes(month_clean)
246
247 webhook_url = os.environ.get('SLACK_URL')
248
249 find_outliers(all_years_stats, month_count)
```