



ECE 356 • Crime Statistics Database Project

```
SELECT
    Student,
    Course,
    Term,
    Institution
FROM WaterlooStudentsView
ORDER BY Student ASC;
```

<i>Student</i>	<i>Course</i>	<i>Term</i>	<i>Institution</i>
Kyle Pinto	ECE 356 - Database Systems	Fall 2021	University of Waterloo
Puranjoy Saha	ECE 356 - Database Systems	Fall 2021	University of Waterloo
Zahin Zaman	ECE 356 - Database Systems	Fall 2021	University of Waterloo

Table of Contents

Introduction

Overview

Source Code

Datasets

London Police Records

NYPD Complaint Data Historic

Crimes in Chicago

LA Crimes

Design

Analysis of Datasets

Design Options

UK & US Data Separation Option

Crimes, Complaints & Stop-and-Searches Option

Crimes Only Option

Decision Matrix

Entity-Relationship Model

Location

Code

Person

Incident

Complaint

Crime	
Search	
Relational Schema	
Tables	
Foreign Key Constraints	
Indexes	
Views	
Client Application	
Overview	
Installation	
Configuration	
Usage Breakdown	
Ideal command line tool	
High level Commands	
Database Creation	
Subcommands	
Walkthrough	
Background check	
Adding a crime to the database	
Filter for crimes at a specific location	
Testing Plan	
Data Mining	
Goal	
Feature Selection	
Results	
Ethnicity	
Gender	
Conclusion	
Challenges	
Tradeoffs	
Future Improvements	
Commands to be implemented in the future	

Introduction

Overview

This project involves the collection of crime records datasets from law enforcement departments in UK and US, and the process of developing an optimally designed database and a client interface for the definition, manipulation and storage of this data.

Source Code

The source code for the project can be found under the following repository:

<https://github.com/alvii147/CrimeStatsAnalysis.git>

Datasets

Datasets used in this project have been collected from [Kaggle](#). Each dataset provides crime data from a different city and state in UK or US.

Name	Location	Link
London Police Records	London, England, UK	[↗]
NYPD Complaint Data Historic	New York City, New York, US	[↗]
Crimes in Chicago	Chicago, Illinois, US	[↗]
LA Crime Data	Los Angeles, California, US	[↗]

London Police Records

This dataset includes crime data from London from late 2014 to mid 2017, held in the following three CSV files:

- [london-outcomes.csv](#)
- [london-street.csv](#)
- [london-stop-and-search.csv](#)

[london-outcomes.csv](#) and [london-street.csv](#) hold data on instances of crime committed in London and their relevant information. [london-stop-and-search.csv](#) holds data on "stop-and-searches" conducted by London police and their relevant information.

The data references the location of each incident by the Lower Layer Super Output Area (LSOA) code of the neighborhood, which can be mapped to specific area names using the [Lookup Table of UK Local Government Areas](#) dataset.

NYPD Complaint Data Historic

This dataset includes records of complaints of incidents reported to the New York City Police Department (NYPD) from 2006 to the end of 2017, in CSV file [NYPD_Complaint_Data_Historic.csv](#). The data contained in this file also includes the NYPD crime code corresponding to each complaint, which is unique to a specific type of crime.

Crimes in Chicago

This dataset includes crimes reported and committed from the records of the Chicago Police Department (CPD) between 2001 and 2017, divided into four CSV files:

- [Chicago_Crimes_2001_to_2004.csv](#)
- [Chicago_Crimes_2005_to_2007.csv](#)
- [Chicago_Crimes_2008_to_2011.csv](#)
- [Chicago_Crimes_2012_to_2017.csv](#)

The columns of this dataset includes the Illinois Uniform Crime Reporting (IUCR) code corresponding to the committed crime, which can be extracted from the [IUCR](#) dataset hosted on the City of Chicago website.

LA Crimes

This dataset includes crimes reported and committed from the records of the Los Angeles Police Department (LAPD) between 2010 and mid-2021, divided into two CSV files:

- `Crime_Data_from_2010_to_2019.csv`
- `Crime_Data_from_2020_to_Present.csv`

The columns of this dataset includes the LAPD crime code corresponding to the committed crime, which is based on the FBI Uniform Crime Reporting codes. These codes can be compiled from the [FBI UCR Handbook](#).

Design

Analysis of Datasets

The first step of the design process is a thorough investigation of the given datasets and their attributes in order to prepare practical options for merging the datasets. By examining each dataset, we discover the following:

- Every dataset entry has a location associated with it. Usually that location is specified by latitude and longitude coordinates, along with a few other address parameters that depend on the area, such as ward, precinct, LSOA code, borough, city, state, country, etc. These address parameters, however, vary significantly between UK and US datasets.
- Every dataset entry also has a few attributes that are common between all or most of the datasets. These attributes include date of occurrence, type of crime, description of crime, and other similar attributes that describe a general incident.
- Every US dataset entry has a unique crime code defined that describes the category of the crime. The uniqueness of these codes also depend on the organization that reports this crime data. For instance, NYPD and IUCR crime codes are not identical.
- Some datasets include victim information, and some do not. Because information related to individual people may lead to privacy issues, the datasets omit personal information such as names, contact numbers etc. and only store their ages (or age ranges), genders and ethnicities.
- The `london-stop-and-searches.csv` file from the London Police Records dataset contains information that is slightly different from the rest of the datasets. It includes information about stop-and-searches conducted by London police, which may or may not have resulted in the discovery of criminal activity. While there are common attributes between this and other datasets, such as location and date, the context of this information is different.

Design Options

From the analysis of the datasets, the most obvious design options we can draw are the definitions of separate entities for crime codes and for individual people.

Each crime code entry should be unique depending on the code and the organization that reports that crime data. Each code should also have a category definition and a description of the crime.

Individual people can also be considered a separate entity. A person could be described as the victim of a crime, the perpetrator of a crime, or the suspect of a stop-and-search. An issue that can be identified here is that the datasets have omitted personal information to avoid privacy issues, which includes information that is typically present in a police department database. Thus, for the sake of completeness, it may make

sense to generate fake names and phone numbers to go with the rest of the information about individual people.

The rest of the information from the datasets mostly consist of partially overlapping attributes, which necessitates the exploration of the differences between the datasets and the design choices to accommodate these differences.

UK & US Data Separation Option

There is a significant difference between location-based attributes of the UK and US datasets. UK addresses use attributes such as LSOA codes and regions for location which usually is not relevant for US addresses, and US addresses use attributes such as precinct, ward, district and state for location, which is irrelevant for UK addresses.

Additionally, UK datasets do not report any crime codes that correspond to the crimes. This is only relevant to crimes in US.

Considering these points, we can outline the separation of UK and US datasets into individual UK and US based addresses and crimes as a practical design option.

