

FILP 00 FINAL REPORT

MINGXI WANG

ABSTRACT. In this report, I will talk about my work achievements, including the learning and practice of GIT uploading, downloading and remote connection, as well as the topic selection and programming practice of Kaggle

CONTENTS

1. Research motivation and context	2
1.1. Project Objectives	2
1.2. Project background	2
2. Research contents and methods	3
2.1. Evaluation metric	3
2.2. Data cleaning	4
2.3. Characteristics of the engineering	4
2.4. Information coding	6
2.5. Training model	6
3. Architecture of the model	7
4. Conclusion	8

Date: 2021-04-25.

2020 Mathematics Subject Classification. Artificial Intelligence.

1. RESEARCH MOTIVATION AND CONTEXT

1.1. Project Objectives. The project analyzed 12 years of crime reports from all of San Francisco's neighborhoods to create a model that can predict crime categories at a given time and place.

1.2. Project background. Crime is a kind of behavior that does great harm to the society. It brings great loss of life and property to human beings. In modern society, human beings prevent crime through strict judicial means, and research on crime prevention based on sociology, psychology and law.

The objective of this project is to quantitatively analyze a data set of nearly 12 years of crime reports from all San Francisco neighborhoods by computer means and to create a model that predicts the type of crime that will occur given information such as time and location.

2. RESEARCH CONTENTS AND METHODS

2.1. Evaluation metric.

- 878,049 samples, a total of 9 features.
- Date, category, description, day of the week, name of the police district, solution, approximate street address of the crime, longitude, latitude.

- First date: 2003-01-06 00:01:00
- Last date: 2015-05-13 23:53:00
- Test data shape (878049, 9)

```
print('First date: ', str(train.Dates.describe()['first']))
print('Last date: ', str(train.Dates.describe()['last']))
print('Test data shape ', train.shape)
```

FIGURE 1. preview I

	Dates	Category	Descript	DayOfWeek	PdDistrict	Resolution	Address	X	Y
0	2015-05-13 23:53:00	WARRANTS	WARRANT ARREST	Wednesday	NORTHERN	ARREST, BOOKED	OAK ST / LAGUNA ST	-122.425892	37.774599
1	2015-05-13 23:53:00	OTHER OFFENSES	TRAFFIC VIOLATION ARREST	Wednesday	NORTHERN	ARREST, BOOKED	OAK ST / LAGUNA ST	-122.425892	37.774599
2	2015-05-13 23:33:00	OTHER OFFENSES	TRAFFIC VIOLATION ARREST	Wednesday	NORTHERN	ARREST, BOOKED	VANNESS AV / GREENWICH ST	-122.424363	37.800414
3	2015-05-13 23:30:00	LARCENY/THEFT	GRAND THEFT FROM LOCKED AUTO	Wednesday	NORTHERN	NONE	1500 Block of LOMBARD ST	-122.426995	37.800873
4	2015-05-13 23:30:00	LARCENY/THEFT	GRAND THEFT FROM LOCKED AUTO	Wednesday	PARK	NONE	100 Block of BRODERICK ST	-122.438738	37.771541

FIGURE 2. preview II

2.2. Data cleaning.

- There are 2323 duplicate items that need to be removed

`train.duplicated().sum()`

- Throw away samples in a range of longitude and latitude (e.g., 50)

2.3. Characteristics of the engineering.

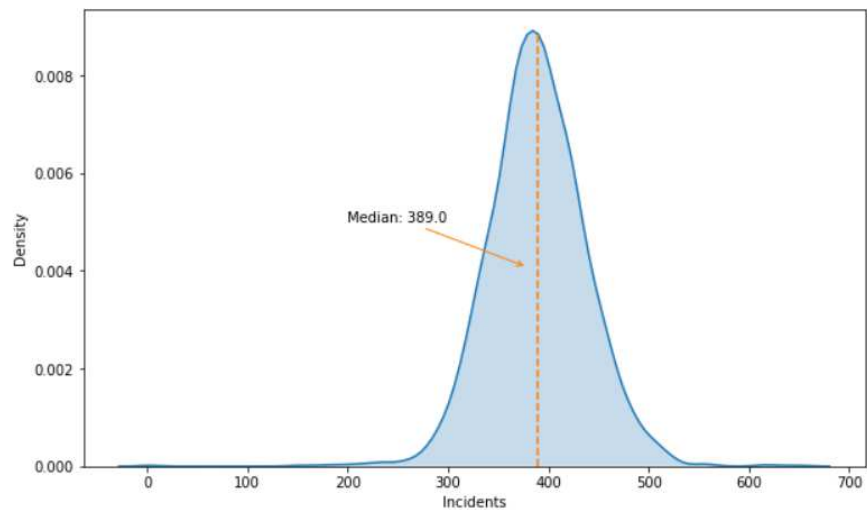


FIGURE 3. Distribution of number of incidents per day

Similarly, there was no significant deviation in the frequency of events over the course of the week.

