
Admission Data Prediction Using Machine Learning Methods

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The

2 1 Introduction

3 NeurIPS requires electronic submissions. The electronic submission site is

4 <https://cmt3.research.microsoft.com/NeurIPS2020/>

5 Please read the instructions below carefully and follow them faithfully.

6 2 Methodology

7 The dataset we chose is created for prediction of Graduate Admissions from an Indian perspective,
8 which predicting admission from 7 important parameters with 500 students. The output is a number
9 from 0 to 100, which represents the probability a student being admitted. Therefore, we consider it as a
10 regression problem.

11 2.1 Preprocessing

12 Firstly, we do the data splitting process, we randomly divide 500 input data into three parts: 320 train
13 data, 80 validation data, and 100 testing data. Secondly, we do subset selection to find the best subset
14 of 7 feature parameters. Based on the RSS loss of linear regression, we find out that the best subset is
15 the total set, we do not need to filter any feature parameter. Then, we normalize the input data before
16 loading them into algorithm models. Additionally, some algorithm may not support regression task,
17 such as logistic regression, LDA, and Naive Bayes, so for these algorithm, we change the regression
18 task into classification task by approximating the output number into 10 neighbor classes: 0, 10, 20,
19 and so on.

20 2.2 Algorithm

21 In order to solve this problem, we considering both classification methods and regression methods.
22 The original data is continuous and we will preprocessing it if we want to use classification methods.

23 **Least square** Fit a linear model with coefficients $w = (w_1, \dots, w_p)$ to minimize the sum of squared
24 residuals between the actual observed data and the predicted data (estimated values) of the data set:
25 $\min_w ||Xw - y||_2^2$.

26 **Ridge regression** Ridge regression solves some problems of ordinary least squares by penalizing
27 the size of the coefficients. What minimizes the ridge coefficient is the sum of squared residuals with
28 penalties: $\min_w ||Xw - y||_2^2 + \alpha ||w||_2^2$.

29 **Lasso regression** Lasso regression consists of a linear model with regular terms of l_1 -norm. Its
30 minimized objective function is: $\min_w \frac{1}{2n_{samples}} \|Xw - y\|_2^2 + \alpha \|w\|_1$.
31 **Knn** Knn is also a regression method, it is used when the data labels are continuous variables rather
32 than discrete variables. The label assigned to the query point is calculated from the average of its
33 nearest neighbor labels.
34 **Decision tree** The nearest neighbor regression is used when the data labels are continuous variables
35 rather than discrete variables. The label assigned to the query point is calculated from the average of
36 its nearest neighbor labels.
37 **SVM** It is very efficient in high-dimensional space, and different kernel functions have a one-to-one
38 correspondence with specific decision functions. Common kernels are already provided, and custom
39 kernels can also be specified.
40 **Boosting** The goal of the boosting method is to combine the prediction results of multiple base esti-
41 mators constructed using a given learning algorithm to obtain better generalization ability/robustness
42 than a single estimator. We mainly focus on Random Forest and AdaBoost.
43 **LDA** This is a classification method. It is derived from simple probability models, and these models
44 can be obtained by Bayes' theorem for the relevant distribution $P(X|y = k)$ of each category k .
45 **Naïve Bayes** This is a classification method. Naïve Bayes methods are a set of supervised learning
46 algorithms based on applying Bayes' theorem with the "naïve" assumption of conditional indepen-
47 dence between every pair of features given the value of the class variable.
48 **Logistic** This is a classification method. Logistic regression is a generalized linear model, so it has
49 many similarities with multiple linear regression analysis. Their model form is basically the same. It
50 gets dependent variable value by logistic function.

51 3 Experiment

52 3.1 Preprocessing

53 In the preprocessing process, after splitting data into training, validation and testing data, and before
54 normalization, we do subset selection to find the best subset of 7 feature parameters. During subset
55 selection, For each $s \in 0, 1, \dots, p$, find the subset in size of s that gives lowest RSS, and use cross-
56 validation to estimate prediction error and select s . Then, we can select the optimal variables. The
57 result is showed below. We can learn from the result that the best subset is the total set, we do
58 not need to filter any feature parameter. In addition, it is quite friendly for us to use the method to
59 selection the best subset, since it need a lot of computation and a lot of time when p is too large, but
60 our p is 7, and the dataset is small, so the running time is not too long.

