
Statistic Machine Learning Project

Anonymous Author(s)

Affiliation

Address

email

Abstract

The abstract paragraph should be indented 1/2 inch (3 picas) on both the left- and right-hand margins. Use 10 point type, with a vertical spacing (leading) of 11 points. The word **Abstract** must be centered, bold, and in point size 12. Two line spaces precede the abstract. The abstract must be limited to one paragraph.

1 Study case and Data

To further understand and practice the knowledge we have obtained from lectures. A case of classification problem was given that we can apply different methods on the data and understand the characteristics and behaviours of these methods.

1.1 Case and data description

To help the city of Washington, D.C. understand whether increasing the number of public bikes is necessary at some certain hour, a machine learning model is expected to predict the public bike demand for given temporal and meteorological features.

A data set contains 1600 random observations is provided for the model training. The target variable is binary for "high" or "low" demand for increasing bikes' number. The description of the features are given in Table ??

1.2 Exploratory data analysis

Exploratory data analysis has been conducted to gain the knowledge of relations among features in the dataset. Also, following questions in the project introduction are answered in this section.

1. Which are the numerical features and which are the categorical features?

2. Is there any trend to need increase in the availability of

Out of that the ways to treat categorical features are different from numerical features, the features are sorted into two groups, numerical and categorical. This sort of them can be identified according to Table ?. The two groups are shown in Table ?, which also answers Question ?.

For categorical features, first step done was to see the label balance. The result is shown in Figure ?. One thing to be noted is that the feature "snow" has only one label. A flat feature will not have input to the model, thus it was excluded from training set. This will be also seen in the correlation analysis.

Similarly, histogram plot was generated to visualize the distribution of numerical features (shown in Figure ?).

Table 1: Labels and features in the data set (?)

Feature Name	Description
increase_stock (prediction label)	low_bike_demand – no need to increase the number of bikes high_bike_demand – the number of bikes needs to be increased
hour_of_day	Hour of the day (from 0 to 23)
day_of_week	Day of the week (from 0 – Monday to 6 – Sunday)
month	Month (from 0 – January to 12 – December)
holiday	If it is a holiday or not (0 – no holiday, 1 – holiday)
weekday	If it is a weekday or not (0 – weekend, 1 – weekday)
summertime	If it is summertime or not (0 – no summertime, 1 – summertime)
temp	Temperature in Celsius degrees
dew	Dew point in Celsius degrees
humidity	Relative humidity (percentage)
precip	Precipitation in mm
snow	Amount of snow in the last hour in mm
snow_depth	Accumulated amount of snow in mm
windspeed	Wind speed in km/h
cloudcover	Percentage of the city covered in clouds
visibility	Distance in km at which objects or landmarks can be clearly seen and identified

Table 2: Categorical and numerical features.

Categorical features	Numerical features
holiday	hour_of_day
weekday	day_of_week
summertime	month
snow	temp
increase_stock	dew
	humidity
	precip
	snow_depth
	windspeed
	cloudcover
	visibility

31 To understand the correlation among features, mostly between targets and other features, correlation
32 analysis has been done plotting correlation maps, shown in Figure ??

33 Dropping out the smaller correlated features ($correlation \geq 0.1$), the left ones are: "temp", "humidity",
34 "hour_of_day", "summertime", "dew", "weekday", and "visibility". To more intuitively see if any
35 features can help separate the target label. Figure ?? was plotted.

36 Finally, with all analysis and figures above, the Question ?? can be answered.

37 It is seen that the label "1" ("high_bike_demand") is concentrated around daytime, growth starts
38 from early morning, drops till late evening after reaching its peak at 15:00 to 16:00. Weekday has
39 more "high_bike_demand" than weekends.

40 Temperature also has impact on bike demand, more bike are needed when temperature is in a comfort-
41 able region, 20 – 30 degC in data. This can also interpret the trend in the summertime, temperature is
42 higher in summer, similarly for dew point.

