# Formatting instructions for NeurIPS 2018

**Rui Cao**
University of British Columbia
`rui.cao@alumni.ubc.ca`

**Li Wang**
University of British Columbia
`lwang@math.ubc.ca`

## Abstract

We use various machine learning models study the shopping behavior people based on a data set of 537577 observation about the black Friday in a retail store.

## 1   Introduction

Our data is from `https://www.kaggle.com/mehdidag/black-friday`, which the feature columns consists of 7 categories of the consumers, such as ages, residence in the city, job position, etc. We would like to explore to the correlations among the above features, so from where we can use machine learning models to make predictions.
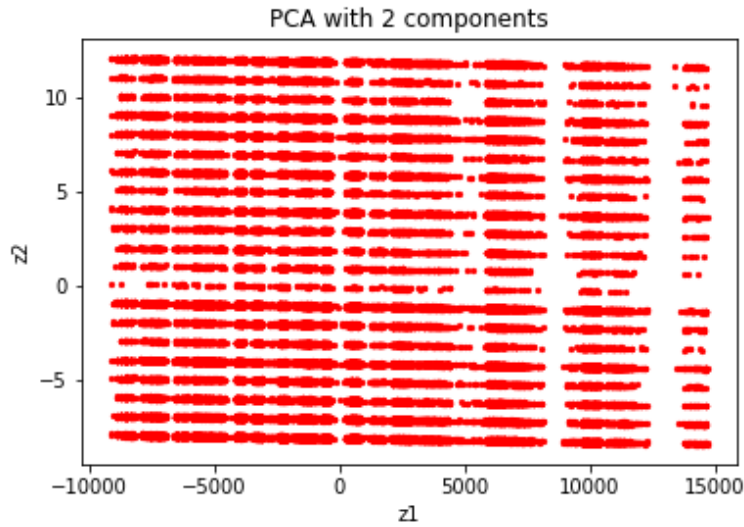
## 2   Related Work

Need reference here....

## 3   Data Analysis

We first analysis the data set, which has total 537577 data. There are seven features corresponding each data, and they are age groups, occupation, city category, residence in the current city, marital status, three product categories, purchase amounts. There are some sparsity in the product category 2,and 3. So we fill those blank as 0. For each customers, it is possible to purchase different items in one category, so one customer may appear several times in the data set. Considering that, there are in total 5891 customers.

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 537577 entries, 0 to 537576
Data columns (total 12 columns):
User_ID                     537577 non-null int64
Product_ID                  537577 non-null object
Gender                      537577 non-null object
Age                         537577 non-null object
Occupation                  537577 non-null int64
City_Category               537577 non-null object
Stay_In_Current_City_Years  537577 non-null object
```

**PCA with 2 components**



