



Figure 1: *San Francisco 5* (2015)

Kaggle-Challenge: San Francisco Crime Classification

Nagaoka University of Technology

Date: January 10, 2019
Instructor: Professor Yukawa

Raphael Emberger

Contents

1. Preface	3
2. Abstract	4
3. Introduction	5
3.1. Initial situation	5
3.2. Objective	5
4. Theoretical Principles	7
4.1. Loss Function	7
5. Methods	8
5.1. Dataset	8
5.2. First Approach	11
5.2.1. Pre-Processing	11
5.2.2. Keras Model	11
5.2.3. Classification process	12
5.2.4. Results	13
5.3. Second Approach	14
5.3.1. Pre-processing	14
5.3.2. Results	16
5.4. Third approach	17
5.4.1. Results	17
6. Results	18
7. Conclusion	19
7.1. Future Considerations	19
8. Listings	20
A. Appendix	22

1. Preface

Firstly, I want to express my gratitude to Professor Yukawa for guiding me in this project and to the Kokusaika staff members to arrange my stay here at the Nagaoka University of Technology(subsequently referred to as "NUT"). I was given the generous opportunity to study at the NUT for one semester, for which I am very grateful. During that time I could choose from the following six Kaggle challenges to work on as project work:

- Toxic Comment Classification Challenge ([Kaggle 2017b](#))
- TalkingData AdTracking Fraud Detection Challenge ([Kaggle 2018b](#))
- Quora Question Pairs ([Kaggle 2017a](#))
- Expedia Hotel Recommendations ([Kaggle 2016](#))
- San Francisco Crime Classification ([Kaggle 2015](#))
- Inclusive Images Challenge ([Kaggle 2018a](#))

Of those, I was most interested in the classification of reported crimes ([Kaggle 2015](#)), as in my opinion this was an interesting challenge, given the dataset to be only consisting of time and spatial data. As such, this report is dedicated to take on this challenge.

2. Abstract

The first attempt to classify crimes based on date-time, district and coordinates was to build a neural network using [Keras \(n.d.\)](#). This approach failed by remaining on the same level as always guessing the most prominent label("Larceny/Theft") - 20%. This reached rank 1058 out of 2335.

The second attempt reached better results. This time, finished projects for the same challenge were used as reference to find problems with the first approach. This time, a Bernoulli Naïve Bayes classifier was used on the binarized dataset and reached a log-loss of 2.464 or 26.02%, which raised the rank up to 675.

The third and last attempt consisted of integrating the first attempt of using Keras into the second attempt. After some adjustment, the rank could be slightly improved once more: With a log-loss of 2.456(26.39%), rank 664 could be claimed.

The time for more improvement was not available because too much time was wasted on the first attempt, but the next steps would have included data enrichment and data manipulation.

3. Introduction

3.1. Initial situation

The challenge has been out since roughly 3 years and since then, many teams have participated and submitted their results. This lead the leader-board to fill up with 2335 submissions which were ranked and their results displayed online(see "Leaderboard" at [Kaggle \(2015\)](#)). The results vary from 34.53877 up to 1.95936, where the sample submission with a value of 32.89183 reaches rank 2241(see [4.1](#) for the ranking principle).

When searching on the internet for documents about that challenge, there are multiple such projects to be found. For example:

- A paper from [Darekar et al. \(2016\)](#). 2 Naïve Bayes, Decision Tree, Random Forest and Support Vector Machines classifiers were used. Reached highest accuracy of 23.16% with a Decision Tree.
- A blog post from [Ramunno-Johnson \(2015\)](#). In that project, a Bernoulli Naïves Bayes classifier was used. Reached a log-loss score of 2.58.
- A blog post from [Murray \(n.d.\)](#). AdaBoost, Bagging, Extra Trees, Gradient Boosting, K-Nearest Neighbors, Random Forest classifiers and Logistic Regression were used in this project. The dataset was enriched greatly by adding 9 other datasets (features like house prices, income, police and public transportation stations, healthcare center and homeless shelter locations, altitudes). Highest accuracy achieved with Gradient Boosted Trees, resulting in 45.7%.

3.2. Objective

The objective of this project is to produce a system that is capable of classifying the type of crime based off of the provided data consisting of date time stamps, the name of the district and street as well as the coordinates of the registered report. To quote [Kaggle \(2015\)](#):

From 1934 to 1963, San Francisco was infamous for housing some of the world's most notorious criminals on the inescapable island of Alcatraz.

Today, the city is known more for its tech scene than its criminal past. But, with rising wealth inequality, housing shortages, and a proliferation of expensive digital toys riding BART to work, there is no scarcity of crime in the city by the bay.

From Sunset to SOMA, and Marina to Excelsior, this competition's dataset provides nearly 12 years of crime reports from across all

of San Francisco's neighborhoods. Given time and location, you must predict the category of crime that occurred.

We're also encouraging you to explore the dataset visually. What can we learn about the city through visualizations like this Top Crimes Map? The top most up-voted scripts from this competition will receive official Kaggle swag as prizes.

Although the Kaggle challenge includes submitting an softmax array of the predictions of the test data, this objectives shifts towards self evaluation on the training set. The reason for this is that the challenge is already over and self evaluation was considered an easier approach to measure the success of the system.

