

---

# Do we need more bikes?

## Project in Statistical Machine Learning

---

**Anonymous Author(s)**

Affiliation

Address

email

### Abstract

1 In this project we develop, and study different statistical machine learning models  
2 for predicting whether the number of available bikes at a given hour should be  
3 increased, a project by the District Department of Transportation in Washington  
4 D.C. The training data set consists of 1600 instances of hourly bike rentals, and  
5 a test set of 400 instances. The models for prediction we have used are: *Logistic*  
6 *regression*, *Discriminant methods: LDA, QDA, k- Nearest Neighbour*, and *Tree*  
7 *Based Methods*. We have found that THE MODEL gives best prediction, with  
8 accuracy ??????

9 **1 Plan**

10 **1.1 From Intro**

11 (i) Explore and preprocess data

12 (ii) try some or all classification methods, which are these?

13     • Logistic Regression

14     • Discriminant analysis: LDA, QDA

15     • K-nearest neighbor

16     • Tree-based methods: classification trees, random forests, bagging

17     • Boosting

18 (iii) Which of these are to be "put in production"?

19 **1.2 From Data analysis task**

20     • Can any trend be seen comparing different hours, weeks, months?

21     • Is there any difference between weekdays and holidays?

22     • Is there any trend depending on the weather?

23 **1.3 From Implementation of methods**

24 Each group member should implement one family each, who did what shall be clear!

25 DNNs are encouraged to be implemented, do this if there is time. (DNN is not a thing a group

26 member can claim as their family.)

27 Implement a naive version, let's do: *Always low\_bike\_demand*

28 **1.3.1 What to do with each method**

29     1. Implement the method (each person individually)

30     2. Tune hyper-parameters, discuss how this is done (each person individually)

31     3. Evaluate with for example cross-validation. Don't use  $E_{k-fold}$  (what is that?) (need to do

32         together)

33     4. (optional) Think about input features, are all relevant? (together)

34 Before training, unify pre-processing FOR ALL METHODS and choose ONE OR MULTIPLE

35 metrics to evaluate the model. (is it necessary to have the same for all?, is it beneficial?) Examples:

36     • accuracy

37     • f1-score

38     • recall

39     • precision

40 Use same test-train split for ALL MODELS

## 41 2 Theoretical background

### 42 2.1 Mathematical Overview of the Models

#### 43 2.1.1 Logistic Regression

44 The backbone of logistic regression is linear regression, i.e. finding the least-squares solution to an  
45 equation system

$$X\theta = b \quad (1)$$

46 given by the normal equations

$$X^T X \theta = X^T b \quad (2)$$

47 where  $X$  is the training data matrix,  $\theta$  is the coefficient vector and  $b$  is the training output. The  
48 parameter vector is then used in the sigmoid function:

$$\sigma(z) = \frac{e^z}{1 + e^z} : \mathbb{R} \rightarrow [0, 1], \quad (3)$$

$$z = x^T \theta, \quad (4)$$

49 where  $x$  is the testing input. This gives a statistical interpretation of the input vector. In the case of a  
50 binary True/False classification, the value of the sigmoid function then determines the class.

### 51 2.2 Input Data Modification

52 By plotting the data and analyzing the .csv file, some observations were made. The different inputs  
53 were then changed accordingly:

- 54 • *Kept as-is:* weekday, windspeed, visibility, temp
- 55 • *Modified:*
  - 56 – month - split into two inputs, one cosine and one sine part. This make the new inputs
  - 57 linear and can follow the fluctuations of the year. The original input was discarded.
  - 58 – hour\_of\_day - split into three boolean variables: demand\_day, demand\_evening,
  - 59 and demand\_night, reflecting if the time was between 08-14, 15-19 or 20-07 respec-
  - 60 tively. This was done because plotting the data showed three different plateaues of
  - 61 demand for the different time intervals. The original input was discarded.
  - 62 – snowdepth, precip were transformed into booleans, reflecting if it was raining or
  - 63 if there was snow on the ground or not. This was done as there was no times where
  - 64 demand was high when it was raining or when there was snow on the ground.
- 65 • *Removed:* cloudcover, day\_of\_week, snow, dew, holiday, summertime, humidity.
- 66 These were removed due to being redundant (e.g. summertime), not showing a clear trend
- 67 (e.g. cloudcover) or both (e.g. day\_of\_week).