```
Marital_Status          537577 non-null int64
Product_Category_1      537577 non-null int64
Product_Category_2      370591 non-null float64
Product_Category_3      164278 non-null float64
Purchase                537577 non-null int64
```
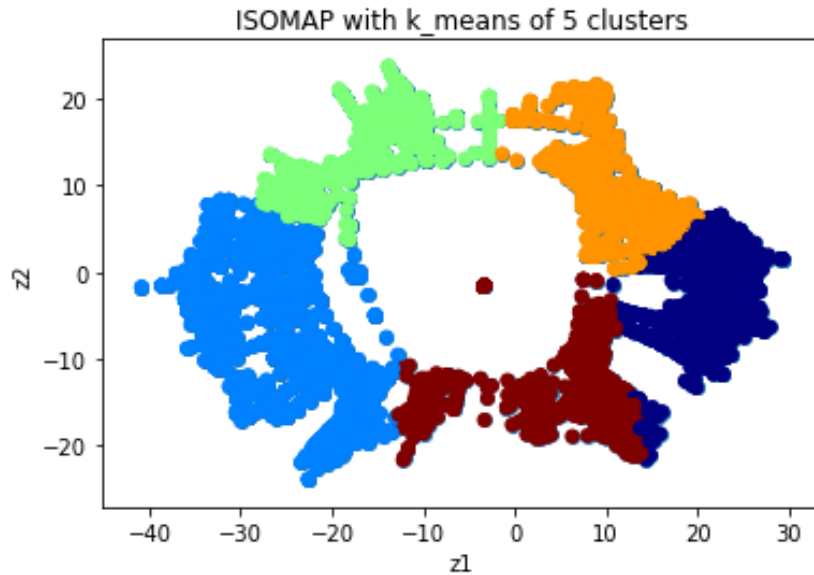
### 3.1 Data Visualization

We first try to visualize data by compressing data using both PCA and ISOMAP. Because of the sparsity of product category 2, and 3. We removed these two features in our data. By randomly choosing the 10000 data from our data set, the PCA variance of the first two components are 9.99997580e-01 1.72858695e-06, which we can see in the figure that the vector is pointing one directions. If we apply the ISOMAP components and K means, we can understand the data from a different perspective. ISOMAP provides us some information of the hidden nonlinearity of the data. However both methods does not give us good clustering results. Thus, we expect that clustering would not be a good model for our data.
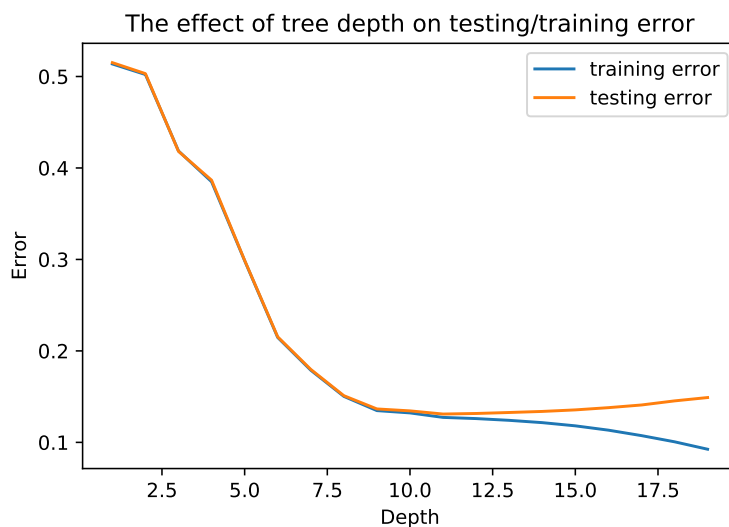
### 3.2 Prediction

To understand which product in category 1 are mostly bought by the customers based on their other information. We try to predict the product category 1 by first dropping the feature of category 2 and 3, and based on the other features of data. We use 5 folds cross validation here. We have tried the following models : decision tree and random forests, and KNN, and softmax. Decision tree and random forests do a better job than the other twos. So from the plot of test error and training error vs depth of tree, we choose decision tree as depth 11, and we have the training error as 0.127 and testing error: 0.131. Random forest with 5 features and 5 trees. has Training error 0.015, and testing error 0.154. However, KNN with metric of cosine is not a good model for our case, which has training error 0.363, and validation error 0.496. Softmax has Training error 0.651 and Validation error 0.652. We also tried polynomial kernel, and Gaussian kernel, both does not give us satisfying results.

From our prediction, it provides information to the retail business that what kind of products in category 1 should they recommend to the customers based on their demographic and social and economic information.

clusters.png



## 4   Citations, figures, tables, references

These instructions apply to everyone.

### 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,

Figure 1: Sample figure caption.

```
\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_2018` package:

```
\PassOptionsToPackage{options}{natbib}
```

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

```
\usepackage[nonatbib]{neurips_2018}
```

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

### 4.2 Footnotes

Footnotes should be used sparingly. If you do require a footnote, indicate footnotes with a number[1] in the text. Place the footnotes at the bottom of the page on which they appear. Precede the footnote with a horizontal rule of 2 inches (12 picas).

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

### 4.3 Figures

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

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

### 4.4 Tables

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

---

[1]Sample of the first footnote.
[2]As in this example.

Table 1: Sample table title

| | Part | | Size ($\mu$m) |
|---|---|---|---|
| Name | Description | | |
| Dendrite | Input terminal | | $\sim$100 |
| Axon | Output terminal | | $\sim$10 |
| Soma | Cell body | | up to $10^6$ |

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

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

$$\texttt{https://www.ctan.org/pkg/booktabs}$$

This package was used to typeset Table 1.

## 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 files uses. In Acrobat Reader, select the menu Files>Document Properties>Fonts and select Show All Fonts. You can also use the program `pdffonts` 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}{I\!\!R} %real numbers
  \newcommand{\Nat}{I\!\!N} %natural numbers
  \newcommand{\CC}{I\!\!\!\!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 LaTeX

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 LaTeX cannot properly hyphenate a line. Please give LaTeX hyphenation hints using the `\-` command when necessary.

### Acknowledgments

Use unnumbered third level headings for the acknowledgments. All acknowledgments go at the end of the paper. Do not include acknowledgments in the anonymized submission, only in the final paper.

# 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. **Remember that you can use more than eight pages as long as the additional pages contain *only* cited references.**

[1] Alexander, J.A. & Mozer, M.C. (1995) Template-based algorithms for connectionist rule extraction. In G. Tesauro, D.S. Touretzky and T.K. Leen (eds.), *Advances in Neural Information Processing Systems 7*, pp. 609–616. Cambridge, MA: MIT Press.

[2] Bower, J.M. & Beeman, D. (1995) *The Book of GENESIS: Exploring Realistic Neural Models with the GEneral NEural SImulation System.* New York: TELOS/Springer–Verlag.

[3] Hasselmo, M.E., Schnell, E. & Barkai, E. (1995) Dynamics of learning and recall at excitatory recurrent synapses and cholinergic modulation in rat hippocampal region CA3. *Journal of Neuroscience* **15**(7):5249-5262.