Crimes, Complaints & Stop-and-Searches Option

Another practical design option is to divide the datasets into three separate entities: complaints, which would include information from the NYPD Complaints Data Historic dataset, stop-and-searches, which would include information from the `london-stop-and-searches.csv` file of the London Police Records dataset, and crimes, which would include crime data from all other datasets.

This design option considers the different entities that the datasets provide and try to minimize the number of irrelevant attributes, while also providing a generalized division between the dataset attributes.

Crimes Only Option

In order to avoid making the entities of the database too specific, we can also consider further generalizing all the information from all the datasets into a single entity. The advantage of this option is avoidance of over-specialization, but has a possible disadvantage of too many irrelevant attributes caused by over-generalization.

Decision Matrix

We can consider our outlined options and construct a weighted decision matrix to decide between our design choices. We will weigh our alternatives in terms of the following criteria:

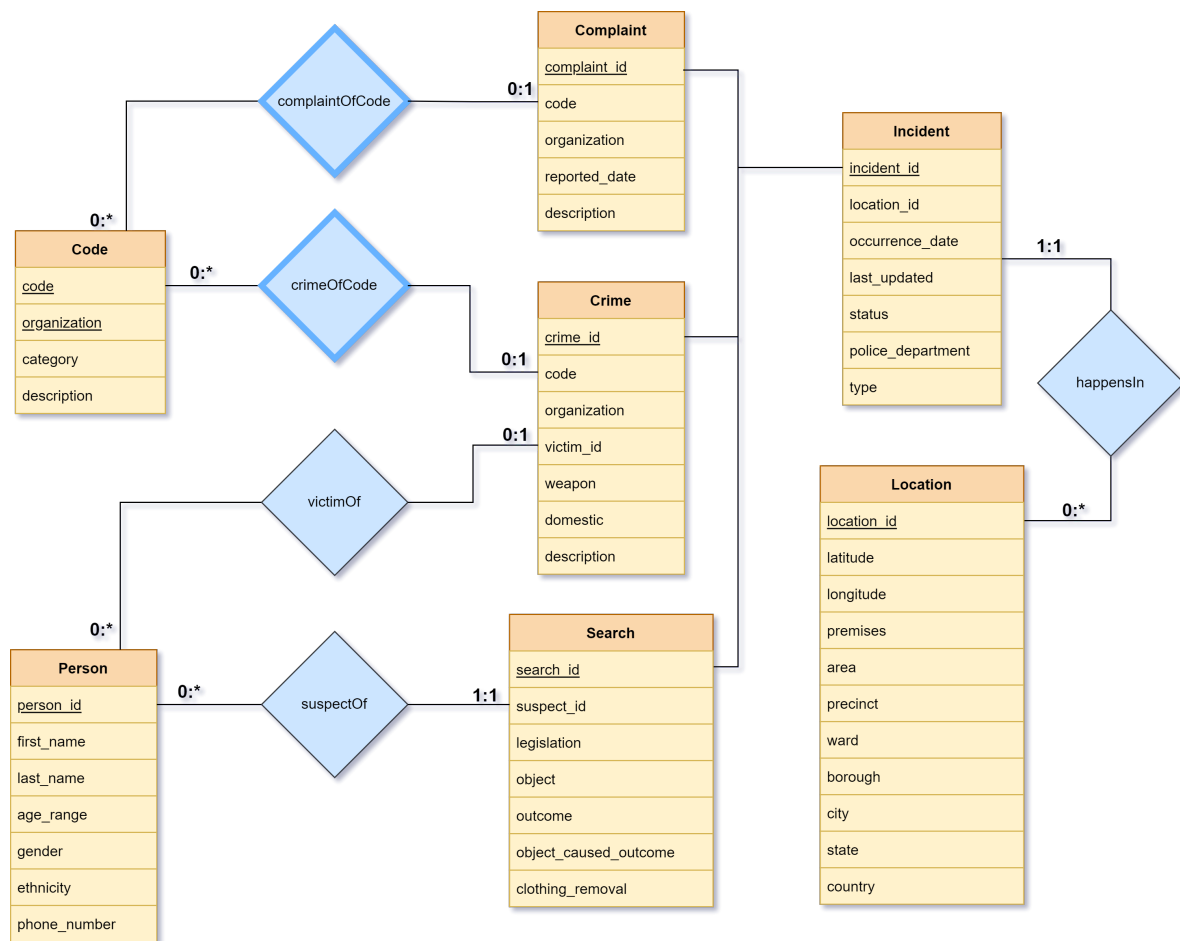
- **Generalization:** database should be a good generalization of the combined datasets
- **Simplicity:** entities and relationships within the database should not be too complex
- **Minimized Data Loss:** loss of data due to merging should be reduced
- **Attribute Relevance:** attributes within the entities should be relevant for most of the database entries
- **Accommodation of Needs:** database should perform operations based on the users' (i.e. police station employees) needs

	Generalization	Simplicity	Minimized Data Loss	Attribute Relevance	Accommodation of needs	Score
Weights	2	1	3	4	5	
UK & US Data Separation	1	1	4	5	2	45
Crimes, Complaints & Stop-and-Searches	3	4	3	4	5	60
Crimes Only	5	5	4	1	3	46

By analyzing our decision options using a decision matrix, we can conclude that separating crimes, complaints and stop-and-searches while merging the UK and US datasets is the most optimal solution given the context of our project.

Entity-Relationship Model

We can now build an entity-relationship model based on our selected design options.



Location

The *Location* entity is a combination of all the location-based attributes from all datasets. This includes attributes that are common between UK and US datasets, such as `latitude` , `longitude` , `premises` , `city` and `country` , as well other attributes that are only unique to either UK or US, such as `area` , `precinct` , `ward` , `borough` , and `state` . Note that LSOA code was omitted from this entity since it was uniquely functionally dependent on `borough` , and not particularly relevant if we already have the name of the borough.

The primary key for this entity is an artificial primary key, `location_id` .

Attribute	Description
<u><code>location_id</code></u>	Artificial primary key
<code>latitude</code>	Latitude coordinate where the incident took place
<code>longitude</code>	Longitude coordinate where the incident took place
<code>premises</code>	Contextual premises of the incident
<code>area</code>	Area where the incident took place
<code>precinct</code>	Police precinct where the incident was reported
<code>ward</code>	Ward where the incident took place
<code>borough</code>	Borough where the incident took place
<code>city</code>	City where the incident took place
<code>state</code>	State/province where the incident took place
<code>country</code>	Country where the incident took place

Code

The *Code* entity contains all the different US crime codes as reported by the NYPD, IUCR and the LAPD, along with additional information regarding the codes, including `category` and `description` .

