
Statistic Machine Learning Project

Anonymous Author(s)

Affiliation

Address

email

Abstract

To fulfill the demand of public bike that more fossil based transportation can be replaced by bike and contribute to alleviating climate change, the city of Washington D.C. has recorded the observations high bike demand along with temporal and meteorological features. With this data, this project study is dedicated to help the city to predict the necessity of increasing number of bikes at certain hours. To achieve this goal, methods including Logistic Regression, KNN, bagging and boosting are deployed to predict the target label with given feature.

Number of group members: 3.

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

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 bicycles?

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 1. The two groups are shown in Table 2, which also answers Question 1.

For categorical features, first step done was to see the label balance. The result is shown in Figure 1. One thing to be noted is that the feature "snow" has only one label. A flat feature will not have input

Table 1: Labels and features in the data set [1]

Feature Name	Description
midrule 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
midrule 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	month
weekday	temp
summertime	dew
snow	humidity
increase_stock	precip
hour_of_day	snow_depth
day_of_week	windspeed
	cloudcover

31 to the model, thus it was excluded from training set. This will be also seen in the correlation analysis.
 32 Similarly, histogram plot was generated to visualize the distribution of numerical features (shown in
 33 Figure 2).

34 To understand the correlation among features, mostly between targets and other features, correlation
 35 analysis has been done plotting correlation maps, shown in Figure 3

36 Dropping out the smaller correlated features ($correlation \geq 0.1$), the left ones are: "temp", "humidity",
 37 "hour_of_day", "summertime", "dew", "weekday", and "visibility". To more intuitively see if any
 38 features can help separate the target label. Figure 4 was plotted.