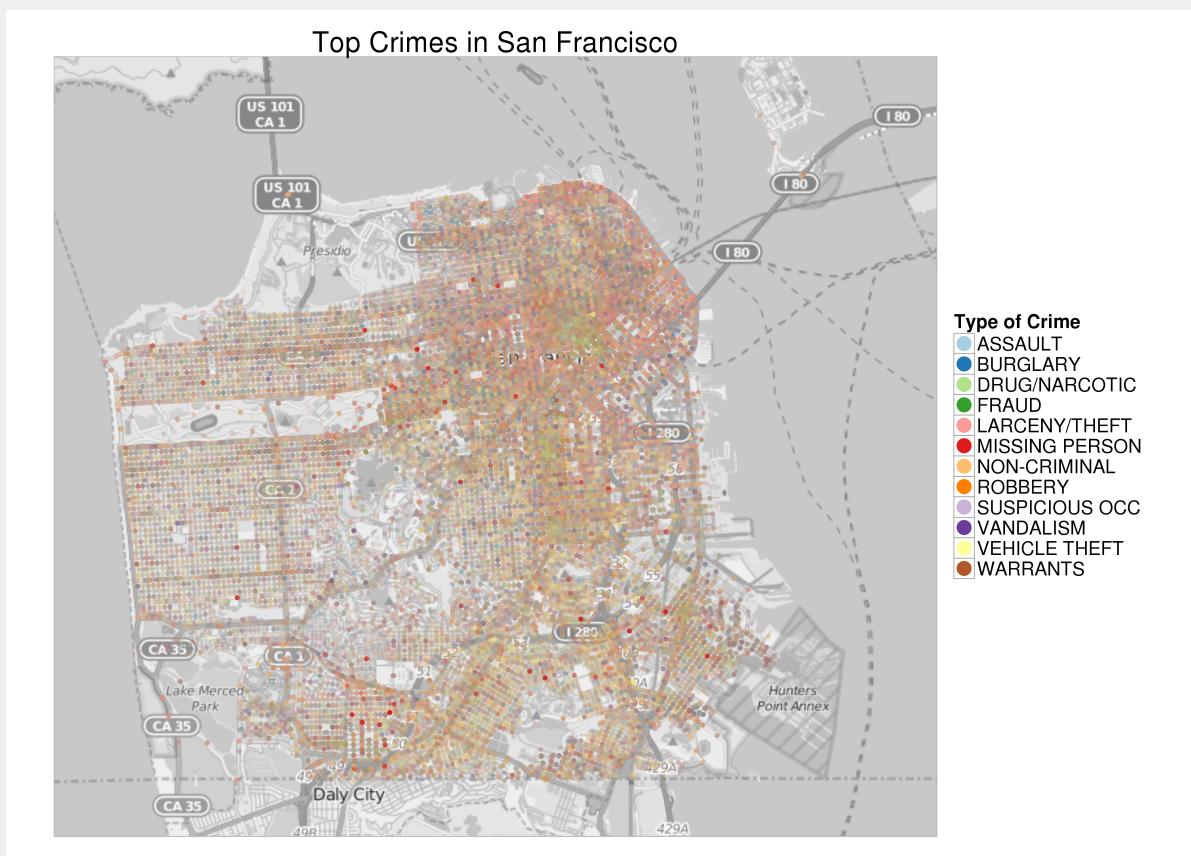


Figure 2: *Top Crimes in San Francisco* (2015)

4. Theoretical Principles

4.1. Loss Function

The ranking of the results on the Kaggle leader board([Kaggle 2015](#)) are based on the multi-class logarithmic loss function¹:

$$\text{loss} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M y_{ij} \log(p_{ij}) \quad (\text{I})$$

N : Number of cases in dataset.

M : Number of classes.

y_{ij} : Label for class. 1 if i is in j . Otherwise 0.

p_{ij} : Predicted probability that i belongs to j .

This basically boils down to a format as follows:

Class 1	Class 2	Class 3
0.24	0.48	0.38

With the labels being:

Class 1	Class 2	Class 3
0.00	1.00	0.00

When those values are applied to I, we get a value of 0.49548. Of course, the closer the prediction is to the actual labels, the smaller the loss value will be.

To calculate examples quickly on the python console, the following code can be used:

```
1 import numpy as np
2 from sklearn.metrics import log_loss
3 labels = np.array([0.0, 1.0, 0.0])
4 prediction = np.array([0.04, 0.78, 0.18])
5 print(log_loss(labels, prediction))
```

Listing 1: Quick Log Loss Calculation in python

¹The framework of [Keras \(n.d.\)](#) refers to the multi-class logarithmic loss function also as "categorical cross-entropy".

5. Methods

5.1. Dataset

The Kaggle challenge ([Kaggle 2015](#)) provides 3 files on their site under "Data":

- **sampleSubmission.csv**(884k x 40): A sample file, demonstrating the expected format for submissions to the challenge. Consists of an array of the softmax prediction of each sample in the **test.csv**.
- **test.csv**(884k x 7): The unlabeled sample subset of the data.
- **train.csv**(878k x 9): The labeled sample subset of the data.

The data itself consists of the gathered crime reports of the San Francisco Police Department from January 1st 2003 through May 13th 2015, where the odd weeks belong to the **test.csv** and the even weeks to **train.csv**.

Here are the 10 first rows of the respective data files:

Id	<i>one column per class</i>
0	<i>zeros, except the second last columns being all ones</i>
1	:
2	:
:	:

Table 1: sampleSubmission.csv(first 10 rows)

The two datasets differ slightly in their columns. The training dataset has added the labels(Category) but also Descript and Resolution, which will be ignored for this project.

Id	Dates	DayOfWeek	PdDistrict	Address	
				X	Y
0	2015-05-10 23:59:00	Sunday	BAYVIEW	2000 Block of THOMAS AV	
1	2015-05-10 23:51:00	Sunday	BAYVIEW	3RD ST / REVERE AV	
2	2015-05-10 23:50:00	Sunday	NORTHERN	2000 Block of GOUGH ST	
3	2015-05-10 23:45:00	Sunday	INGLESIDE	4700 Block of MISSION ST	
4	2015-05-10 23:45:00	Sunday	INGLESIDE	4700 Block of MISSION ST	
5	2015-05-10 23:40:00	Sunday	TARAVAL	BROAD ST / CAPITOL AV	
6	2015-05-10 23:30:00	Sunday	INGLESIDE	100 Block of CHENERY ST	
7	2015-05-10 23:30:00	Sunday	INGLESIDE	200 Block of BANKS ST	
8	2015-05-10 23:10:00	Sunday	MISSION	2900 Block of 16TH ST	
9	2015-05-10 23:10:00	Sunday	CENTRAL	TAYLOR ST / GREEN ST	

Table 2: test.csv(first 10 rows)

Dates	Category	Descript	DayOfWeek	PdDistrict
2015-05-13 23:53:00	WARRANTS	WARRANT ARREST	Wednesday	NORTHERN
2015-05-13 23:53:00	OTHER OFFENSES	TRAFFIC VIOLATION ARREST	Wednesday	NORTHERN
2015-05-13 23:33:00	OTHER OFFENSES	TRAFFIC VIOLATION ARREST	Wednesday	NORTHERN
2015-05-13 23:30:00	LARCENY/THEFT	GRAND THEFT FROM LOCKED AUTO	Wednesday	NORTHERN
2015-05-13 23:30:00	LARCENY/THEFT	GRAND THEFT FROM LOCKED AUTO	Wednesday	PARK
2015-05-13 23:30:00	LARCENY/THEFT	GRAND THEFT FROM UNLOCKED AUTO	Wednesday	INGLESIDE
2015-05-13 23:30:00	VEHICLE THEFT	STOLEN AUTOMOBILE	Wednesday	INGLESIDE
2015-05-13 23:30:00	VEHICLE THEFT	STOLEN AUTOMOBILE	Wednesday	BAYVIEW
2015-05-13 23:00:00	LARCENY/THEFT	GRAND THEFT FROM LOCKED AUTO	Wednesday	RICHMOND
2015-05-13 23:00:00	LARCENY/THEFT	GRAND THEFT FROM LOCKED AUTO	Wednesday	CENTRAL
Resolution	Address	X	Y	
"ARREST, BOOKED"	OAK ST / LAGUNA ST	-122.425891675136	37.7745985956747	
"ARREST, BOOKED"	OAK ST / LAGUNA ST	-122.425891675136	37.7745985956747	
"ARREST, BOOKED"	VANNESS AV / GREENWICH ST	-122.42436302145	37.8004143219856	
NONE	1500 Block of LOMBARD ST	-122.42699532676599	37.80087263276921	
NONE	100 Block of BRODERICK ST	-122.438737622757	37.771541172057795	
NONE	0 Block of TEDDY AV	-122.40325236121201	37.713430704116	
NONE	AVALON AV / PERU AV	-122.423326976668	37.7251380403778	
NONE	KIRKWOOD AV / DONAHUE ST	-122.371274317441	37.7275640719518	
NONE	600 Block of 47TH AV	-122.508194031117	37.776601260681204	
NONE	JEFFERSON ST / LEAVENWORTH ST	-122.419087676747	37.8078015516515	