Since the uniqueness of the entries of this entity is dependent on both the crime code and the reporting organization, the primary key for this entity is a composite key made up of attributes `code` and `organization` .

Attribute	Description
<u><code>code</code></u>	Crime code
<u><code>organization</code></u>	Reporting organization
<code>category</code>	Category of crime code
<code>description</code>	Description of crime code

Person

The *Person* entity holds information about individual people that are relevant to the database (including information that may have to be auto-generated, such as `first_name` , `last_name` and `phone_number`).

The primary key for this entity is an artificial primary key, `person_id` .

Attribute	Description
<code><u>person_id</u></code>	Artificial primary key
<code>first_name</code>	First name of person
<code>last_name</code>	Last name of person
<code>age_range</code>	Age range of person
<code>gender</code>	Gender of person
<code>ethnicity</code>	Ethnicity of person
<code>phone_number</code>	Phone number of person

Incident

The *Incident* entity represents an incident reported in any of the datasets and is meant to hold general information that is relevant to all specific incidents. This includes attributes `occurrence_date` , `type` , `status` , `police_department` and `last_updated` .

This entity also contains a `location_id` attribute which is related to the *Location* entity through relation *happensIn* . Ideally, this is a many-to-one relation between *Incident* and *Location* (i.e. multiple incidents can happen in the same location).

The primary key for this entity is an artificial primary key, `incident_id` .

Attribute	Description
<code><u>incident_id</u></code>	Artificial primary key
<code>location_id</code>	Primary key of <i>Location</i> entity describing location of the incident
<code>occurrence_date</code>	Date of occurrence of incident
<code>last_updated</code>	Date when this record was last updated
<code>status</code>	Current status of the incident
<code>police_department</code>	Police department the incident was reported to
<code>type</code>	Type of crime

Complaint

The *Complaint* entity is a specialization of the *Incident* entity and provides additional information about incidents that are complaints about crimes, through attributes such as `reported_date` and `description`.

The `code` and `organization` attributes of the *Complaint* entity are related to the *Code* entity that describes the specific crime code referenced by the complaint, through the relation *complaintOfCode*.

Note

This is defined as a weak entity set because a complaint about a crime cannot be valid if there exists no criminal law (i.e. crime code) that prohibits that action.

The primary key for this entity is an artificial primary key, `complaint_id`.

Attribute	Description
<code>complaint_id</code>	Artificial primary key
<code>code</code>	Crime code
<code>organization</code>	Reporting organization
<code>reported_date</code>	Date of complaint report
<code>description</code>	Description of complaint

Crime

The *Crime* entity is a specialization of the *Incident* entity and provides additional information about incidents that are reported crimes, through attributes such as `weapon`, `domestic` and `description`.

Like those of the *Complaint* entity, the `code` and `organization` attributes of the *Crime* entity are related to the *Code* entity that describes specific crime code of the crime that's reportedly committed, through the relation *crimeOfCode*. Note that, once again, this is a weak entity set.

The `victim_id` attribute is related to the *Person* entity and describes the information related to the victim of the crime, through relation *victimOf*.

The primary key for this entity is an artificial primary key, `crime_id`.

Attribute	Description
<code>crime_id</code>	Artificial primary key
<code>code</code>	Crime code
<code>organization</code>	Reporting organization
<code>victim_id</code>	Primary key of <i>Person</i> entity victim to crime
<code>weapon</code>	Weapon used
<code>domestic</code>	Indicates whether the crime was domestic
<code>description</code>	Description of the crime committed

Search

The *Search* entity is a specialization of the *Incident* entity and provides additional information about incidents that are stop-and-searches conducted by the police, through attributes such as `legislation`, `object`, `outcome`, `object_caused_outcome` and `clothing_removal`.

The `suspect_id` attribute of the *Search* entity is related to the *Person* entity and describes the information related to the suspect of the stop-and-search, through relation *suspectOf*.

The primary key for this entity is an artificial primary key, `search_id`.

Attribute	Description
<code>search_id</code>	Artificial primary key
<code>victim_id</code>	Primary key of <i>Person</i> entity suspect to search
<code>legislation</code>	Legislation under which search was conducted
<code>object</code>	Object of search
<code>outcome</code>	Outcome of search
<code>object_caused_outcome</code>	Indicates whether outcome was a result of search object
<code>clothing_removal</code>	Indicates whether clothing of suspect was removed for search

Relational Schema

Tables

Now we can convert our entity-relationship model to relational schema. We can start by constructing the `CREATE TABLE` commands for our required tables.

```
CREATE TABLE Incident (  
    incident_id INT NOT NULL AUTO_INCREMENT,  
    location_id INT,  
    occurrence_date DATE,  
    last_updated DATE,  
    status VARCHAR(128),  
    police_department VARCHAR(256),  
    type VARCHAR(128),  
    PRIMARY KEY(incident_id)  
);
```

```
CREATE TABLE Location (  
    location_id INT NOT NULL AUTO_INCREMENT,  
    latitude DECIMAL(11, 8),  
    longitude DECIMAL(11, 8),  
    premises VARCHAR(128),  
    area VARCHAR(256),  
    precinct DECIMAL(4),  
    ward DECIMAL(3),  
    borough VARCHAR(64),  
    city VARCHAR(64),
```

```
state VARCHAR(64),
country VARCHAR(64),
PRIMARY KEY(location_id)
);
```

```
CREATE TABLE Crime (
  crime_id INT NOT NULL AUTO_INCREMENT,
  incident_id INT,
  code VARCHAR(4),
  organization VARCHAR(16),
  victim_id INT,
  weapon VARCHAR(256),
  domestic BOOL,
  description VARCHAR(256),
  PRIMARY KEY(crime_id)
);
```

```
CREATE TABLE Complaint (
  complaint_id INT NOT NULL AUTO_INCREMENT,
  incident_id INT,
  code VARCHAR(4),
  organization VARCHAR(16),
  reported_date DATE,
  description VARCHAR(256),
  PRIMARY KEY(complaint_id)
);
```

```
CREATE TABLE Search (
  search_id INT NOT NULL AUTO_INCREMENT,
  incident_id INT,
  suspect_id INT,
  legislation VARCHAR(256),
  object VARCHAR(256),
  outcome VARCHAR(256),
  object_caused_outcome BOOL,
  clothing_removal BOOL,
  PRIMARY KEY(search_id)
);
```

```
CREATE TABLE Person (
  person_id INT NOT NULL AUTO_INCREMENT,
  first_name VARCHAR(64),
  last_name VARCHAR(64),
  age_range VARCHAR(16),
  gender VARCHAR(16),
  ethnicity VARCHAR(64),
  phone_number VARCHAR(16),
  PRIMARY KEY(person_id)
);
```

```
CREATE TABLE Code (
  code VARCHAR(4) NOT NULL,
  organization VARCHAR(16) NOT NULL,
  category VARCHAR(256),
  description VARCHAR(256),
  PRIMARY KEY(code, organization)
);
```

Note

A difference between the entity-relationship model and our defined tables is that the `Crime`, `Complaint` and `Search` tables now have an additional attribute, `incident_id` which was not reflected on the entity-relationship model. This is because the relationship between `Crime`, `Complaint` and `Search` are specializations of `Incident`, and so must be related using a foreign key attribute, i.e. `incident_id`.