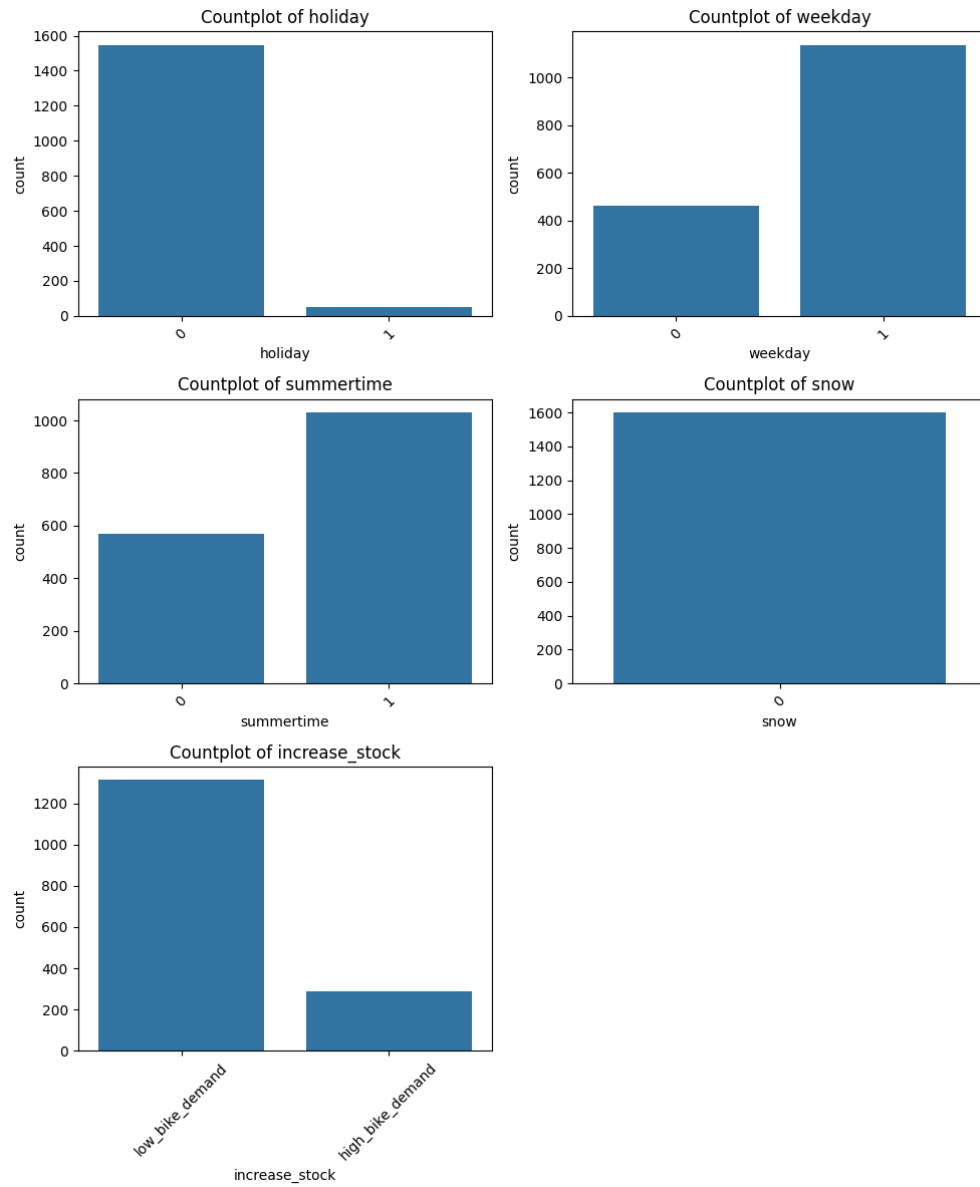


Figure 1: Label counts for categorical features

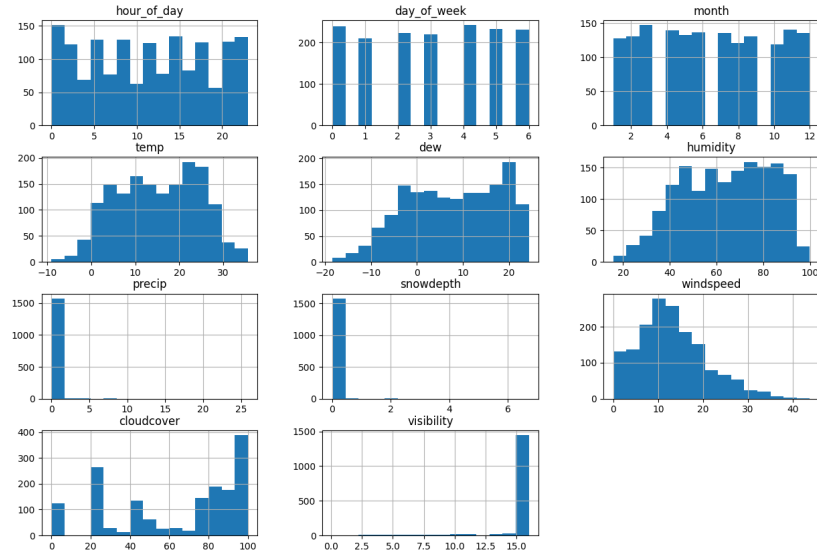


Figure 2: histograms of numerical features

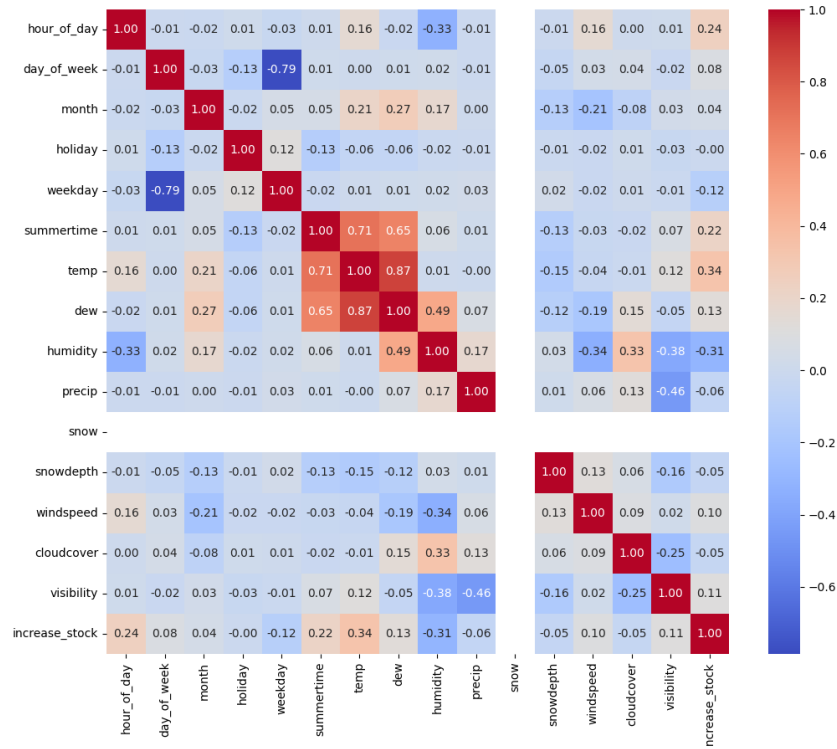


Figure 3: Correlation among features

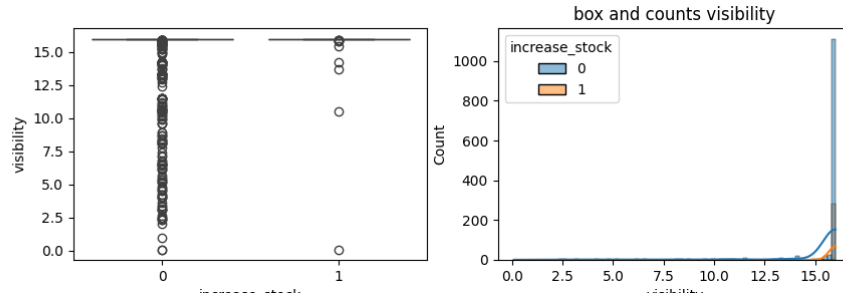


Figure 4: Boxs and Histograms for features with different targets laebls

Table 3: Boosting models perfomance

Model	Training time (s)	Precision (Weighted)	Recall (Weighted)	F1 score (Weighted)
AdaBoostingClassifier	0.008	0.643	0.630	0.633
GradientBoostingClassifier	0.109	0.773	0.655	0.708
CatBoostClassifier	0.593	0.783	0.693	0.734

2 Methodology

2.1 Machine learning models

2.2 Validation

2.2.1 Cross-validation

2.2.2 Metrics

2.2.3 Hyperparameter Optimisation

3 Results

3.1 Boosting

Three boosting models were apply in this study, AdaBoostingClassifier and GradientBoostingClassifier from *Scikit-learn*, and CatBoostClassifier. Their performance is shown in Table ??