Table 3: train.csv(first 10 rows)

ARSON	ASSAULT	BAD CHECKS
BRIBERY	BURGLARY	DISORDERLY CONDUCT
DRIVING UNDER THE INFLUENCE	DRUG/NARCOTIC	DRUNKENNESS
EMBEZZLEMENT	EXTORTION	FAMILY OFFENSES
FORGERY/COUNTERFEITING	FRAUD	GAMBLING
KIDNAPPING	LARCENY/THEFT	LIQUOR LAWS
LOITERING	MISSING PERSON	NON-CRIMINAL
OTHER OFFENSES	PORNOGRAPHY/OBSCENE MAT	PROSTITUTION
RECOVERED VEHICLE	ROBBERY	RUNAWAY
SECONDARY CODES	SEX OFFENSES FORCIBLE	SEX OFFENSES NON FORCIBLE
STOLEN PROPERTY	SUICIDE	SUSPICIOUS OCC
TREA	TRESPASS	VANDALISM
VEHICLE THEFT	WARRANTS	WEAPON LAWS

Table 4: Crime classes

The classes occur in an unbalanced matter: "Larceny/Theft" is the most predominant recorded crime, taking up nearly 19.92% of the dataset. For this reason, 19.92% is considered the bottom line of accuracy.

5.2. First Approach

For the first approach a **Keras** (n.d.) model on top of **Tensorflow** (n.d.) was chosen. For this, the first step was to pre-process the dataset to standardize it and properly feed it to the neural network.

5.2.1. Pre-Processing

To handle CSV files properly, a `CsvFile` class was created that represents a single csv file. When instantiated, it loads the csv file as a Pandas `DataFrame`. Apart from an abstract `def parse(self)`, it implements per-data-field functions that prepare the respective column for a conversion to a numerical representation(i.e. `def _prepare_date(self, date: datetime)`). It also defines an export function `def toNpArray(self) -> ndarray`, which allows to access the data as a `numpy array`.

From this basic class, three other classes were derived:

- `class TestDataCsvFile`
This class represents the `test.csv` file. It implements the missing `parse(self)` method.
- `class TrainDataCsvFile`
This class represents the sample part of the `train.csv` file. It implements the `parse(self)` method.
- `class TrainLabelsCsvFile`
This class represents the label part of the `train.csv` file. When instantiating, it can make a link to an already existing `TrainDataCsvFile` class, to prevent loading the same csv file a second time. It implements the `parse(self)` method.

All the `CsvFile` derived classes support loading their already preprocessed data from the disk. This was implemented to reduce calculation time when performing multiple runs.

Feature Scaling and Mean Normalization Apart from converting the data into integers or floats, the `_prepare_X` functions also perform feature scaling(usually to a $[-1, 1]$ set - sometimes $[0, 1]$) and mean normalization(i.e. the coordinates X and Y).

5.2.2. Keras Model

To build the model, a dedicated `Model` class was created. This class operates as a Keras model factory, using the `get_model` to either create and

train a model or load it's weights and parameters from a file and return the model.

The layers of the model changed greatly over time. This is the the last version of the model:

```
28     model = keras.Sequential([
29         keras.layers.Dense(16, input_shape=(train_data.shape[1],),
30         activation='relu'),
31         keras.layers.Dense(64, activation='relu'),
32         keras.layers.Dense(39, activation='softmax')
33     ])
34     self.log.info("Constructed model")
35     optimizer = keras.optimizers.Adam(lr=0.04)
36     model.compile(optimizer=optimizer,
37                     loss='sparse_categorical_crossentropy',
38                     metrics=['accuracy'])
39     self.log.info("Compiled model")
40
41     model.fit(train_data, train_labels, epochs=5, batch_size=200)
42     self.log.info("Trained model")
```

Listing 2: Keras model - model.py

5.2.3. Classification process

The classification process was written using the classes created in [5.2.1](#) and [2\(5.2.2\)](#):

```
34 trainsamplesfile = TrainDataCsvFile()
35 trainlabelfile = TrainLabelsCsvFile(trainsamplesfile)
36 testsamplesfile = TestDataCsvFile()
37 for file in [trainsamplesfile, trainlabelfile, testsamplesfile]:
38     if args.prep_data or not file.prep_file_exists():
39         file.parse()
40         file.save()
41     else:
42         file.load()
43
44 if args.train:
45     mdl = Model().get_model(
46         trainsamplesfile.toNpArray(),
47         trainlabelfile.toNpArray()
48     )
49 else:
50     mdl = Model().get_model()
51
52 predictions = mdl.predict(trainsamplesfile.toNpArray())
53 print("LogLoss: {}".format(log_loss(trainlabelfile.toNpArray(),
54         predictions)))
54 predicted_crime = np.argmax(predictions, axis=1)
55 print("Accuracy: {}%".format(accuracy_score(trainlabelfile.toNpArray(),
56         predicted_crime) * 100))
```

```
56  
57 for i in range(0, 19):  
58     predicted = trainlabelfile.CATEGORIES[np.argmax(predictions[i])]  
59     actual = trainlabelfile.get(i)  
60     logging.info("{} ?= {}".format(actual, predicted))
```

Listing 3: Classification process(first approach) - main.py

5.2.4. Results

When using this setup, various settings did not raise the accuracy in any way. The accuracy value usually hovered barely below 20%, which coincides with the bottom line described in [5.1](#) and could at best get to rank 1058. It was suspected that bad preprocessing of the dataset as well as the lack of knowledge and experience in neural networks was the cause for the lack of progress.

After various unfruitful tries, this approach had to be abandoned because of the lack in progress due to lack of knowledge and experience with Keras.

5.3. Second Approach

After an input from Prof. Yukawa, a new approach was established. This step consisted of using finished projects for this challenge as reference. For this, the project of [Ramunno-Johnson \(2015\)](#) was chosen.