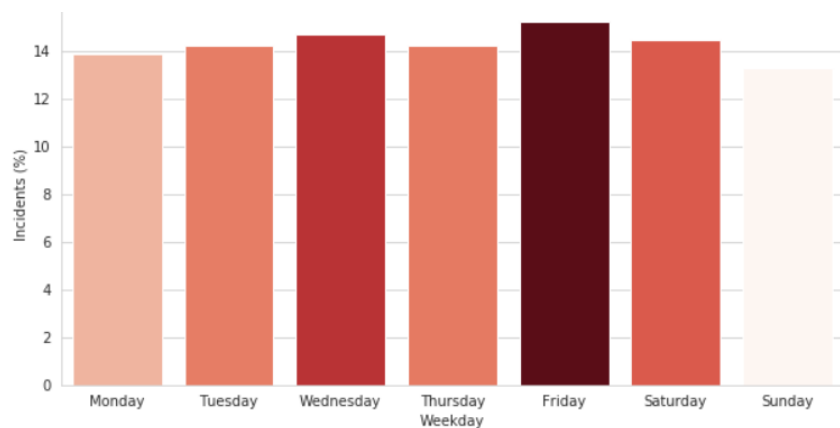


FIGURE 4. incidents per Weekday

A total of 39 discrete categories of incidents were recorded by police stations, the most common being theft (19.91%), non-criminal cases (10.50%) and assault (8.77%)

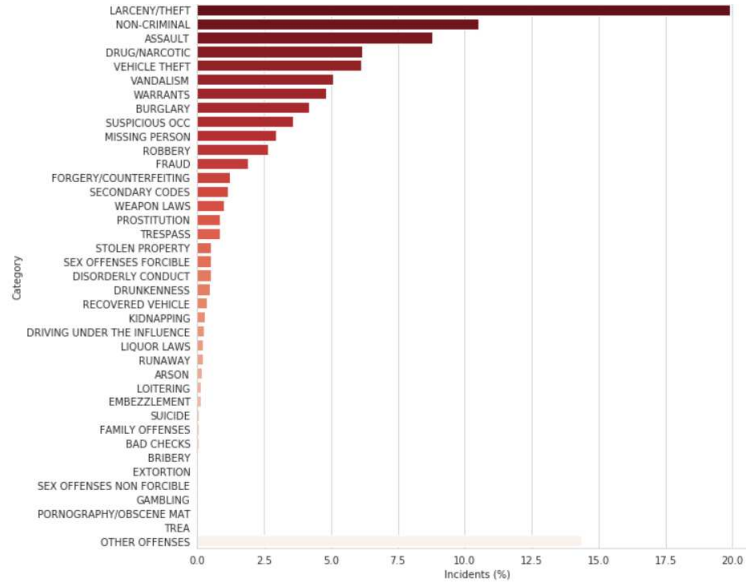


FIGURE 5. incidents per Crime Category

The chart below shows the average number of incidents per hour for the five crime categories. Obviously, different crimes occur with different frequencies at different times of the day. Prostitution, for example, takes place mostly at night, gambling incidents take place from late at night until morning, and burglaries from early morning until afternoon. As before, this is clear evidence that time parameters will also play an important role.

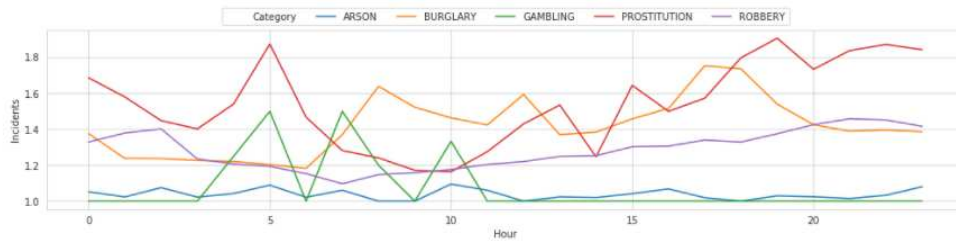


FIGURE 6. average number of incidents per hour

- From the "Date" field, we extracted the date, month, year, hours, minutes, days, and days since the first day.
- We extracted from the "Address" field whether the event occurred at the intersection or on a building component.

2.4. Information coding. The discrete data is converted to a number between 0 and N-1, where N is the number of different values in a list. You can think of it as the number of different values of a feature.

2.5. Training model.

- Train 23 Epochs.
- Use cross validation to evaluate the quality of our model.
After training, the model obtained:
- a cross validation score of 2.46.
- it got 2.49 on the test set.

3. ARCHITECTURE OF THE MODEL

The model architecture consists of six steps, as follows:

- 1. Make the decision tree fit the data
- 2. Evaluation model
- 3. Add weight to incorrect samples.
- 4. Select the leaf nodes with the greatest incremental loss for growth.
- 5. Generate the tree at the node in the previous step.
- 6. Go to Step 2 for the cycle

4. CONCLUSION

Conclusion: Based on data analysis, processing and utilization, our model can effectively predict crime types. By using the multi-classification algorithm to analyze the input features and update the model parameters through iterative training, our model can be trained more and more accurately. So as to complete the task of classification prediction.

(A. 1) SCHOOL OF COMPUTER SCIENCE,, JILIN UNIVERSITY, CHANGCHUN 130012, CHINA
Email address, A. 1: mxwang@tulip.academy