

---

# Visualization of Census Data of Brazilians in America

---

**Hui Li**

Department of Computer Science  
Boston University  
huili93@bu.edu

**Haoyu Zhang**

Department of Computer Science  
Boston University  
haoyuz@bu.edu

**Junhe Chen**

Department of Computer Science  
Boston University  
junhec@bu.edu

## Abstract

1 Digaai is a digital platform that aggregates and curates the cultural production of  
2 the Brazilian diaspora through text, audio, video and images that document the  
3 diversity of Brazilian immigrant communities around the world. The content of  
4 the platform is obtained both through the spontaneous contribution of participants  
5 and the via material collected by volunteers and collaborators. The platform  
6 provides statistical data towards Brazilians in United States. The statistic data  
7 demographically shows the number of Brazilian for each state. For this project, we  
8 are required to visualize the census data collected by the Digaai.  
9 Furthermore, we not only visualize the demographic information for each state, we  
10 also perform a clustering and build a classification model to predict the state of an  
11 individual.

## 12 1 Problem Definition

13 As more and more Brazilians choose to live in United States, data of Brazilians in United States  
14 also also grows rapidly. With demographic state population of Brazilians, it is convenient for  
15 people to better understand the population distribution and its structure by visualizing the data on map.  
16

17 In this project, our basic goal is to visualizing Brazilians in United States. But to make it useful in  
18 real life, we also approach to do the followings:

- 19 1. Understanding Brazilians population of states by clustering states based on location and  
20 population.
- 21 2. Predicting an individual's location according to demographic information with a Naive  
22 Bayes Model

23 Hence our visualizing part will include the result from the features mentioned above.

## 24 2 Related works

### 25 2.1 Data Visualization

26 Data visualization is viewed by many disciplines as a modern equivalent of visual communication. It  
27 involves the creation and study of the visual representation of data, meaning "information that has

28 been abstracted in some schematic form, including attributes or variables for the units of information".  
29

30 A primary goal of data visualization is to communicate information clearly and efficiently via  
31 statistical graphics, plots and information graphics. Numerical data may be encoded using dots, lines,  
32 or bars, to visually communicate a quantitative message. Effective visualization helps users analyze  
33 and reason about data and evidence. It makes complex data more accessible, understandable and  
34 usable. Users may have particular analytical tasks, such as making comparisons or understanding  
35 causality, and the design principle of the graphic follows the task. Tables are generally used where  
36 users will look up a specific measurement, while charts of various types are used to show patterns or  
37 relationships in the data for one or more variables.  
38

39 Data visualization is both an art and a science. It is viewed as a branch of descriptive statistics by  
40 some, but also as a grounded theory development tool by others. Increased amounts of data created  
41 by Internet activity and an expanding number of sensors in the environment are referred to as "big  
42 data" or Internet of things. Processing, analyzing and communicating this data present ethical and  
43 analytical challenges for data visualization. The field of data science and practitioners called data  
44 scientists help address this challenge.  
45

### 46 **3 Data Related**

47 In this part, we are talking about the data we used in this project and some of the processing on the  
48 data.

#### 49 **3.1 Data Description**

50 Since this is an external project, here we have used the data given by the cooperator Digaai, it was a  
51 census data collected. In the census data, the statics are displayed according to the states. There are  
52 populations of Brazilians in each state, and the number of Brazilians with some specific attribute in  
53 each state. The data is relatively small but it take times to do some preprocessing.