The system now relies on a Bernoulli Naïve Bayes classifier, which in contrast to the first approach [5.2](#) does not require the system to build up layers to construct. It also uses a Logistic Regression algorithm to have a bar for comparison to be able to determine progress.

5.3.1. Pre-processing

In this approach, the data was processed differently. To prevent the algorithm to search for patterns according to the principle of locality and simplify the data by binarization, columns were split up into several different dummy columns, namely:

- Minutes(m0 - m59)
- Hours(H1 - H12)
- Days(D1 - D31)
- Months(M1 - M12)
- Years(Y2003 - Y2015)
- Weekdays(Monday - Sunday)
- Districts(Bayview - Tenderloin)

```
30 def preprocess_dataframe(data: pd.DataFrame) -> Tuple[pd.DataFrame, list]:  
31     print("Binarize data")  
32     minute =  
33         → pd.get_dummies(data.Dates.dt.minute).rename(columns=MINUTECOLUMNS)  
34     hour = pd.get_dummies(data.Dates.dt.hour).rename(columns=HOURCOLUMNS)  
35     day = pd.get_dummies(data.Dates.dt.day).rename(columns=DAYCOLUMNS)  
36     month =  
37         → pd.get_dummies(data.Dates.dt.month).rename(columns=MONTHCOLUMNS)  
38     year = pd.get_dummies(data.Dates.dt.year).rename(columns=YEARCOLUMNS)  
39     weekdays = pd.get_dummies(data.DayOfWeek)  
40     districts = pd.get_dummies(data.PdDistrict)  
41     x = data.X  
42     y = data.Y  
43     print("Assemble new array")  
44     new_data = pd.concat([minute, hour, day, month, year, weekdays,  
→ districts, x, y], axis=1)  
45     columns = new_data.keys().tolist()  
46     return new_data, columns
```

Listing 4: Pre-processing method - newmain.py

```

47 def evaluate(prediction, labels):
48     print("LogLoss: {}".format(log_loss(labels, prediction)))
49     predicted_crime = np.argmax(prediction, axis=1)
50     print("Accuracy: {}%".format(accuracy_score(labels, predicted_crime)
51         * 100))

```

Listing 5: Evaluation method - newmain.py

```

53 print("Load Data with pandas, and parse the first column into datetime")
54 train = pd.read_csv('train.csv', parse_dates=['Dates'])
55 test = pd.read_csv('test.csv', parse_dates=['Dates'])
56
57 print("Convert crime labels to numbers")
58 le_crime = preprocessing.LabelEncoder()
59 crime = le_crime.fit_transform(train.Category)
60
61 print("Build training data")
62 train_data, features = preprocess_dataframe(train)
63 train_data['crime'] = crime
64
65 print("Features[{}]: {}".format(len(features), np.array(features)))
66
67 print("Split up training data")
68 # training, validation = train_test_split(train_data, test_size=.20)
69 training = train_data
70 validation = train_data
71
72 # Bernoulli Naïve Bayes
73 print("Train Bernoulli Naïve Bayes classifier")
74 air_bnb = BernoulliNB()
75 air_bnb.fit(training[features], training['crime'])
76
77 print("Predict labels")
78 predicted = air_bnb.predict_proba(validation[features])
79
80 print("Validate prediction")
81 evaluate(predicted, validation['crime'])

```

Listing 6: Bernoulli Naïve Bayes fitting- newmain.py

```

114 lr = LogisticRegression(C=0.1, solver='lbfgs', multi_class='multinomial')
115 lr.fit(training[features], training['crime'])
116
117 print("Predict labels")
118 predicted = np.array(lr.predict_proba(validation[features]))
119
120 print("Validate prediction")
121 evaluate(predicted, validation['crime'])

```

Listing 7: Logistic Regression fitting - newmain.py

This lead to an expansion of the number of columns from 6 to 159, dropping column "Address" and binarizing all the other columns except for the coordinates, which remain floats.

5.3.2. Results

This new classification system reached a value of 2.464 with the Bernoulli Naïve Bayes classifier. This log-loss value corresponds to an accuracy of 26.02%. For comparison, the Logistic Regression reached 2.591 with an accuracy of 24.43%. This result would correspond to rank 675 on the Kaggle leaderboard.

To finalize this project, it was decided to combine the findings of those two approaches to further advance in this challenge.

5.4. Third approach

Although the second approach 5.3 was at an acceptable level, the failed first attempt was decided to be merged into the second one as comparison and to draw conclusions about what went wrong the first time. For this the classes written in the first approach were completely abandoned and only the core code of the first approach was amended and integrated into the script of 5.3.

5.4.1. Results

The amended Keras model showed results on the same level as the Bernoulli Naïves Bays classifier from 5.2.3. After a few tries, occasionally even slightly better results were reached:

```
95 model = keras.Sequential([
96     keras.layers.Dense(80, input_dim=len(features), activation='relu'),
97     keras.layers.Dense(118, activation='relu'),
98     keras.layers.Dense(39, activation='softmax')
99 ])
100 optimizer = keras.optimizers.Adam(lr=0.01)
101 model.compile(optimizer=optimizer,
102                 loss='sparse_categorical_crossentropy',
103                 metrics=['sparse_categorical_accuracy'])
104 model.fit(training[features], training['crime'], epochs=5,
105             batch_size=1024)
106
107 print("Predict labels")
108 predicted = model.predict_proba(validation[features])
109
110 print("Validate prediction")
111 evaluate(predicted, validation['crime'])
```

Listing 8: Keras model integrated - newmain.py

With this configuration the model reached a log-loss of 2.456 (accuracy: 26.39%), which was equivalent to rank 664. Which satisfies our definition of Done.

6. Results

Although the first approach lead to some measurable results(rank 1058), there were still many improvements to be made. Some of them were successfully implemented in the second approach and lead to an advance in rank to 675(log-loss: 2.464). After applying those findings into the first approach by merging both attempts, a further - although small - improvement could be made: The rank improved to 664(log-loss: 2.456).

7. Conclusion

During this project, a number of mistakes have been committed. The biggest one was investing too much time into the first attempt, instead of looking up finished projects for reference. However, this is only clear now in hindsight, as at the point in time, no disadvantage could be identified. Nevertheless did this loss in time cause the project to halt at rank 664(log-loss: 2.456), instead of advancing further due to further research.

Another problem was the lack of knowledge on the part of constructing neural networks with APIs like Keras. This lead to a number of unnecessary attempts on improvement on the wrong front-line. Much time was spent on implementing the conversion from data into integer representation, feature scaling and mean normalization, even though later in the process(see [5.3.1](#)) it turned out that binarizing the data was much faster in both programming and execution and brought better results.