Our relational schema is almost nearly similar to our entity-relationship model. We set the `AUTO_INCREMENT` option for all our artificial primary keys so that they are sequentially incremented automatically on insertion.

Foreign Key Constraints

Now we can define our foreign key constraints based on the entity-relationship model relations.

```
ALTER TABLE Complaint ADD CONSTRAINT Complaint_Incident
    FOREIGN KEY (incident_id) REFERENCES Incident(incident_id);
```

```
ALTER TABLE Crime ADD CONSTRAINT Crime_Incident
    FOREIGN KEY (incident_id) REFERENCES Incident(incident_id);
```

```
ALTER TABLE Search ADD CONSTRAINT Search_Incident
    FOREIGN KEY (incident_id) REFERENCES Incident(incident_id);
```

```
ALTER TABLE Incident ADD CONSTRAINT happensIn
    FOREIGN KEY (location_id) REFERENCES Location(location_id);
```

```
ALTER TABLE Complaint ADD CONSTRAINT complaintOfCode
    FOREIGN KEY (code, organization) REFERENCES Code(code, organization);
```

```
ALTER TABLE Crime ADD CONSTRAINT crimeOfCode
    FOREIGN KEY (code, organization) REFERENCES Code(code, organization);
```

```
ALTER TABLE Crime ADD CONSTRAINT victimOf
    FOREIGN KEY (victim_id) REFERENCES Person(person_id);
```

```
ALTER TABLE Search ADD CONSTRAINT suspectOf
    FOREIGN KEY (suspect_id) REFERENCES Person(person_id);
```

Indexes

Based on our primary and foreign keys, we can also set up appropriate indexes in order to improve the performance of our database.

```
CREATE INDEX Incident_PK_IDX
ON Incident (incident_id);
```

```
CREATE INDEX Incident_Location_FK_IDX
ON Incident (location_id);
```

```
CREATE INDEX Location_PK_IDX
ON Location (location_id);
```

```
CREATE INDEX Crime_PK_IDX
ON Crime (crime_id);
```

```
CREATE INDEX Crime_Incident_FK_IDX
ON Crime (incident_id);
```

```
CREATE INDEX Crime_Code_FK_IDX
ON Crime (
    code,
    organization
);
```

```
CREATE INDEX Crime_Person_FK_IDX
ON Crime (victim_id);
```

```
CREATE INDEX Complaint_PK_IDX
ON Complaint (complaint_id);
```

```
CREATE INDEX Complaint_Incident_FK_IDX
ON Complaint (incident_id);
```

```
CREATE INDEX Complaint_Code_FK_IDX
ON Complaint (
    code,
    organization
);
```

```
CREATE INDEX Search_PK_IDX
ON Search (search_id);
```

```
CREATE INDEX Search_Incident_FK_IDX
ON Search (incident_id);
```

```
CREATE INDEX Search_Person_FK_IDX
ON Search (suspect_id);
```

```
CREATE INDEX Person_PK_IDX
ON Person (person_id);
```

```
CREATE INDEX Code_PK_IDX
ON Code (
    code,
    organization
);
```

Views

Finally, we can also create database views based on joins that can provide information that go together. This is to help make it easier to query for information from the client application. We create the following views (the SQL code for creation of views is omitted here due to length):

View Name	View Description
CrimeView	Attributes from a join between <code>Crime</code> , <code>Incident</code> , <code>Location</code> , <code>Code</code> and <code>Person</code> (<code>victim_id</code>)
ComplaintView	Attributes from a join between <code>Complaint</code> , <code>Incident</code> , <code>Location</code> and <code>Code</code>
SearchView	Attributes from a join between <code>Search</code> , <code>Incident</code> , <code>Location</code> , and <code>Person</code> (<code>suspect_id</code>)

Client Application

Overview

The client application serves as a front end for the database. It uses an interactive command line interface to maintain and look up records, as well as create, delete, and populate the database.

Installation

Clone the repository:

```
git clone https://github.com/alvii147/CrimeStatsAnalysis.git
```

Navigate into repository directory:

```
cd CrimeStatsAnalysis/
```

Create and activate Python virtual environment (optional):

```
python3 -m venv env
# Linux & MacOS
source env/bin/activate
# Windows
source env/Scripts/activate
```

Install dependencies:

```
pip3 install -r requirements.txt
```

Configuration

In order to run the program, you will need access to a database on Marmoset. For best results, use a empty database. You can store the database credentials in a configuration data file under `src/MySQLUtils/config.ini` .

Note

Using `config.ini` is not required but is recommended. If you choose to not use the configuration file, then the client application will prompt you for your credentials each time it connects to the database.

To store your credentials in a `config.ini` file:

```
echo [mysqlconfig] > src/MySQLUtils/config.ini
echo host = <hostname> >> src/MySQLUtils/config.ini
echo user = <username> >> src/MySQLUtils/config.ini
echo password = <password> >> src/MySQLUtils/config.ini
echo database = <databasename> >> src/MySQLUtils/config.ini
```

`<hostname>` is the hostname of the database, e.g. `marmoset04.shoshin.uwaterloo.ca`.

`<username>` and `<password>` are your MySQL username and password.

`<databasename>` is the name of your database on the server.

The file should look like this:

```
[mysqlconfig]
host = marmoset04.shoshin.uwaterloo.ca
user = waterlooid
password = mylittlepony
database = db356_waterlooid
```

Usage Breakdown

The client supports a variety of commands which can be broken down into multiple categories.

Once the environment is configured, you can run the client application `crime.py`. First navigate to the source directory:

```
cd src/
```

Once inside `src`, you can run the client application as follows:

```
python3 crime.py
```

This should display a list of the top-level commands. Please refer to the following sections for more information about how to use the client.

Ideal command line tool

Ideal list of features we want implemented:

1. Add information to the database
2. Update information from the database
3. Delete information from the database
4. Run a background check on someone
5. Able to query data based on location, code, and time
6. Allow the ability for the user to specify certain flags for information rather than prompting question
7. Allow functionality to sort data in ascending or descending order

8. Incorporate flags that allow user to use the “group by” feature in SQL
9. Command that allow you to gather statistics on the data ex(what fraction of suspects are found guilty in a given region for a particular type of crime)

The features we did **not** get to implement are: 6, 7, 8 and 9

For us the focus was really on ensuring that the core functionalities of a command line tool were available in the database. This included being able to add, update, delete and query data from the database. As we implemented the client, we had to decide regarding the ease of use of the application. While we do see that the ideal client application would have both flags (ex. Put -d to modify the date) and prompt the user for questions, we ended up implementing the prompting of information. This is because when it comes to our use cases a variety of different people would use the client such as law enforcement and government officials. So, ease of use was more important than the flexibility.