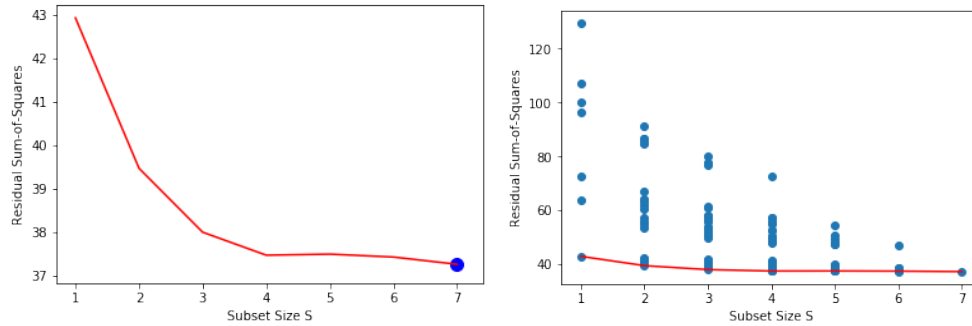


Figure 1: Subset Selection Result.

61 3.2 Algorithm

62 3.2.1 Regression

63 First, we use regression algorithm to fit the admission rate. We use without shrinkage, lasso and ridge
64 models to find the best model, by optimizing the parameter α through the analysis of RSS error.
65 Alpha indicates the degree of shrinkage. When α approaches 1, it indicates that the degree of
66 shrinkage reaches its maximum; when α approaches 0, it indicates that there is no shrinkage. The

67 RSS of the three methods varies with alpha are shown as follows. It can be seen that the smallest
 68 RSS is at lasso regression, and the alpha at this time is 0.3419, the accuracy is 0.8706.

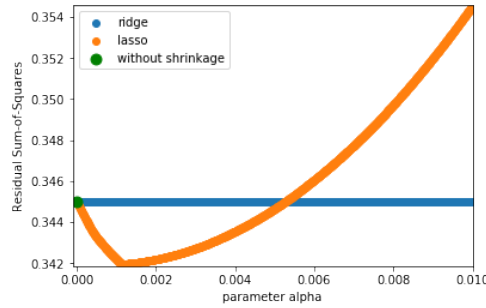


Figure 2: Subset Selection Result.

68

69 3.2.2 Decision Tree

70 In decision tree algorithm, we find the optimal model by finding at which depth we will get the lowest
 71 residual sum-of-squares in validation set. From Figure 2 we can know that we can get the lowest
 72 residual sum-of-squares at $depth = 4$, and the RSS value is 0.4518. The accuracy of this method is 0.8386.

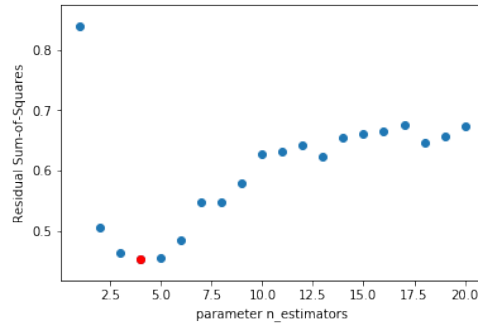


Figure 3: Subset Selection Result.

73

74 3.2.3 KNN

75 Using KNN regression, we need to find the optimal k value. We want to find at which k value we
 will get the lowest residual sum-of-squares in validation set.

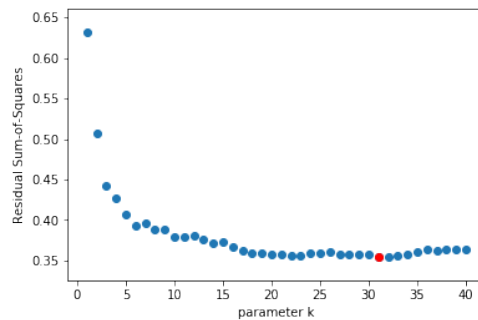


Figure 4: Knn Result.

76

77 From Figure 2 we can know that we can get the lowest residual sum-of-squares at $k = 31$, and the
 78 RSS value is 0.3546. The accuracy of this method is 0.8706.

79 3.2.4 SVM

80 In SVM regression method, we can apply different kernel on it. The min error without any kernel is
 81 0.5596.

- 82 • **rbf kernel:** We will find the lowest residual sum-of-squares in validation set at $gamma = 5.0351e - 05$, and the RSS value is 0.4495. The accuracy of this method is 0.6179.