7.1. Future Considerations

For future projects of similar nature, a better theoretical foundation is required. However, the mistakes made in this project were very insightful and helped in understanding general processes of creating neural networks.

If more time was available, the next step would have been enriching the dataset by relying on additional statistics and data from the internet. Other cited projects have relied on this method as well and just careless perusal makes this step seem very promising. The provided dataset by Kaggle after all only included date-time, district name and coordinates(address as well, but this was ignored in this project).

Another point to consider would be to process the data more: For example evening out distributions that otherwise are heavily uneven could lead to the neural network to not just abandon those samples in favour of sample types that are more prominent.

8. Listings

References

- Darekar, R., Dandona, R. & Sureshbabu, V. (2016), ‘Predicting and Analysis of Crime in San Francisco’.
URL: <https://www.slideshare.net/SameerDarekar1/san-francisco-crime-analysis-classification-kaggle-contest> 5
- Kaggle (2015), ‘San Francisco Crime Classification’.
URL: <https://www.kaggle.com/c/sf-crime> 3, 5, 7, 8
- Kaggle (2016), ‘Expedia Hotel Recommendations’.
URL: <https://www.kaggle.com/c/expedia-hotel-recommendations> 3
- Kaggle (2017a), ‘Quora Question Pairs’.
URL: <https://www.kaggle.com/c/quora-question-pairs> 3
- Kaggle (2017b), ‘Toxic Comment Classification Challenge’.
URL: <https://www.kaggle.com/c/jigsaw-toxic-comment-classification-challenge> 3
- Kaggle (2018a), ‘Inclusive Images Challenge’.
URL: <https://www.kaggle.com/c/inclusive-images-challenge> 3
- Kaggle (2018b), ‘TalkingData AdTracking Fraud Detection Challenge’.
URL: <https://www.kaggle.com/c/talkingdata-adtracking-fraud-detection> 3
- Keras (n.d.), ‘The Python Deep Learning library’.
URL: <https://keras.io/> 4, 7, 11
- Murray, M. (n.d.), ‘Classifying San Francisco Crime Incidents’.
URL: <http://mattmurray.net/classifying-san-francisco-crime-incidents/> 5
- Ramunno-Johnson, D. (2015), ‘Machine learning to predict San Francisco crime’.
URL: <http://efavdb.com/predicting-san-francisco-crimes/> 5, 14
- San Francisco 5 (2015).
URL: <https://allhdwallpapers.com/san-francisco-new-hd-wallpapers-high-resolution/> 1, 21

Tensorflow (n.d.), ‘An open source machine learning framework for everyone’.

URL: <https://www.tensorflow.org/> 11

Top Crimes in San Francisco (2015).

URL: <https://www.kaggle.com/benhamner/san-francisco-top-crimes-map> 6, 21

List of Figures

1. <i>San Francisco 5</i> (2015)	1
2. <i>Top Crimes in San Francisco</i> (2015)	6

List of Tables

1. sampleSubmission.csv(first 10 rows)	8
2. test.csv(first 10 rows)	9
3. train.csv(first 10 rows)	9
4. Crime classes	10

Listings

1. Quick Log Loss Calculation in python	7
2. Keras model - model.py	12
3. Classification process(first approach) - main.py	12
4. Pre-processing method - newmain.py	14
5. Evaluation method - newmain.py	15
6. Bernoulli Naïve Bayes fitting- newmain.py	15
7. Logistic Regression fitting - newmain.py	15
8. Keras model integrated - newmain.py	17
9. main.py	22
10. model.py	23
11. preprocessor.py	24
12. newmain.py	31

A. Appendix

```
1 #!/bin/env python3
2
3 import argparse
4 import logging
5
6 import tensorflow as tf
7 from sklearn.metrics import log_loss, accuracy_score
8 from tensorflow import keras
9
10 from model import Model
11 from preprocessor import *
12
13 logging.basicConfig(
14     format='%(asctime)s %(levelname)-8s %(name)20s: %(message)s',
15     datefmt='%Y-%m-%d %H:%M:%S',
16     level=logging.DEBUG
17 )
18
19 logging.info("Tensorflow: {}".format(tf.__version__))
20 logging.info("Keras: {}".format(keras.__version__))
21 logging.info("Numpy: {}".format(np.version.full_version))
22
23 parser = argparse.ArgumentParser(description="Classifies crimes of San
24     ↪ Francisco")
25 parser.add_argument('-t', '--train', action='store_true',
26                     help="Create and train the model(default: load model
27     ↪ from disk).")
28 parser.add_argument('-p', '--prep-data', action='store_true',
29                     help="Preprocess data files even if the preprocessed
30     ↪ data already exists.")
31 args = parser.parse_args()
32 if args.prep_data:
33     logging.debug("Preparing data from csv files")
34 if args.train:
35     logging.debug("Training model")
36
37 trainsamplesfile = TrainDataCsvFile()
38 trainlabelfile = TrainLabelsCsvFile(trainsamplesfile)
39 testsamplesfile = TestDataCsvFile()
40 for file in [trainsamplesfile, trainlabelfile, testsamplesfile]:
41     if args.prep_data or not file.prep_file_exists():
42         file.parse()
43         file.save()
44     else:
45         file.load()
46
47 if args.train:
48     mdl = Model().get_model(
49         trainsamplesfile.toNpArray(),
50         trainlabelfile.toNpArray()
```

```

48     )
49 else:
50     mdl = Model().get_model()
51
52 predictions = mdl.predict(trainsamplesfile.toNpArray())
53 print("LogLoss: {}".format(log_loss(trainlabelfile.toNpArray(),
54                                predictions)))
54 predicted_crime = np.argmax(predictions, axis=1)
55 print("Accuracy: {}%".format(accuracy_score(trainlabelfile.toNpArray(),
56                                predicted_crime) * 100))
57
58 for i in range(0, 19):
59     predicted = trainlabelfile.CATEGORIES[np.argmax(predictions[i])]
60     actual = trainlabelfile.get(i)
61     logging.info("{} == {} ".format(actual, predicted))
62
63 # for key in trainlabelfile.stats:
64 #     trainlabelfile.stats[key] /= trainlabelfile.count
65 #     print("{} : {}%".format(trainlabelfile.stats[key] * 100,
66 #                             trainlabelfile.CATEGORIES[key]))
67
68 # for i in range(0, 10):
69 #     # (date, _, _, address, lat, long) = data.get(i)
70 #     print(trainfile.get(i), trainlabelfile.get(i))
71 #     webbrowser.open(
72 #         "https://www.google.ch/maps/"
73 #         "@{}, {}, 58m/data=!3m1!1e3".format(lat, long)
74 #     )