39 Finally, with all analysis and figures above, the Question 2 can be answered.

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

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

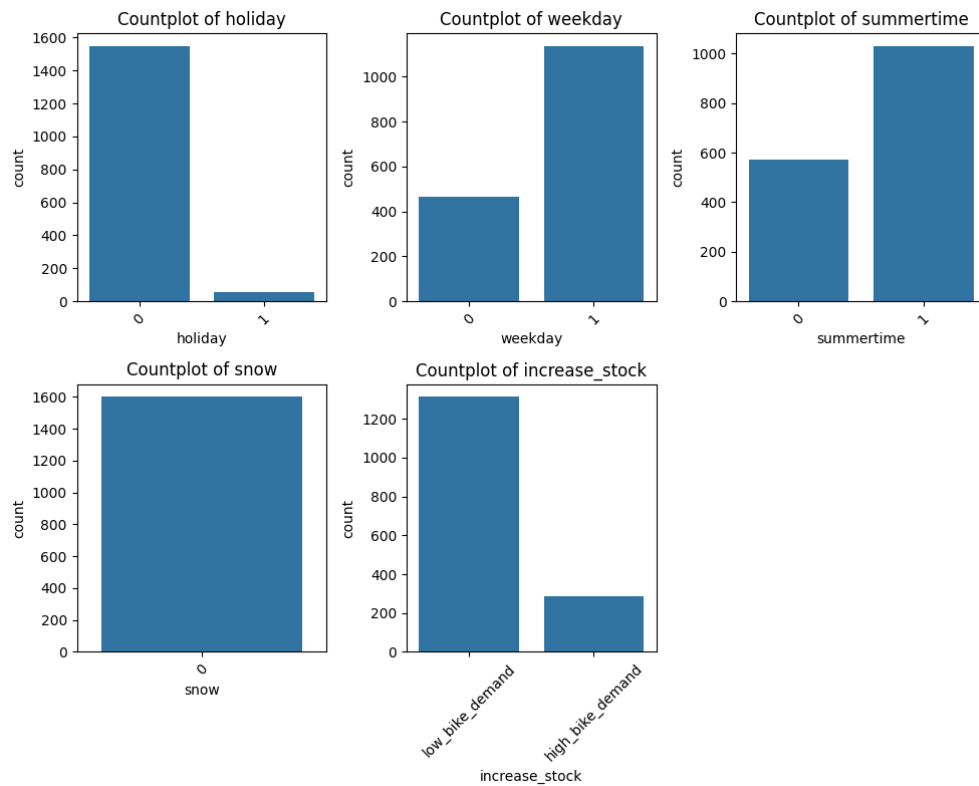


Figure 1: Label counts for categorical features

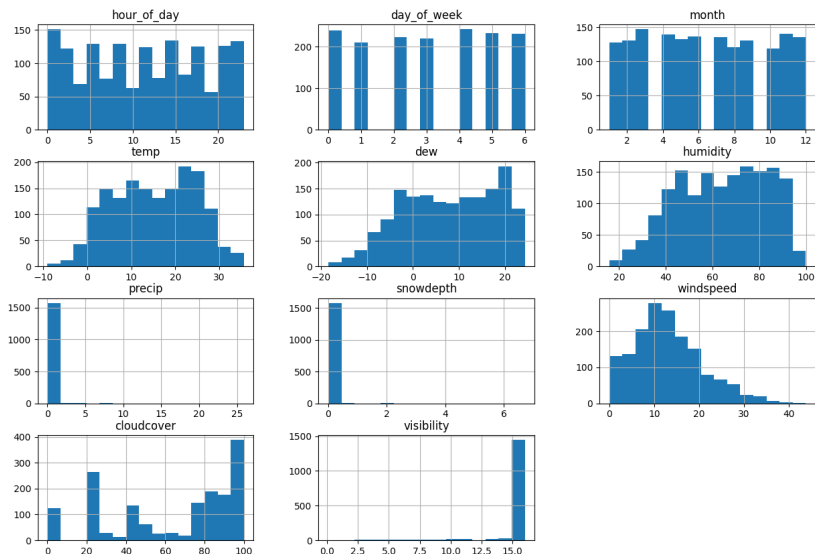


Figure 2: histograms of numerical features

46 2 Methodology

47 2.1 Machine learning models

48 2.1.1 Boosting

49 Boosting is an emsenble method which is built on the idea that even a week model can capture some
50 patterns between target variables and given features. Starting from a base model, optmizing it based
51 on the returned error can produce a new model which would become the base for next model. By
52 ensembling the predictions from these models, the intention of boosting is to reduce bias[2].

53

54 In this study, three maintream boosting algorithm were deployed, **AdaBoost**[2], **GradientBoost**[2],
55 and **CatBoost**[3]. The methods can expressed as following equations.

$$\text{AdaBoost: } \hat{y}_{\text{boost}}^{(B)}(\mathbf{x}) = \text{sign} \left(\sum_{b=1}^B \alpha^{(b)} \hat{y}^{(b)}(\mathbf{x}) \right) \quad (1)$$

$$\text{GradientBoost: } f^{(B)}(\mathbf{x}) = \sum_{b=1}^B \alpha^{(b)} f^{(b)}(\mathbf{x}), \quad (2)$$

$$\text{CatBoost: } f^{(B)}(\mathbf{x}) = \sum_{b=1}^B \eta^{(b)} T^{(b)}(\mathbf{x}) \quad (3)$$

56 2.1.2 Logistic Regression

57 Logistic Regression is one of the most common supervised machine learning models for binary
58 classification problems. Unlike Linear Regression, it uses a sigmoid function, where the raw output
59 value is passed through and converted into a probability between 0 and 1. Since our problem has two
60 possible outcomes, this model is specifically a binomial logistic regression.

61

62 Logistic Regression first calculates the linear combination of the input features:

$$z = w_0 + w_1x_1 + w_2x_2 + \cdots + w_nx_n$$

63 Here, x_1, x_2, \dots, x_n are the input feature values, w_0, w_1, \dots, w_n are the model's learned weights,
64 and z is the model's internal score. This value is then passed through the sigmoid function to convert
65 it into a probability between 0 and 1.

66

67 2.1.3 K-Nearest Neighbours

68 K-Nearest Neighbours (KNN) is a supervised machine learning algorithm used for both classification
69 and regression problems. It works by finding the K nearest points to a given data point and making
70 predictions based on the majority class for classification or the average value for regression.

71 KNN is a simple, non-parametric, and lazy learning algorithm because it performs computation
72 during prediction rather than during training.

73 The distance between two data points, \mathbf{x}_1 and \mathbf{x}_2 , is measured using the Euclidean distance formula:

$$d(\mathbf{x}_1, \mathbf{x}_2) = \sqrt{(x_{11} - x_{21})^2 + (x_{12} - x_{22})^2 + \cdots + (x_{1n} - x_{2n})^2} \quad (4)$$

74 For a new data point \mathbf{X}_{new} , the algorithm calculates its distance from all instances in the dataset and
75 selects the K instances with the smallest distances. The final predicted class is determined based on
76 the majority vote among these K nearest neighbours.

77 It is important to note that if K is too small, the model becomes overfitted. If K is too large, the
78 model becomes underfitted. In our project, we selected the optimal K value using the GridSearchCV
79 function.

80 2.2 Validation

81 K-fold method with 5 subset was used for cross-validation to the deployed models. Considering that
82 label are imbalanced in the target variable, metrics like accuracy would not be a good choice for
83 performance measurement. F-score, given by equation 5, is used for measuring the performance in this
84 case.

$$F_1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}} = \frac{2TP}{2TP + FP + FN} \quad (5)$$

85 3 Results

86 As mentioned in Section 2, different models of four family were deployed. The result is shown in
87 Table 3.

Table 3: f1-score of best models

model	Cat Boost		KNN		Logistic Regression		Random Forest	
	0	1	0	1	0	1	0	1
F1-score	0.93	0.61	0.92	0.57	0.94	0.62	0.96	0.82
Support	270	50	270	50	270	50	262	58

88 4 Discussion

89 5 Conclusion

90 References

- 91 [1] Do we need more bikes? project in machine learning. Technical report, Department of Informa-
92 tion Technology, Uppsala University, November 2024. URL N/A. Course: Statistical Machine
93 Learning, 1RT700.
- 94 [2] Andreas Lindholm, Niklas Wahlström, Fredrik Lindsten, and Thomas B. Schön. *Machine*
95 *Learning: A First Course for Engineers and Scientists*. Cambridge University Press, 2022.
96 Pre-publication draft, July 8, 2022. Cambridge University Press. Available at <http://smlbook.org>.
- 97 [3] Liudmila Prokhorenkova, Gleb Gusev, Aleksandr Vorobev, Anna Veronika Dorogush, and Andrey
98 Gulin. Catboost: unbiased boosting with categorical features. In *Advances in Neural Information*
99 *Processing Systems (NeurIPS)*, volume 31, 2018.

100 A Appendix