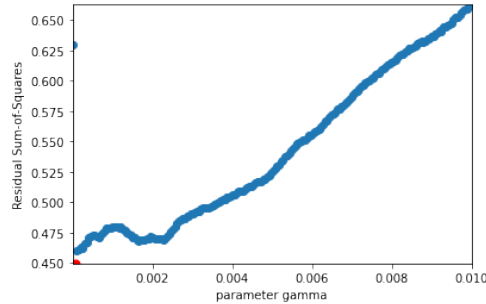


Figure 5: rbf kernel

83

- 84 • **linear kernel:** We will find the lowest residual sum-of-squares in validation set at $C = 0.0918$, and the RSS value is 0.4476. The accuracy of this method is 0.6221.

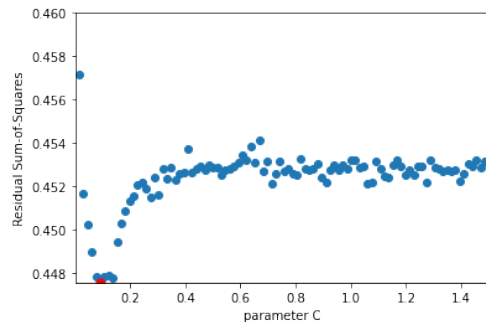


Figure 6: linear kernel

85

- 86 • **poly kernel:** We will find the lowest residual sum-of-squares in validation set at $degree = 1$,
 87 and the RSS value is 0.4532. The accuracy of this method is 0.6211.

88 3.2.5 AdaBoost

89 When we using AdaBoost method, we are using a lot of weak estimators to regress this problem. So
 90 we want to find the optimal number of the weak estimators.

91 We will find the lowest residual sum-of-squares in validation set at the $n_estimators = 10$, and the
 92 RSS value is 0.3951. The accuracy of this method is 0.6836.

93 At the same time, we realized that the running time may have some relation with the number of the
 94 weak estimators, and then we record the running time of this method with different number of the
 95 weak estimators.

96 The min time is 11035 microseconds at $n_estimators = 1$, and the running time increases as the
 97 number of the weak estimators increases. We can get the conclusion that the optimal RSS may not
 98 correspond with the shortest running time.

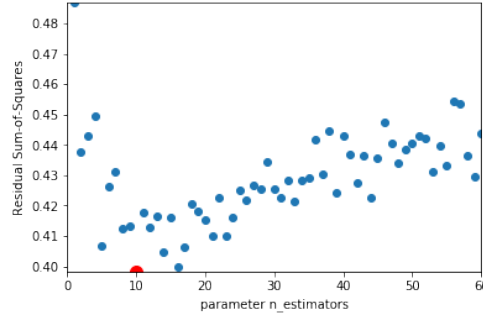


Figure 7: AdaBoost: estimator number and RSS

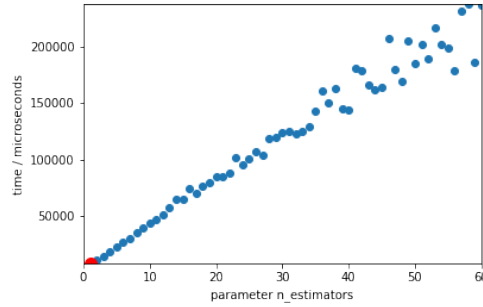


Figure 8: AdaBoost: estimator number and RSS

3.2.6 Random Forest

In random forest algorithm, we need to find the optimal model by finding the optimal $n_estimators$ and optimal depth, by get the lowest residual sum-of-squares in validation set.

Firstly, we choose best $n_estimators$ by running model with diffrent $n_estimators$ and choose the best. But during this process, we find out as $n_estimators$ increase, the time that cost to run the algorithm increase linearly, so its not elegant to choose large $n_estimators$. After finding the optimal $n_estimators$, we start finding the optimal depth with optimal $n_estimators$, the thild figure in Figure 4 shows that at $depth = 4$, the validation error reach maximum. Also, the model start overfitting at around $depth = 3$.

From Figure 3 we can know that we can get the lowest residual sum-of-squares in the validation set at $depth = 4$, $n_estimators = 13$, and the RSS value is 0.3797. The accuracy of this method is 0.8303.

3.2.7 Other

Besides the regression method above, we also use some classification methods to solve this problem.

- **Logistic** When we try different penalty function, we find that l_1 function is a little bit better than l_2 . The lowest residual sum-of-squares in validation set is 0.7880, and the accuracy is 0.4810.
- **Naïve Bayes** Comparing Gaussian NB and Bernoulli NB, we find that Gaussian NB has better effect. The lowest residual sum-of-squares in validation set is 0.7, and the accuracy is 0.5952.
- **LDA** The LDA method with default solver *svd* can reach the accuracy 0.5833, and the lowest residual sum-of-squares in validation set is 0.5440.