```

Listing 9: main.py

```

1 import logging as log
2
3 # from tensorflow import keras
4 import keras
5 from numpy import ndarray
6
7
8 class Model:
9     file = "model.h5"
10    log = None
11
12    def __init__(self):
13        self.log = log.getLogger(self.__class__.__name__)
14
15    def _train(self, train_data, train_labels, test_data=None,
16              test_labels=None):
17        r"""Construct and train model
18
19        :param train_data: The training data
20        :type train_data: ndarray
21        :param train_labels: The training labels
22        :type train_labels: ndarray

```

```

22     :param test_data: The test data for verification
23     :type test_data: ndarray
24     :param test_labels: The test labels for verification
25     :type test_labels: ndarray
26     :return:
27     """
28     model = keras.Sequential([
29         keras.layers.Dense(16, input_shape=(train_data.shape[1],),
30     ↪ activation='relu'),
31         keras.layers.Dense(64, activation='relu'),
32         keras.layers.Dense(39, activation='softmax')
33     ])
34     self.log.info("Constructed model")
35     optimizer = keras.optimizers.Adam(lr=0.04)
36     model.compile(optimizer=optimizer,
37                     loss='sparse_categorical_crossentropy',
38                     metrics=['accuracy'])
39     self.log.info("Compiled model")
40
41     model.fit(train_data, train_labels, epochs=5, batch_size=200)
42     self.log.info("Trained model")
43
44     if test_data is not None and test_labels is not None:
45         test_loss, test_acc = model.evaluate(test_data, test_labels)
46         self.log.info("Tested model")
47         self.log.info("Test accuracy: {}".format(test_acc))
48
49     model.save(self.file)
50     self.log.info("Saved model")
51     return model
52
53     def get_model(self, train_data=None, train_labels=None,
54     ↪ test_data=None, test_labels=None):
55         if train_data is not None and train_labels is not None:
56             return self._train(train_data, train_labels, test_data,
57     ↪ test_labels)
58         else:
59             self.log.debug("Loading model")
60             return keras.models.load_model(self.file, compile=False)

```

Listing 10: model.py

```

1 #!/bin/env python3
2
3 import logging as log
4 import os
5 import sys
6 from datetime import datetime
7
8 import numpy as np
9 import pandas as pd
10 import pytz
11

```

```

12
13 class CsvFile:
14     filename = ''
15     df: pd.DataFrame = None
16     df_orig: pd.DataFrame = None
17     log = None
18     COLS = [
19         'Id',
20         'Dates',
21         'Year',
22         'Month',
23         'Day',
24         'Hour',
25         'Minute',
26         'Weekday',
27         'Season',
28         'Daynight',
29         'DayOfWeek',
30         'PdDistrict',
31         'Address',
32         'X',
33         'Y'
34     ]
35     DAYS = ['Monday', 'Tuesday', 'Wednesday', 'Thursday', 'Friday',
36             ↪ 'Saturday', 'Sunday']
37     DISTRICTS = ['BAYVIEW', 'NORTHERN', 'INGLESIDE', 'TARAVAL',
38             ↪ 'MISSION', 'CENTRAL', 'TENDERLOIN', 'RICHMOND',
39             ↪ 'SOUTHERN', 'PARK']
40     min_date = pd.Timestamp('1/1/2003 00:00:00')
41     max_date = pd.Timestamp('1/1/2016 00:00:00')
42     sf_tz = pytz.timezone('US/Pacific')
43
44     def __init__(self, filename, csvfile=None):
45         self.log = log.getLogger(self.__class__.__name__)
46         self.filename = filename
47         file = self.filename + '.csv'
48         if csvfile is not None:
49             self.df_orig = csvfile.df_orig
50             self.log.info("Linked to data frame of '{}' with shape"
51             ↪ "{}".format(csvfile.filename, self.df_orig.shape))
52         else:
53             self.log.debug("Reading csv file '{}'".format(file))
54             self.df_orig = self._read_file(file)
55             self.log.info("Read data frame of shape"
56             ↪ "{}".format(self.df_orig.shape))
57
58     def _read_file(self, file):
59         raise NotImplementedError()
60
61     def _read_prep_file(self, file):
62         return pd.read_csv(
63             file,

```

```

60         index_col='Id',
61         keep_date_col=True,
62         parse_dates=[ "Dates" ]
63     )
64
65     def parse(self):
66         raise NotImplemented()
67
68     def _prepare_date(self, date: datetime):
69         # date = datetime.strptime(date, '%Y-%m-%d %H:%M:%S')
70         # date = date.astimezone(self.sf_tz)
71         max_delta = self.max_date - self.min_date
72         delta = date - self.min_date
73         return int(delta.total_seconds()) / max_delta.total_seconds()
74
75     def _transform_date(self):
76         def get_season(month: int):
77             if month < 3 or month == 12:
78                 return 0 # December - February: Winter
79             if month < 6:
80                 return 1 # March - May: Spring
81             if month < 9:
82                 return 2 # June - August: Summer
83             if month < 12:
84                 return 3 # September - November: Autumn
85             raise Exception()
86
87         def get_day_night(hour: int):
88             if hour <= 7 or hour >= 7:
89                 return 0 # Night
90             return 1 # Day
91
92         self.log.debug("Parsing Year")
93         self.df['Year'] = self.df['Dates'].apply(lambda x: x.year)
94         self.log.debug("Parsing Month")
95         self.df['Month'] = self.df['Dates'].apply(lambda x: x.month)
96         self.log.debug("Parsing Day")
97         self.df['Day'] = self.df['Dates'].apply(lambda x: x.day)
98         self.log.debug("Parsing Hour")
99         self.df['Hour'] = self.df['Dates'].apply(lambda x: x.hour)
100        self.log.debug("Parsing Minute")
101        self.df['Minute'] = self.df['Dates'].apply(lambda x: x.minute)
102        self.log.debug("Parsing Weekday")
103        self.df['Weekday'] = self.df['Dates'].apply(lambda x:
104            ↪ x.isoweekday())
105        self.log.debug("Parsing Season")
106        self.df['Season'] = self.df['Dates'].apply(lambda x:
107            ↪ get_season(x.month))
108        self.log.debug("Parsing Daynight")
109        self.df['Daynight'] = self.df['Dates'].apply(lambda x:
110            ↪ get_day_night(x.day))
111        self.log.debug("Parsing Dates")