The additional features of being able to implement commands that allow you to group and sort data would be very beneficial. Especially from the standpoint of understanding the criminal activity better in a particular region. This can help bring to light many of the crimes taking place in the community and can maybe help you see what attributes or factors contribute to such crimes.

High level Commands

The client can be categorized as two sets of commands. The first are the high levels commands that allow you to do the following:

- adding information
- updating information
- deleting information
- showing / querying the database

To see the high level commands in the client type in the following :

```
python3 crime.py help
```

The following are some of the high level commands available in the client application. This can be seen from the command above.

Command	Description
<code>help</code>	Show this message
<code>create</code>	Create all tables
<code>load [num_of_entries]</code>	Load data from CSVs into tables (defaults to 1000 entries)
<code>clear</code>	Delete all entries in tables
<code>clean</code>	Drop all tables from database
<code>add</code>	Add entries to the database
<code>delete</code>	Delete entries from the database
<code>update</code>	Update entries in the database
<code>background</code>	Run background check on person
<code>show</code>	Show detailed record information
<code>filter</code>	Filter records based on location, date, and code

Example of running one of these commands in the terminal :

```
python3 crime.py add
```

Database Creation

Once the repository is configured, you can use the client to create the database. The client cannot interact with the database until it is created.

Start by creating the database tables, foreign keys, and views:

```
python3 crime.py create
```

Once the tables are created, load the data from the CSVs. This may take some time depending on the server load and how many records are being loaded:

```
python3 crime.py load <num_of_entries>
```

At this point, the database should be fully created and loaded with preliminary data from the CSVs. Note that some supplementary data is obtained from additional CSVs which are located under `src/codes` . These contain the crime codes for various crimes that are found in the database.

If you would like to delete all entries from the database tables, use the following `clear` command. This will **not** drop the tables from the database:

```
python3 crime.py clear
```

Alternatively, if you would like to remove all the tables from the database, use the `clean` command. This will drop the tables along with their entries, and is useful when you would like to rebuild the database from scratch or if you would like to reload the records from the CSVs:

```
python3 crime.py clean
```

Subcommands

`create`, `load`, `clear`, and `clean` are used for database creation and deletion only and do not have subcommands. The commands for the client can be further broken down to allow the user greater control over the database. For example the `add` command allows you to specify which tables you want to add the data to.



Note:

In the event that the client encounters an error, any modifications to the database will not be committed. The database will only be updated in the event of a successful command run.

To see all the specific commands type the command followed by `help` :

```
python3 crime.py add help
```

The following commands allow you to modify or query the information in the database. When the specific command is run the user will be prompted for a set of attributes through questions. The user can choose to set these attributes or leave them blank.

- `python3 crime.py add help`

Command	Description
<code>add code</code>	Add a crime code from a crime enforcement organization
<code>add location</code>	Add a location of a complaint, crime, or search
<code>add person</code>	Add a suspect or victim
<code>add complaint</code>	Add a complaint that was made to the police
<code>add crime</code>	Add a crimes that has taken place
<code>add search</code>	Add a suspect search record
<code>add help</code>	Show this message

- `python3 crime.py delete help`

Command	Description
<code>delete code <code> <organization></code>	Delete a code from a crime enforcement organization
<code>delete location <location_id></code>	Delete a location
<code>delete person <person_id></code>	Delete a person
<code>delete complaint <complaint_id></code>	Delete a police complaint record
<code>delete crime <crime_id></code>	Delete a crime record
<code>delete search <search_id></code>	Delete a suspect search record
<code>delete help</code>	Show this message

- `python3 crime.py update help`

Command	Description
<code>update code <code> <organization></code>	Update information about a crime code
<code>update location <location_id></code>	Update location details
<code>update person <person_id></code>	Update information about a person
<code>update complaint <complaint_id></code>	Update details of a police complaint
<code>update crime <crime_id></code>	Update details of a crime record
<code>update search <search_id></code>	Update details of a suspect search
<code>update help</code>	Show this message

- `python3 crime.py background help`

Command	Description
<code>background id</code>	Run background check based on person ID
<code>background name</code>	Run background check based on first and last names

- `python3 crime.py show help`

Command	Description
<code>show code</code>	Show information about an organization's crime code
<code>show person <person_id></code>	Show information about a person
<code>show location <location_id></code>	Show information about a location
<code>show complaint <complaint_id></code>	Show information about a complaint
<code>show search <search_id></code>	Show information about a search
<code>show crime <crime_id></code>	Show information about a crime
<code>show help</code>	Show this message

- `python3 crime.py filter help`

Command	Description
<code>filter complaint</code>	Filter complaint by location, date, and code
<code>filter crime</code>	Filter crimes by location, date, and code
<code>filter search</code>	Filter search by location, date, and code
<code>filter help</code>	Show this message

Walkthrough

Lets run through a few commands to get a better idea of how the command line tool works.