4 Conclusion

These instructions apply to everyone.

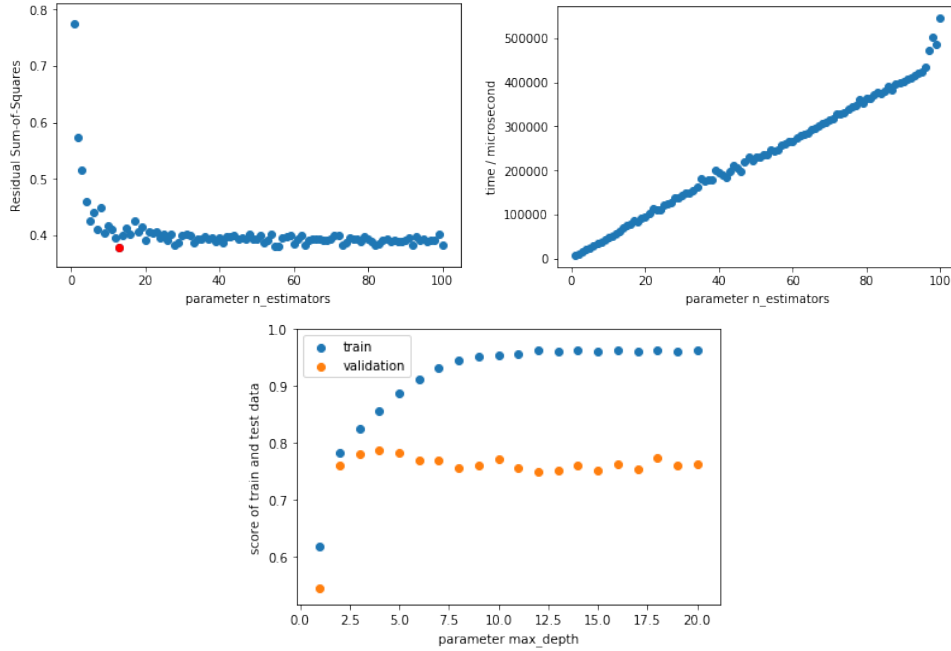


Figure 9: Subset Selection Result.

Table 1: Sample table title

Algorithm	RSS error	Accuracy
regression(Lasso)	0.3419	87.057%
KNN	0.3546	81.570%
Decision Tree	0.4518	83.863%
SVM(Linear)	0.4476	62.206%
AdoBoost	0.3982	69.178%
Random Forest	0.3797	83.029%
LDA	0.5440	58.333%
Naive Bayes(Gaussian)	0.7000	59.5238%
Logistic(l1-penalty)	0.7880	48.0952%

4.1 Citations within the text

The natbib package will be loaded for you by default. Citations may be author/year or numeric, as long as you maintain internal consistency. As to the format of the references themselves, any style is acceptable as long as it is used consistently.

The documentation for natbib may be found at

<http://mirrors.ctan.org/macros/latex/contrib/natbib/natnotes.pdf>

Of note is the command `\citet`, which produces citations appropriate for use in inline text. For example,

```
\citet{hasselmo} investigated\dots
```

produces

Hasselmo, et al. (1995) investigated...

If you wish to load the natbib package with options, you may add the following before loading the neurips_2020 package:

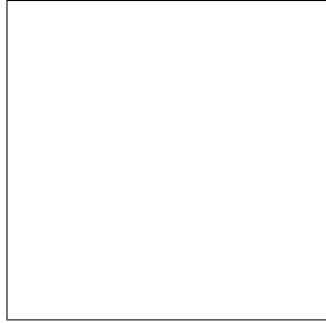


Figure 10: Sample figure caption.

136 `\PassOptionsToPackage{options}{natbib}`

137 If `natbib` clashes with another package you load, you can add the optional argument `nonatbib`
138 when loading the style file:

139 `\usepackage[nonatbib]{neurips_2020}`

140 As submission is double blind, refer to your own published work in the third person. That is, use “In
141 the previous work of Jones et al. [4],” not “In our previous work [4].” If you cite your other papers
142 that are not widely available (e.g., a journal paper under review), use anonymous author names in the
143 citation, e.g., an author of the form “A. Anonymous.”

144 4.2 Footnotes

145 Footnotes should be used sparingly. If you do require a footnote, indicate footnotes with a number¹
146 in the text. Place the footnotes at the bottom of the page on which they appear. Precede the footnote
147 with a horizontal rule of 2 inches (12 picas).