```

```

109     self.df['Dates'] = self.df['Dates'].apply(self._prepare_date)
110
111     def _prepare_day(self, daystr):
112         return self.DAYS.index(daystr)/(len(self.DAYS)/2)-1
113
114     def _prepare_district(self, daystr):
115         return self.DISTRICTS.index(daystr)/(len(self.DISTRICTS)/2)-1
116
117     @staticmethod
118     def _prepare_address(addressstr):
119         return hash(addressstr)/(sys.maxsize/2)-1
120
121     @staticmethod
122     def _prepare_latitude(latitudestr):
123         return float(latitudestr)/180
124
125     @staticmethod
126     def _prepare_longitude(longitudestr):
127         return float(longitudestr)/180
128
129     def save(self):
130         if self.df is None:
131             self.log.error("Data not yet parsed")
132             return
133         file = self._prep_file()
134         self.log.debug("Writing csv file '{}'".format(file))
135         self.df.to_csv(
136             file,
137             index_label='Id'
138         )
139         self.log.info("Wrote data frame of shape
140             → {}".format(self.df.shape))
141
142     def _prep_file(self):
143         return self.filename + '_prep.csv'
144
145     def prep_file_exists(self):
146         return os.path.isfile(self._prep_file())
147
148     def load(self):
149         file = self._prep_file()
150         self.log.debug("Reading csv file '{}'".format(file))
151         self.df = self._read_prep_file(file)
152         self.log.info("Read data frame of shape
153             → {}".format(self.df_orig.shape))
154
155     def get(self, index):
156         if self.df is None:
157             self.log.error("Data not yet parsed")
158             return
159         date = self.df_orig['Dates'][index]
160         # date = datetime.strptime(self.df_old.at[index], 'Dates'],

```

```

159     ↳ '%Y-%m-%d %H:%M:%S')
160     day = self.df_orig['DayOfWeek'][index]
161     # day = self.df_old.at[index, 'DayOfWeek']
162     district = self.df_orig['PdDistrict'][index]
163     # district = self.df_old.at[index, 'DayOfWeek']
164     address = self.df_orig['Address'][index]
165     latitude = float(self.df_orig['Y'][index])
166     longitude = float(self.df_orig['X'][index])
167     return date, day, district, address, latitude, longitude
168
169 def toNpArray(self):
170     return np.reshape(self.df.values, self.df.shape)
171
172 class TestDataCsvFile(CsvFile):
173
174     def __init__(self):
175         super().__init__("test")
176
177     def _prep_file(self):
178         return self.filename + '_samples.csv'
179
180     def _read_file(self, file):
181         return pd.read_csv(
182             file,
183             index_col='Id',
184             keep_date_col=True,
185             parse_dates=["Dates"]
186         )
187
188     def parse(self):
189         self.df = self.df_orig.copy()
190         self.log.debug('Parsing Dates')
191         self._transform_date()
192         self.log.debug('Parsing Day of the week')
193         self.df['DayOfWeek'] =
194             ↳ self.df['DayOfWeek'].apply(self._prepare_day)
195             self.log.debug('Parsing District')
196             self.df['PdDistrict'] =
197             ↳ self.df['PdDistrict'].apply(self._prepare_district)
198             self.log.debug('Parsing Address')
199             self.df['Address'] =
200             ↳ self.df['Address'].apply(self._prepare_address)
201             self.log.debug('Parsing Longitude')
202             self.df['X'] = self.df['X'].apply(self._prepare_longitude)
203             self.log.debug('Parsing Latitude')
204             self.df['Y'] = self.df['Y'].apply(self._prepare_latitude)
205             self.log.info('Parsed dataframe')
206

```

```

207     def __init__(self, csvfile=None):
208         super().__init__("train", csvfile)
209
210     def _prep_file(self):
211         return self.filename + '_samples.csv'
212
213     def _read_file(self, file):
214         return pd.read_csv(
215             file,
216             keep_date_col=True,
217             parse_dates=["Dates"]
218         )
219
220     def parse(self):
221         self.df = self.df_orig.copy()
222         # print(self.df.shape)
223         # print(self.df['Dates'].shape)
224         # print(self.df[['Dates']].shape)
225         # print(self.df.at[1, 'Dates'])
226         # print(self.df.at[1, 'Dates'].minute)
227         # exit(0)
228         self.log.debug('Parsing Dates')
229         self._transform_date()
230         self.log.debug('Deleting Category')
231         self.df = self.df.drop('Category', axis=1)
232         # self.log.debug('Parsing Category')
233         # self.df['Category'] =
234         ↳ self.df['Category'].apply(self._prepare_category)
235             self.log.debug('Deleting Descript')
236             self.df = self.df.drop('Descript', axis=1)
237             self.log.debug('Parsing Day of the week')
238             self.df['DayOfWeek'] =
239             ↳ self.df['DayOfWeek'].apply(self._prepare_day)
240                 self.log.debug('Parsing District')
241                 self.df['PdDistrict'] =
242                 ↳ self.df['PdDistrict'].apply(self._prepare_district)
243                     self.log.debug('Deleting Resolution')
244                     self.df = self.df.drop('Resolution', axis=1)
245                     self.log.debug('Parsing Address')
246                     self.df['Address'] =
247                     ↳ self.df['Address'].apply(self._prepare_address)
248                         self.log.debug('Parsing Longitude')
249                         self.df['X'] = self.df['X'].apply(self._prepare_longitude)
250                         self.log.debug('Parsing Latitude')
251                         self.df['Y'] = self.df['Y'].apply(self._prepare_latitude)
252                         self.log.info('Parsed dataframe')
253
254
255 class TrainLabelsCsvFile(CsvFile):
256     CATEGORIES = ['ARSON', 'ASSAULT', 'BAD CHECKS', 'BRIBERY',
257     ↳ 'BURGLARY', 'DISORDERLY CONDUCT',
258                 'DRIVING UNDER THE INFLUENCE', 'DRUG/NARCOTIC',

```

```

254     ↳ 'DRUNKENNESS', 'EMBEZZLEMENT', 'EXTORTION',
255         'FAMILY OFFENSES', 'FORGERY/COUNTERFEITING', 'FRAUD',
256     ↳ 'GAMBLING', 'KIDNAPPING', 'LARCENY/THEFT',
257         'LIQUOR LAWS', 'LOITERING', 'MISSING PERSON',
258     ↳ 'NON-CRIMINAL', 'OTHER OFFENSES',
259         'PORNOGRAPHY/OBSCENE MAT', 'PROSTITUTION', 'RECOVERED
260     ↳ VEHICLE', 'ROBBERY', 'RUNAWAY',
261         'SECONDARY CODES', 'SEX OFFENSES FORCIBLE', 'SEX
262     ↳ OFFENSES NON FORCIBLE', 'STOLEN PROPERTY', 'SUICIDE',
263         'SUSPICIOUS OCC', 'TREA', 'TRESPASS', 'VANDALISM',
264     ↳ 'VEHICLE THEFT', 'WARRANTS', 'WEAPON LAWS']
265     stats = {}
266     count = 0
267
268     def __init__(self, csvfile=None):
269         super().__init__("train", csvfile)
270
271     def _prep_file(self):
272         return self.filename + '_labels.csv'
273
274     def _read_file(self, file):
275         return pd.read_csv(file)
276
277     def _read_prep_file(self, file):
278         return pd.read_csv(
279             file,
280             index_col='Id'
281         )
282
283     def _prepare_category(self, category):
284         cat = self.CATEGORIES.index(category)
285         if not cat in self.stats:
286             self.stats[cat] = 0.0
287         self.stats[cat] += 1.0
288         self.count += 1
289         return cat
290
291     def parse(self):
292         self.df = self.df_orig.copy()
293         # self.log.debug('Deleting Id')
294         # self.df = self.df.drop('Id', axis=1)
295         self.log.debug('Deleting Dates')
296         self.df = self.df.drop('Dates', axis=1)
297         self.log.debug('Parsing Category')
298         self.df['Category'] =
299             ↳ self.df['Category'].apply(self._prepare_category)
300             self.log.debug('Deleting Descript')
301             self.df = self.df.drop('Descript', axis=1)
302             self.log.debug('Deleting DayOfWeek')
303             self.df = self.df.drop('DayOfWeek', axis=1)
304             self.log.debug('Deleting PdDistrict')
305             self.df = self.df.drop('PdDistrict', axis=1)