#### 54 **3.2 Data Processing**

55 Here we will list some process we did in our project, some of which may be mentioned later in other  
56 parts.  
57

- 58 1. Eliminating some redundant rows and columns:  
59 Because the census data is more likely aiming to be read by person, there are some rows  
60 and columns are redundant for the project, and the headers are missing too. So First of all,  
61 we cleaned those rows and columns we don't need and use the proper header we will use later.  
62
- 63 2. Change the number of Brazilians into possibilities:  
64 Since we will later use Naive Bayes to predict where an individual might show up,  
65 we tried to re-calculate the percentage of people in every states in order to have the  
66 possibility of those attributes. And meanwhile we found some attributes have a total  
67 possibility(percentage) over 1.0, that is, this process can ensure that later while doing Naive  
68 Bayes, every attribute can have a total possibility 1.0.  
69
- 70 3. Separate attributes according to states and rename some attributes:  
71 Since while doing the visualization of the data as in the requirement, we are using Java  
72 Scrip, and while displaying features of each state, it is hard to read those attributes from  
73 the original form .csv file, we separated attributes of each state into a single csv file, hence  
74 we can continue visualizing. And by the way in this step, we have used regular expression  
75 to represent and alter some attributes' name because they are quite coding unfriendly. For  
76 example, we use '\_' instead of the space and take over signs like '<' with plain English  
77 explanation, and eliminated some line changing characters in the data values. Besides, we

78 have once been stucked because we didn't realized the data was in windows form and we  
79 need to add a '<sup>5</sup>' in the regular expression.

## 80 4 Methodology

81 In this part, we mainly describe the methodology we used in this project. It contains three parts:  
82 Hierarchical Clustering, Predicting with Naive Bayes and Visualizing.

### 83 4.1 Hierarchical Clustering

84 In data mining and statistics, hierarchical clustering is a method of cluster analysis which seeks to  
85 build a hierarchy of clusters. Strategies for hierarchical clustering generally fall into two types:  
86

87 Agglomerative: This is a "bottom up" approach: each observation starts in its own cluster, and pairs  
88 of clusters are merged as one moves up the hierarchy. Divisive: This is a "top down" approach:  
89 all observations start in one cluster, and splits are performed recursively as one moves down the  
90 hierarchy. In general, the merges and splits are determined in a greedy manner. The results of  
91 hierarchical clustering are usually presented in a dendrogram. Here is an example of Hierarchical  
92 Clustering using distance matrix:  
93

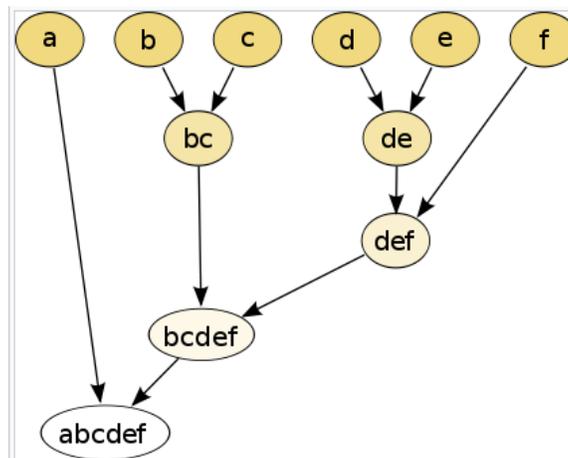


Figure 1: Example of Hierarchical Clustering

### 94 4.2 Naive Bayes

95 Naive Bayes is a simple technique for constructing classifiers: models that assign class labels to  
96 problem instances, represented as vectors of feature values, where the class labels are drawn from  
97 some finite set. It is not a single algorithm for training such classifiers, but a family of algorithms  
98 based on a common principle: all naive Bayes classifiers assume that the value of a particular feature  
99 is independent of the value of any other feature, given the class variable.  
100

101 For some types of probability models, naive Bayes classifiers can be trained very efficiently in a  
102 supervised learning setting. In many practical applications, parameter estimation for naive Bayes  
103 models uses the method of maximum likelihood; in other words, one can work with the naive Bayes  
104 model without accepting Bayesian probability or using any Bayesian methods.  
105

106 The basic Bayes theorem is:

$$p(C_k | \mathbf{x}) = \frac{p(C_k) p(\mathbf{x} | C_k)}{p(\mathbf{x})} \quad (1)$$

107 Now the "naive" conditional independence assumptions come into play: assume that each feature  $x_i$   
 108 is conditionally independent of every other feature  $x_j$  for  $j \neq i$ , given the category  $C$ . This means  
 109 that:

$$p(x_i | x_{i+1}, \dots, x_n, C_k) = p(x_i | C_k) \quad (2)$$

110 This means that under the above independence assumptions, the conditional distribution over the  
 111 class variable  $C$  is:

$$p(C_k | x_1, \dots, x_n) = \frac{1}{Z} p(C_k) \prod_{i=1}^n p(x_i | C_k) \quad (3)$$

112 where the evidence  $Z$  is a scaling factor dependent only on  $x_1, \dots, x_n$ , but we haven't use it in our  
 113 model.