Background check

For the background check you can specify by which factors you want to lookup the person by:

```
python3 crime.py background
```

```
$ python3 crime.py background
```

```
> Search for person by:
```

```
> [1] ID
```

```
> [2] First & Last Names
```

```
> [3] <quit>
```

```
Enter selection: 1
```

```
Enter ID: 1
```

```
> Search returned 1 results
```

```
> Select person to run a background check on
```

```
> [1] Bob Jones
```

```
> [2] <quit>
```

```
Enter selection: 1
```

```
> -----
```

```
> Person information for Bob Jones:
```

```
> -----
```

```
> person_id: 1
```

```
> first_name: Bob
```

```
> last_name: Jones
```

```
> age_range: 56-57
```

```
> gender: Male
```

```
> ethnicity: White
```

```
> phone_number: (416) 456-2836
```

```
> -----
```

```
> -----
```

```
> Bob Jones was a victim of the following crimes:
```

```
> -----
```

```
> crime_id: 1898
```

```
> incident_id: 2098
```

```
> occurrence_date: 2017-01-01
```

```
> code: 110
```

```
> organization: IUCR
```

```
> category: HOMICIDE
```

```
> code_description: FIRST DEGREE MURDER
```

```
> type: sometupe
```

```
> status: something
```

```
> police_department: Chicago Police Department
```

```
> weapon: Knife
```

```
> domestic: 1
```

```
> description: Victim was attacked by a dangerous gang
```

```
> location_id: 1
```

```
> latitude: 51.45327300
```

```
> longitude: -0.00089000
```

```
> premises: None
> area: None
> precinct: None
> ward: None
> borough: None
> city: Chicago
> state: None
> country: United States
> victim_id: 1
> victim_first_name: Bob
> victim_last_name: Jones
> victim_age_range: 56-57
> victim_gender: Male
> victim_ethnicity: White
> victim_phone_number: (416) 456-2836
> -----
> crime_id: 1899
> incident_id: 2099
> occurrence_date: 2009-09-09
> code: 420
> organization: NYPD
> category: Drugs
> code_description: Intoxication
> type: None
> status: None
> police_department: None
> weapon: None
> domestic: None
> description: Victim was drugged
> location_id: 1
> latitude: 51.45327300
> longitude: -0.00089000
> premises: None
> area: None
> precinct: None
> ward: None
> borough: None
> city: New York
> state: None
> country: United States
> victim_id: 1
> victim_first_name: Bob
> victim_last_name: Jones
> victim_age_range: 56-57
> victim_gender: Male
> victim_ethnicity: White
> victim_phone_number: (416) 456-2836
> -----
> crime_id: 1900
> incident_id: 2100
> occurrence_date: 2003-05-21
> code: 114
> organization: NYPD
> category: ARSON
> code_description: None
> type: None
> status: Under Investigation
> police_department: New York Police Department
```

```
> weapon: Flamethrower
> domestic: 0
> description: Suspect torched victim\'s car with a flamethrower
> location_id: 5
> latitude: 51.47168500
> longitude: -0.13594000
> premises: None
> area: None
> precinct: None
> ward: None
> borough: None
> city: New York
> state: None
> country: United States
> victim_id: 1
> victim_first_name: Bob
> victim_last_name: Jones
> victim_age_range: 56-57
> victim_gender: Male
> victim_ethnicity: White
> victim_phone_number: (416) 456-2836
> -----
> crime_id: 1901
> incident_id: 2101
> occurrence_date: 2003-05-21
> code: 101
> organization: City of London Police
> category: ASSUALT
> code_description: None
> type: None
> status: Case Resolved
> police_department: City of London Police
> weapon: Handgun
> domestic: 0
> description: Victim was shot by suspect using a handgun
> location_id: 5
> latitude: 51.47168500
> longitude: -0.13594000
> premises: None
> area: None
> precinct: None
> ward: None
> borough: None
> city: London
> state: None
> country: United Kingdom
> victim_id: 1
> victim_first_name: Bob
> victim_last_name: Jones
> victim_age_range: 56-57
> victim_gender: Male
> victim_ethnicity: White
> victim_phone_number: (416) 456-2836
> -----
> -----
> Bob Jones was searched as a suspect of the following crimes:
> -----
```

```
> No stop & search records found
```

The background check can also be run with the `id` or `name` flags that allow you to search for a person based on id or name:

```
python3 crime.py background id
python3 crime.py background name
```

Adding a crime to the database

```
$ python3 crime.py add crime
> Do you know the ID of the victim?
[yes/no]: yes
[INT(10)] victim_id: 1
> Do you know the code and organization for this crime?
[yes/no]: yes
[VARCHAR(4)] code: 101
[VARCHAR(16)] organization: NYPD
> Do you know the location ID for this incident?
[yes/no]: yes
[INT(10)] location_id: 5
> Incident:
[DATE] occurrence_date: 03/05/21
[DATE] last_updated: 03/07/21
[VARCHAR(128)] status: Under Investigation
[VARCHAR(256)] police_department: New York Police Department
[VARCHAR(128)] type:
> Crime:
[VARCHAR(256)] weapon: Handgun
[TINYINT(3)] domestic: 0
[VARCHAR(256)] description: Victim was shot by suspect using a handgun
~ Added crime '1901' to database
```

Note

When adding a crime, the client will prompt you for the necessary information. If you don't know the information for a particular attribute, type `[Enter]` to skip it. This will insert a `NULL` for that field.

Using the SQL `SELECT` command, we can confirm that the crime record was successfully added to the database:

```
mysql> SELECT crime_id, weapon, description FROM Crime WHERE crime_id = 1901;
+-----+-----+-----+
| crime_id | weapon | description |
+-----+-----+-----+
| 1901 | Handgun | Victim was shot by suspect using a handgun |
+-----+-----+-----+
1 row in set (0.00 sec)
```

If you do not know some information that is mandatory, the command will abort, and you will be asked to separately add the corresponding records. For example, you cannot add a crime with a crime code and organization that does not yet exist in the database. In this case you will need to add a new crime code using the `python3 crime.py add code` command:

```
$ python3 crime.py add crime
> Do you know the ID of the victim?
[yes/no]: yes
[INT(10)] victim_id: 2
> Do you know the code and organization for this crime?
[yes/no]: no
> Please add a new code using 'python3 crime.py add code'
```

In other cases, the client will verify that your data is valid. For example, it will reject crime codes that do not exist in the database:

```
> Do you know the ID of the victim?
[yes/no]: yes
[INT(10)] victim_id: 3
> Do you know the code and organization for this crime?
[yes/no]: yes
[VARCHAR(4)] code: 1000
[VARCHAR(16)] organization: LAPD
X 'LAPD' code '1000' not found!
> Please add a new code using 'python3 crime.py add code'
```

Filter for crimes at a specific location

The `filter` command prompts for multiple sections that you can narrow your search by which include the following:

- Location
- Crime codes
- Date

You can specify one or multiple of these in the same filter query. The client will show you the most important information from the results. You can use the resulting codes and the `python3 crime.py show` command to learn more about these records based on their IDs.

Let's say that you wanted to know the crimes that took place in `Los Angeles` after `2020` and reported by the `LAPD` in the `Vehicle stolen` category.

```
python3 crime.py filter
> Do you want to filter by location?
[yes/no]: yes
> Do have a country you want to filter by?
[yes/no]: no
> Do have a city you want to filter by?
[yes/no]: yes
[VARCHAR(64)] city: Los Angeles
> Do have a state you want to filter by?
[yes/no]: no
> Do have a borough you want to filter by?
[yes/no]: no
> Do you want to filter by code?
[yes/no]: yes
> Do you know the crime code?
[yes/no]: yes
[VARCHAR(4)] code: 510
> Do you know the name of the crime enforcement organization?
[yes/no]: yes
[VARCHAR(16)] organization: LAPD
```



```
> Do you want to filter by date?
[yes/no]: yes
> [1] Filter by exact date
> [2] Filter by date range
Enter selection: 2
> Do have a max date you want to filter by?
[yes/no]: no
> Do have a min date you want to filter by?
[yes/no]: yes
[DATE] occurrence_date: 01/09/20
```