According to the performance, CatBoostClassifier was chosen to the model to be Optimized on Hyperparameters, where *GridSearch* was used. The result is shown as a confusion matrix (Figure ?? and also in Table ??).

Table 4: Performance of optimized CatBoostClassifier

Precision (Weighted)	Recall (Weighted)	F1 score (Weighted)
0.85	0.85	0.85

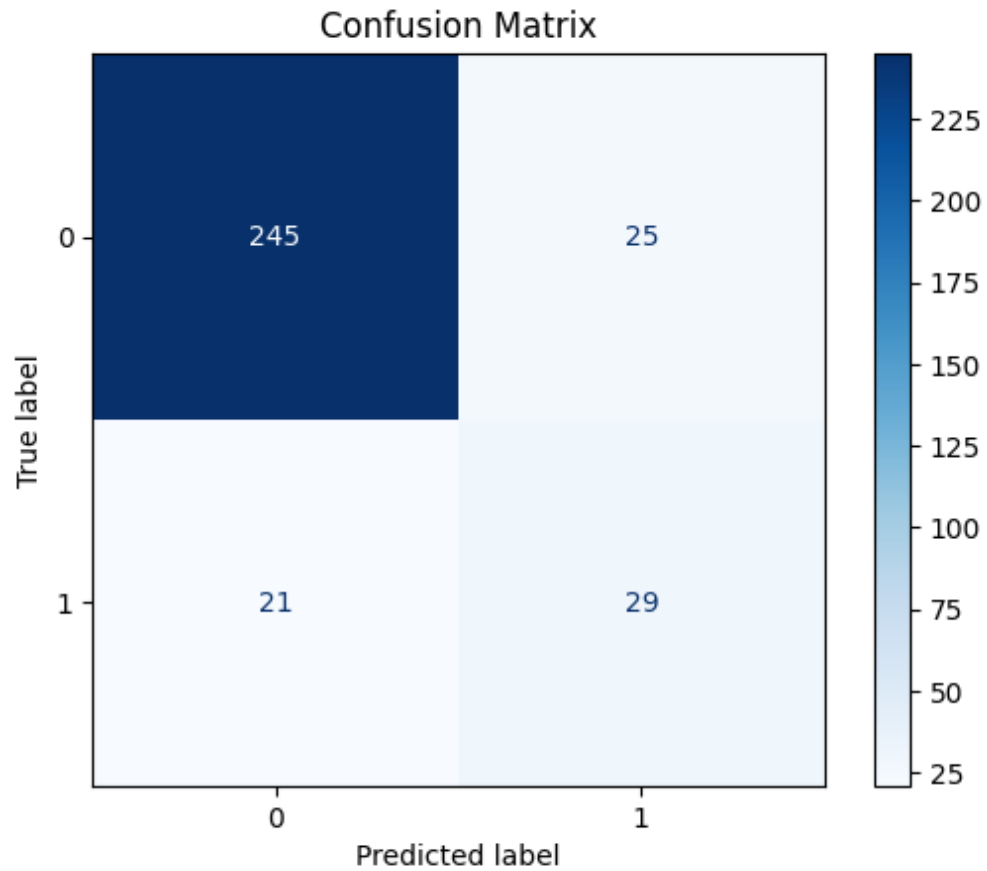


Figure 5: Confusion matrix of Optimized CatBoostClassifier

56 4 Discussion

57 5 Conclusion

58 A Appendix

59 Optionally include extra information (complete proofs, additional experiments and plots) in the
 60 appendix. This section will often be part of the supplemental material.