```

```

299     self.log.debug('Deleting Resolution')
300     self.df = self.df.drop('Resolution', axis=1)
301     self.log.debug('Deleting Address')
302     self.df = self.df.drop('Address', axis=1)
303     self.log.debug('Deleting Longitude')
304     self.df = self.df.drop('X', axis=1)
305     self.log.debug('Deleting Latitude')
306     self.df = self.df.drop('Y', axis=1)
307     self.log.debug('Parsed dataframe')
308
309     def get(self, index):
310         if self.df is None:
311             self.log.error("Data not yet parsed")
312             return
313         district = self.CATEGORIES[self.df.at[index, 'Category']]
314         return district

```

Listing 11: preprocessor.py

```

1 #!/bin/env python3
2
3 from typing import Tuple
4
5 import keras
6 import numpy as np
7 import pandas as pd
8 from sklearn import preprocessing
9 from sklearn.linear_model import LogisticRegression
10 from sklearn.metrics import log_loss, accuracy_score
11 from sklearn.naive_bayes import BernoulliNB
12
13 MINUTECOLUMNS = {}
14 for min_int in range(0, 60):
15     MINUTECOLUMNS[min_int] = "m{}".format(min_int)
16 HOURCOLUMNS = {}
17 for hour_int in range(0, 24):
18     HOURCOLUMNS[hour_int] = "H{}".format(hour_int)
19 DAYCOLUMNS = {}
20 for day_int in range(0, 30):
21     DAYCOLUMNS[day_int] = "D{}".format(day_int + 1)
22 MONTHCOLUMNS = {}
23 for month_int in range(0, 11):
24     MONTHCOLUMNS[month_int] = "M{}".format(month_int + 1)
25 YEARCOLUMNS = {}
26 for year_int in range(2003, 2015):
27     YEARCOLUMNS[year_int] = "Y{}".format(year_int)
28
29
30 def preprocess_dataframe(data: pd.DataFrame) -> Tuple[pd.DataFrame, list]:
31     print("Binarize data")
32     minute =
33         ↪ pd.get_dummies(data.Dates.dt.minute).rename(columns=MINUTECOLUMNS)
34     hour = pd.get_dummies(data.Dates.dt.hour).rename(columns=HOURCOLUMNS)

```

```

34     day = pd.get_dummies(data.Dates.dt.day).rename(columns=DAYCOLUMNS)
35     month =
36         ↪ pd.get_dummies(data.Dates.dt.month).rename(columns=MONTHCOLUMNS)
37     year = pd.get_dummies(data.Dates.dt.year).rename(columns=YEARCOLUMNS)
38     weekdays = pd.get_dummies(data.DayOfWeek)
39     districts = pd.get_dummies(data.PdDistrict)
40     x = data.X
41     y = data.Y
42     print("Assemble new array")
43     new_data = pd.concat([minute, hour, day, month, year, weekdays,
44         ↪ districts, x, y], axis=1)
45     columns = new_data.keys().tolist()
46     return new_data, columns
47
48
49
50
51
52
53 print("Load Data with pandas, and parse the first column into datetime")
54 train = pd.read_csv('train.csv', parse_dates=['Dates'])
55 test = pd.read_csv('test.csv', parse_dates=['Dates'])
56
57 print("Convert crime labels to numbers")
58 le_crime = preprocessing.LabelEncoder()
59 crime = le_crime.fit_transform(train.Category)
60
61 print("Build training data")
62 train_data, features = preprocess_dataframe(train)
63 train_data['crime'] = crime
64
65 print("Features[{}]: {}".format(len(features), np.array(features)))
66
67 print("Split up training data")
68 # training, validation = train_test_split(train_data, test_size=.20)
69 training = train_data
70 validation = train_data
71
72 # Bernoulli Naïve Bayes
73 print("Train Bernoulli Naïve Bayes classifier")
74 air_bnb = BernoulliNB()
75 air_bnb.fit(training[features], training['crime'])
76
77 print("Predict labels")
78 predicted = air_bnb.predict_proba(validation[features])
79
80 print("Validate prediction")
81 evaluate(predicted, validation['crime'])
82

```

```

83 # Predict crimes of test dataset
84 print("Build test data")
85 test_data, _ = preprocess_dataframe(test)
86
87 print("Predict test labels")
88 predicted = air_bnb.predict_proba(test_data[features])
89
90 print("Write results")
91 result = pd.DataFrame(predicted, columns=le_crime.classes_)
92 result.to_csv('testResult.csv', index=True, index_label='Id')
93
94 print("Train Keras")
95 model = keras.Sequential([
96     keras.layers.Dense(80, input_dim=len(features), activation='relu'),
97     keras.layers.Dense(118, activation='relu'),
98     keras.layers.Dense(39, activation='softmax')
99 ])
100 optimizer = keras.optimizers.Adam(lr=0.01)
101 model.compile(optimizer=optimizer,
102                 loss='sparse_categorical_crossentropy',
103                 metrics=['sparse_categorical_accuracy'])
104 model.fit(training[features], training['crime'], epochs=5,
105             batch_size=1024)
106
107 print("Predict labels")
108 predicted = model.predict_proba(validation[features])
109
110 print("Validate prediction")
111 evaluate(predicted, validation['crime'])
112
113 # Logistic Regression
114 print("Train Logistic Regression for comparison")
115 lr = LogisticRegression(C=0.1, solver='lbfgs', multi_class='multinomial')
116 lr.fit(training[features], training['crime'])
117
118 print("Predict labels")
119 predicted = np.array(lr.predict_proba(validation[features]))
120
121 print("Validate prediction")
122 evaluate(predicted, validation['crime'])

```

Listing 12: newmain.py