

Crime Rate Prediction Using Machine Learning

Sirivanth Paladugu^{1*}, Tarun Sai Yakkala², Neeraj Boggarapu³, Sri Krishna Kumar Modekurty⁴

^{1,2,3,4}Student, Department of Computer Science and Engineering, SRM University, Guntur, India

Abstract: The main objective of our project is to predict the Crime Rate in different regions using specific parameters like density, country, crime rate, central etc. In this project we considered a dataset of a particular Country. The region we selected is the United States of America. We used the Linear Regression model to predict the crime rate. Finally, a graph is plotted after implementing Linear Regression. A graph is plotted between density and crime rate to enhance distribution of crime rate in a particular region.

Keywords: Data mining, data preprocessing, linear regression, predictive analysis.

1. Introduction

Analysis of crime is a methodological approach to the identification and assessment of criminal patterns and trends. Before starting this project we have gone through many datasets all over the world, but in the USA a state named North Carolina has given us sufficient data to start off with. The variables with highest correlation in regards with crime rate are urban and density. The reason for highest correlation for urban and density variables is that the urban areas of North Carolina are densely populated. As a result, there is a high probability of multicollinearity between the density and urban features. The combination of density and location features will help us in predicting the crime rate in North Carolina. The columns of wage, wfed and wtrd are positively correlated with the density feature. From this we can instinctively understand that the weekly wages would be higher in urban localities. The features wtrd and wfir are positively correlated with wfed and wloc respectively. Both the wfir and wtrd are moderately correlated. The features crime rate, density, urban, wfed, wtrd and taxpc are very highly correlated with the Crime Rate. The prediction of future crime location can be determined by geographical data mining approaches

2. Problem Survey

In our problem we will be evaluating and examining the large pre-existing databases in order to generate new information which would help us to find the solution. The prediction is based on the extraction of the new information using the existing datasets. The main aim of this problem is to perform the survey on certain algorithms which helps us to analyze the crime rate.

3. Dataset Description

In our dataset we have considered several parameters like country, year, crime rate, density etc. After once we get the raw dataset we then clean the dataset. After we clean the dataset we then preprocess it. We then use linear regression to obtain our results, the results are obtained in the form of graphs. We then use a statistical approach to find values like count, mean, standard deviation, minimum and maximum values.

For each and every parameter we have given a model summary. A combination of density and location (west/ central/ urban) can help aid crime rate prediction. The variables with highest correlation in regards with crime rate are urban and density.

4. Data Preprocessing

The important step in the entire process is data preprocessing. Data preprocessing allows us to remove the unwanted data with the help of data cleaning, this allows the user to have a dataset which contains more valuable information after the preprocessing stage in the mining process.

Data cleaning:

```
1 crimeData = crimeData[crimeData.county != 185]
2 crimeData = crimeData[crimeData.county != 115]
3 crimeData = crimeData[crimeData['prbarr'] < 1]
4 crimeData = crimeData[crimeData['prbconv'] < 1]
5 crimeData = crimeData[crimeData['west']+crimeData['central'] <= 1]
6 crimeData = crimeData.drop('year', axis=1)
7 print (crimeData.shape)
```

(80, 24)

Preprocessing the cleaned data:

```
1 import statsmodels.api as sm
2 y = crimeData['crmrte']
3 X = crimeData['density']
4 X = sm.add_constant(X)
5 model = sm.OLS(y, X).fit()
6 density_pvalue = model.pvalues['density']
7 model.summary()
```

Implementation:

We have implemented the problem using an algorithm called Linear Regression. We have developed a python code based on the algorithm. For this problem we also calculated R Square, P - Value Difference, Number of Observations, Covariance type, etc.

*Corresponding author: sirivanth2000@gmail.com

Parameters crime rate:

```
1 y = crimeData['crrmrte']
2 X = crimeData.drop('crrmrte', axis=1)
3 X = sm.add_constant(X)
4 model = sm.OLS(y, X).fit()
5 model.summary()
```

Parameters density urban:

```
1 y = crimeData['crrmrte']
2 X = crimeData[['density', 'urban']]
3 X = sm.add_constant(X)
4 model = sm.OLS(y, X).fit()
5 density_pvalue_upd = model.pvalues['density']
6 print('Difference in P-Value = ' + str(density_pvalue_upd - density_pvalue))
7 model.summary()
```

Difference in P-Value = 2.2791391426354033e-05

Parameters crime rate urban country:

```
1 y = crimeData['crrmrte']
2 X = crimeData.drop(['crrmrte', 'urban', 'county'], axis=1)
3 X = sm.add_constant(X)
4 model = sm.OLS(y, X).fit()
5 model.summary()
```

Parameters crime rate urban country:

```
1 y = crimeData['crrmrte']
2 X = crimeData.drop(['crrmrte', 'urban', 'county', 'wmfg', 'prbpris', 'wloc', 'west', 'wtuc'], axis=1)
3 X = sm.add_constant(X)
4 model = sm.OLS(y, X).fit()
5 model.summary()
```

Resultant Graphs:

Statistical approach towards crime rate prediction:

```
1 print('Statistics of Crime Rate: \n')
2 print(crimeData['crrmrte'].describe())
3 plt.figure(figsize=(6,6))
4 plt.title('Distribution of Crime Rate Feature')
5 sns.distplot(crimeData['crrmrte'], color='b', bins=100, hist_kws={'alpha': 0.4})
```

Statistics of Crime Rate:

count	80.000000
mean	0.035126
std	0.018846
min	0.010623
25%	0.023359
50%	0.030342
75%	0.041639
max	0.098966
Name: crrmrte, dtype: float64	

Graph Output:

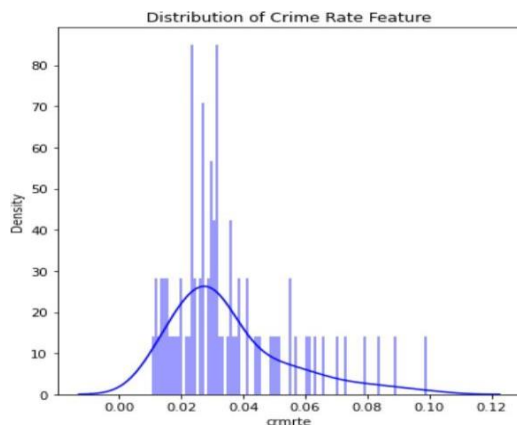


Fig. 1. Distribution of crime rate feature

Crime data mined in different perspectives:

```
1 crimeData.hist(figsize=(18,18), bins=40, xlabelsize=8, ylabelsize=8);
```

Graph Outputs:

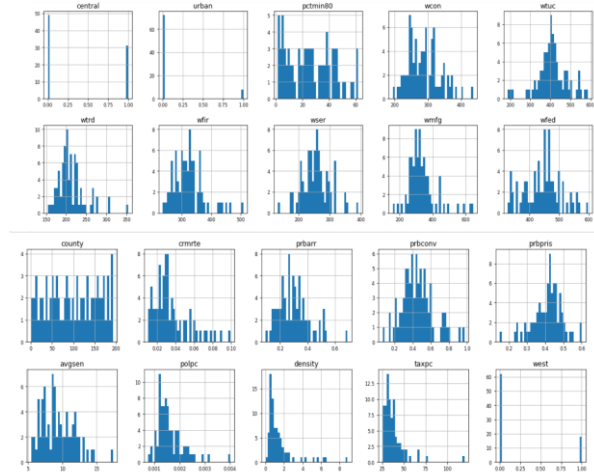


Fig. 2. Output

5. Conclusion

The map of Actual vs. Predicted is linear. This indicates that the forecast is right. The amount of data available for input is small. The plot could be more linear if there was more evidence.

We may also incorporate the Boolean features west, center, and urban into a single feature with categorical values 1, 2, and 3 to create a single feature. A single function like this may be more useful for prediction. On features, functional transformations (such as log, function) can be useful. If there is a possibility to incorporate functionality, using 'unemployment rate' as a proxy for crime rate might be useful.

References

- [1] Y. Xu, C. Fu, E. Kennedy, S. Jiang, and S. Owusu-Agyemang, "The impact of street lights on spatial-temporal patterns of crime in Detroit, Michigan," *Cities*, no. October 2017, pp. 0-1, 2018.
- [2] Shah, N., Bhagat, N. & Shah, M. Crime forecasting: a machine learning and computer vision approach to crime prediction and prevention. *Vis. Comput. Ind. Biomed. Art* 4, 9 (2021).
- [3] Prithi S, Aravindan S, Anusuya E, Kumar AM (2020) GUI based prediction of crime rate using machine learning approach. *Int J Comput Sci Mob Comput* 9(3):221–229
- [4] Pratibha, A. Gahalot, Uprant, S. Dhiman and L. Chouhan, "Crime Prediction and Analysis," *2nd International Conference on Data, Engineering and Applications (IDEA)*, 2020, pp. 1-6.
- [5] B. Sivanagaleela and S. Rajesh, "Crime Analysis and Prediction Using Fuzzy C-Means Algorithm," *2019 3rd International Conference on Trends in Electronics and Informatics (ICOEI)*, 2019, pp. 595-599.
- [6] A. Kumar, A. Verma, G. Shinde, Y. Sukhdeve and N. Lal, "Crime Prediction Using K-Nearest Neighboring Algorithm," *2020 International Conference on Emerging Trends in Information Technology and Engineering (ic-ETITE)*, 2020, pp. 1-4.
- [7] S. Agarwal, L. Yadav and M. K. Thakur, "Crime Prediction Based on Statistical Models," *2018 Eleventh International Conference on Contemporary Computing (IC3)*, 2018, pp. 1-3.
- [8] A. Almaw and K. Kadam, "Crime Data Analysis and Prediction Using Ensemble Learning," *2018 Second International Conference on Intelligent Computing and Control Systems (ICICCS)*, 2018, pp. 1918-1923.
- [9] D. M. Raza and D. B. Victor, "Data mining and Region Prediction Based on Crime Using Random Forest," *2021 International Conference on Artificial Intelligence and Smart Systems (ICAIS)*, 2021, pp. 980-987.