148 Note that footnotes are properly typeset *after* punctuation marks.²

149 4.3 Figures

150 All artwork must be neat, clean, and legible. Lines should be dark enough for purposes of reproduction.
151 The figure number and caption always appear after the figure. Place one line space before the figure
152 caption and one line space after the figure. The figure caption should be lower case (except for first
153 word and proper nouns); figures are numbered consecutively.

154 You may use color figures. However, it is best for the figure captions and the paper body to be legible
155 if the paper is printed in either black/white or in color.

156 4.4 Tables

157 All tables must be centered, neat, clean and legible. The table number and title always appear before
158 the table. See Table 1.

159 Place one line space before the table title, one line space after the table title, and one line space after
160 the table. The table title must be lower case (except for first word and proper nouns); tables are
161 numbered consecutively.

162 Note that publication-quality tables *do not contain vertical rules*. We strongly suggest the use of the
163 `booktabs` package, which allows for typesetting high-quality, professional tables:

164 <https://www.ctan.org/pkg/booktabs>

165 This package was used to typeset Table 1.

¹Sample of the first footnote.

²As in this example.

5 Final instructions

Do not change any aspects of the formatting parameters in the style files. In particular, do not modify the width or length of the rectangle the text should fit into, and do not change font sizes (except perhaps in the **References** section; see below). Please note that pages should be numbered.

6 Preparing PDF files

Please prepare submission files with paper size “US Letter,” and not, for example, “A4.”

Fonts were the main cause of problems in the past years. Your PDF file must only contain Type 1 or Embedded TrueType fonts. Here are a few instructions to achieve this.

- You should directly generate PDF files using `pdflatex`.
- You can check which fonts a PDF file uses. In Acrobat Reader, select the menu Files>Document Properties>Fonts and select Show All Fonts. You can also use the program `pdf fonts` which comes with `xpdf` and is available out-of-the-box on most Linux machines.
- The IEEE has recommendations for generating PDF files whose fonts are also acceptable for NeurIPS. Please see <http://www.emfield.org/icuwb2010/downloads/IEEE-PDF-SpecV32.pdf>
- `xfig` “patterned” shapes are implemented with bitmap fonts. Use “solid” shapes instead.
- The `\bbold` package almost always uses bitmap fonts. You should use the equivalent AMS Fonts:

```
\usepackage{amsfonts}
```

followed by, e.g., `\mathbb{R}`, `\mathbb{N}`, or `\mathbb{C}` for \mathbb{R} , \mathbb{N} or \mathbb{C} . You can also use the following workaround for reals, natural and complex:

```
\newcommand{\RR}{\mathbb{R}} %real numbers
\newcommand{\Nat}{\mathbb{N}} %natural numbers
\newcommand{\CC}{\mathbb{C}} %complex numbers
```

Note that `amsfonts` is automatically loaded by the `amssymb` package.

If your file contains type 3 fonts or non embedded TrueType fonts, we will ask you to fix it.

6.1 Margins in L^AT_EX

Most of the margin problems come from figures positioned by hand using `\special` or other commands. We suggest using the command `\includegraphics` from the `graphicx` package. Always specify the figure width as a multiple of the line width as in the example below:

```
\usepackage[pdftex]{graphicx} ...
\includegraphics[width=0.8\linewidth]{myfile.pdf}
```

See Section 4.4 in the `graphics` bundle documentation (<http://mirrors.ctan.org/macros/latex/required/graphics/grfguide.pdf>)

A number of width problems arise when L^AT_EX cannot properly hyphenate a line. Please give LaTeX hyphenation hints using the `\-` command when necessary.

References

References follow the acknowledgments. Use unnumbered first-level heading for the references. Any choice of citation style is acceptable as long as you are consistent. It is permissible to reduce the font size to `small` (9 point) when listing the references. **Note that the Reference section does not count towards the eight pages of content that are allowed.**

- 207 [1] Alexander, J.A. & Mozer, M.C. (1995) Template-based algorithms for connectionist rule extraction. In
208 G. Tesauro, D.S. Touretzky and T.K. Leen (eds.), *Advances in Neural Information Processing Systems 7*, pp.
209 609–616. Cambridge, MA: MIT Press.
- 210 [2] Bower, J.M. & Beeman, D. (1995) *The Book of GENESIS: Exploring Realistic Neural Models with the*
211 *GEneral NEural Simulation System*. New York: TELOS/Springer-Verlag.
- 212 [3] Hasselmo, M.E., Schnell, E. & Barkai, E. (1995) Dynamics of learning and recall at excitatory recurrent
213 synapses and cholinergic modulation in rat hippocampal region CA3. *Journal of Neuroscience* **15**(7):5249-5262.