A few rows of the command output (some columns have been truncated for readability):

occurrence_date	code	city	victim_first_name	victim_last_name
2019-11-15	510	Los Angeles	Mary	Anzalone
2019-12-11	510	Los Angeles	Carl	Townsend
2019-10-15	510	Los Angeles	Virgil	Thissen

💡 Tip

If you want to look up a particular code to verify its entries before filtering by code, use the `show code` command as shown below.

```
$ python3 crime.py show code
> Do you know the crime code?
[yes/no]: yes
[VARCHAR(4)] code: 510
> Do you know the name of the crime enforcement organization?
[yes/no]: yes
[VARCHAR(16)] organization: LAPD
```

Code	Organization	Category	Description
510	LAPD	NULL	VEHICLE - STOLEN

Testing Plan

In order to test our client and database functionality, we made use of both manual and automatic testing methods. We individually tested each client command and verified the corresponding functionality using MySQL commands on the database. For example, the `add` and `delete` commands were manually verified by crosschecking for deletions and additions via SQL `SELECT` queries from the same tables, both before and after running the commands. Similar testing was conducted for `update`. These three commands make up the DDL and DML.

The remaining commands deal with querying only and not the SQL DDL/DML language such as `INSERT`, `DELETE`, and `UPDATE`. These include the `update`, `show`, `filter`, and `background` commands. These commands actually generate SQL `SELECT` queries in order to obtain the data, so it was easy to verify their output by constructing equivalent queries in the MySQL terminal.

In addition to manually testing our database functionality, we also developed automated test cases that accomplish the same goal. This test code runs various combinations of commands and subcommands and verifies the output by cross-referencing against the results that are directly queried from the database. The coverage for these unit tests are shown below:

```
$ coverage run tests.py -b
test_add_filter_complaint (__main__.TestCLI) ... ok
test_add_show_delete_crime (__main__.TestCLI) ... ok
test_add_show_delete_person (__main__.TestCLI) ... ok
test_background_check (__main__.TestCLI) ... ok
test_isNull (__main__.TestUtils) ... ok
test_isQuoted (__main__.TestUtils) ... ok
test_stripQuotes (__main__.TestUtils) ... ok
```

```
-----
Ran 7 tests in 1614.001s
```

OK

```
$ coverage report
```

Name	Stmts	Miss	Cover
MySQLUtils/MySQLUtils.py	27	3	89%
MySQLUtils/__init__.py	1	0	100%
crime.py	1252	815	35%
db.py	97	18	81%
log.py	19	3	84%
person.py	38	3	92%
tests.py	142	0	100%
transfer.py	418	196	53%
utils.py	101	7	93%
TOTAL	2095	1045	50%

Data Mining

The following sections detail the data mining investigation that we conducted on the crime dataset.

Goal

The goal of our data mining exercise was to answer to following question:

What factors influence the outcomes of stop-and-searches in the London area?

To evaluate this question, we made use of the data records found in `london-stop-and-search.csv`.

Feature Selection

As part of the data mining investigation, we decided to compare the results of London stop-and-searches based on two metrics:

- Suspect Ethnicity
- Suspect Gender

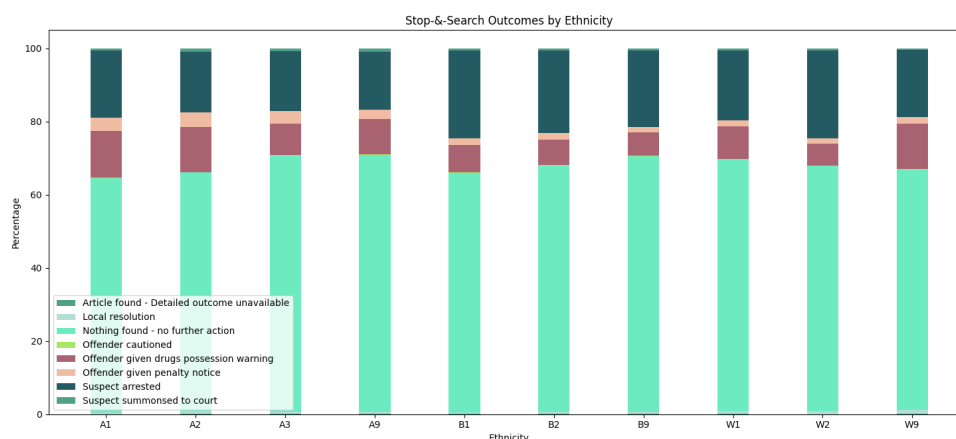
These metrics work well for this purpose as [london-stop-and-search.csv](#) has detailed information for both of these attributes. This allows us to draw stronger conclusions from the data when comparing outcomes with respect to ethnicity and gender. For references, these are the possible values for ethnicity that are found under the London Stop & Search records:

Ethnicity Code	Description
W1	White British
W2	White Irish
W9	Any other White ethnic background
B1	Black Caribbean
B2	Black African
B9	Any other Black ethnic background
A1	Asian Indian
A2	Asian Pakistani
A3	Asian Bangladeshi
A9	Any other Asian ethnic background

Results

The investigation produced results against suspect ethnicities and genders. The graphs show what percentage of people had each outcome, and the data is separated by gender and ethnicity.

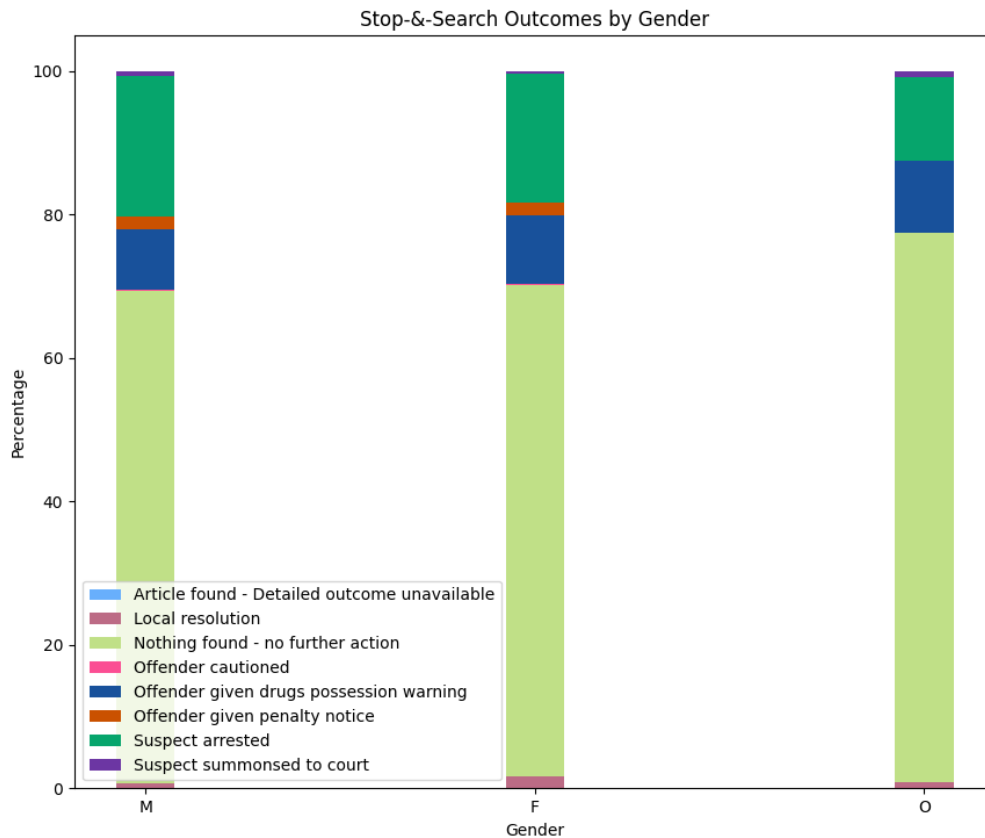
Ethnicity



The first observation is that the majority of suspects are discharged without any further action being taken. This is more or less evenly distributed across the ethnicities, which suggests that in general, stop and searches tend to result in no further action. A more interesting point of comparison is the percentage of suspects that were arrested after being searched. While not overly dramatic, it is apparent that a larger percentage of people from the B1, B2, B9, ethnic groups were arrested compared to those in the A1, A2, A3, and A9 groups. The W1, W2, and W9 ethnic groups have inconsistent results which span

from as low as **A** groups to as high as **B** groups. In simple terms, the results imply that black suspects are the most likely to be arrested, whereas south asian suspects are the least likely to be arrested. On the other hand, white suspects may or may not have a higher chance of being arrested, but the data does not provide a strong enough case to claim one or the other.

Gender



We first analyze the difference between male and female suspects, due to the more or less similar results in both those categories. Among males and females, both are equally likely to be discharged without any further action, as well as given penalties and drug warnings. However, it appears that male suspects are slightly more likely to be arrested than female suspects. Comparing across all three categories makes it quite apparent that those under the "other" category are more likely to get discharged without further action, and nobody under the other category was given a penalty. Note that the "other" category may also contain data for males and females whose gender was not recorded, so these findings might be slightly skewed with respect to the "other" category.

Conclusion

Challenges

There were quite few challenges that we encountered while building this project both from a command line interface and database design aspect.

Starting off with the database design, one of the really tricky aspects revolved around the two subsets of the US and UK database. One of our primary objectives was to develop a database that was generic. In other words it is a database that could be used in any crime system around the world. While there were many similarities between the two datasets, there are also a certain level of differences which makes it quite difficult to generalize.

We also faced challenges when developing our client command line tool. It was difficult to find a balance between usability and flexibility when designing the client commands, especially those which required prompting the user to enter values for attributes. By nature, the client program does not provide the full flexibility of SQL querying, so we needed to decide what the user should and should not be able to do with respect to modifying and querying the database.

Tradeoffs

Throughout this project we were forced to make tradeoffs both from a database design and command line interface perspective.

One of the big tradeoffs we made was in regards to the generalization of the database. We made it a goal to ensure that our crime database and system could be implemented in any crime system. Due to this you lose a certain level of granularity that exists within different regions throughout the world. While these region specific details could be important when it comes to crime, our goal was to ensure that we built a solid foundation for a crime database. Our understanding was that this database and system could later be adapted to a particular region at which point certain nuances could be added in.

Another tradeoff in our database design came from the relationship of linking people to multiple crimes. With the datasets given to us, each crime record corresponded to a single individual. Now in reality this certainly will not be the case as a single person can commit multiple crimes. That is why in reality it would be better to have a mapping between the `crime_id` and the `person_id` in a separate table. This would allow multiple people to correspond to a single crime and reduce the number of duplicates in the database.

In terms of the client tool development, we had to make a tradeoff between the flexibility and usability of the UI. We decided to use a combination of interactive prompting and command line arguments in order to implement various functionality. For example, prompting was used when adding or searching for attribute values, which would be very confusing for the user (and difficult to implement) if it were done via command line arguments, especially since some of the tables have many attributes, each with different formats. Having interactive prompts avoids confusing the user and allowed us to provide them with real-time feedback. To avoid too many levels of prompting, we separated core functions such as `add`, `delete`, `filter`, etc. into their own commands and subcommands, which preserves some granularity without becoming overwhelming for the user.

We also made the decision to restrict the user from attempting to micromanage sensitive attributes, such as primary keys. These attributes were abstracted away behind the user interface to make the database system easier to use. Therefore the user does not need to have an in depth knowledge of how the ER model works in order to use it. It also prevents the user from accidentally breaking constraints in the database, which would be impossible to fix if they did not have access to anything but the client tool. The overall goal was to provide the same core functionality of SQL via the client without the user having to actually write or think about the actual SQL code.

Future Improvements

When it comes to developing the client application one of the biggest challenges is preventing SQL injection. SQL injection is when the user can find vulnerable inputs and create malicious SQL commands to execute in the database. While this may not be a massive issue now considering we allow a certain subset of instructions in the CLI, it could be a massive issue as this application grows. One of the ways to prevent SQL injection is to do extensive input validation. While we do have a level of input validation based on the types of the attributes, we could certainly do more when it comes to adding content to the database. This is certainly something to improve on in the future as it would ensure greater security for this application. Furthermore aside from implementing more input validation in the future, testing for SQL injection should also be a priority. This would ensure that development efforts in the future don't introduce any new security leaks.

Another improvement to our system would be to add more rigorous error checking and unit tests for the client. Given the time constraints, there was only so much of the client that we could verify with regards to stability and robustness. In the future, it would be wise to more thoroughly test the various client commands using invalid data.

Additionally, it would be worthwhile to investigate the functionality of our database with additional datasets. The current datasets are only limited to certain areas and are not a comprehensive representation of all the types of crime records that exist around the world. Testing with additional datasets may reveal necessary design changes that can improve the performance, reliability, and functionality of the crime database design and implementation.

Commands to be implemented in the future

Due to the time restrictions we had to make a call on the features that we implemented for the command line interface. While we were able to implement the major features of add, update, delete and show in the database there were a few features that we simply didn't have the time to implement. A lot of these commands are more centered towards the statistics of the crime. For example being able to see run a command to see what fraction of crimes where the suspect is found guilty vs not guilty. Commands like this can paint a better picture of the criminal activity in an area over a specific period which could be crucial to law enforcement.