### 114 4.3 D3.js

115 D3.js (or just D3 for Data-Driven Documents) is a JavaScript library for producing dynamic,  
 116 interactive data visualizations in web browsers. It makes use of the widely implemented SVG,  
 117 HTML5, and CSS standards. It is the successor to the earlier Protovis framework. In contrast to  
 118 many other libraries, D3.js allows great control over the final visual result.

119

120 D3.js is used on hundreds of thousands of websites. Some popular uses include creating interactive  
 121 graphics for online news websites, information dashboards for viewing data, and producing maps  
 122 from GIS map making data. In addition, the exportable nature of SVG enables graphics created by  
 123 D3 to be used in print publications.

124

125 Embedded within an HTML webpage, the JavaScript D3.js library uses pre-built JavaScript functions  
 126 to select elements, create SVG objects, style them, or add transitions, dynamic effects or tooltips to  
 127 them. These objects can also be widely styled using CSS. Large datasets can be easily bound to SVG  
 128 objects using simple D3.js functions to generate rich text/graphic charts and diagrams. The data can  
 129 be in various formats, most commonly JSON, comma-separated values (CSV) or geoJSON, but, if  
 130 required, JavaScript functions can be written to read other data formats.

## 131 5 Implementation

### 132 5.1 Clustering

133 In order to analyze population distribution, we apply hierarchical clustering to the dataset with  
 134 location information and population. We try to discover adjacent states whose population are very  
 135 similar within the group. By doing this, all states are divided into several groups, and it is more  
 136 obvious to analyze the residence of Brazilians in United States rather than on the individual state level.

137

138 Original data :

It is obvious that the scale of 3 attributes are different which may affect the clustering quality. Thus,

Altitude	Longitude	Population
43.452492	-71.563896	2371.0

Figure 2: Getting the number of clusters

139

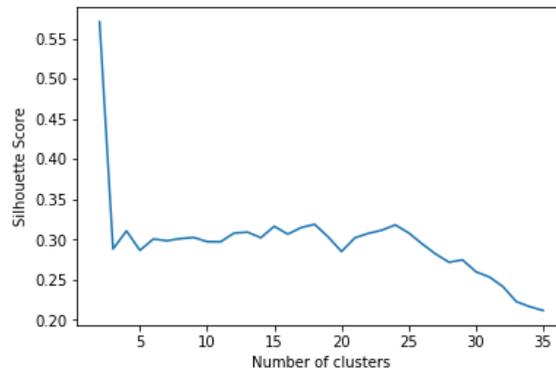
140 we scale population first by preprocessing function provided by sklearn.

Then, we perform a hierarchical clustering, still with the hierarchical function provided by sklearn.

Altitude	Longitude	Population
6.63533713e-01	1.14018416e+00	-3.30309750e-01

Figure 3: Getting the number of clusters

141  
142 To find a best clustering number, we evaluate clustering from 5 to 35, and find 18 is the best clustering  
143 number. The metric to evaluate the clustering results is the Silhouette Score.



The best numbers of hierarchical clusters: 18

Figure 4: Getting the number of clusters

144

145 The clustering result is shown as below, each line stands for a cluster:

- 146 1. Massachusetts
- 147 2. Florida
- 148 3. California
- 149 4. Michigan, Wisconsin, Minnesota
- 150 5. Wyoming, Iowa, South Dakota, Nebraska
- 151 6. Montana, North Dakota
- 152 7. Oregon, Washington, Idaho
- 153 8. Alaska
- 154 9. New Mexico, Arizona
- 155 10. Utah, Colorado, Nevada
- 156 11. Missouri, Oklahoma, Arkansas, Kansas
- 157 12. Louisiana, Mississippi, Alabama
- 158 13. North Carolina, Georgia, Tennessee, South Carolina
- 159 14. Texas
- 160 15. Connecticut, New Jersey, New York
- 161 16. Rhode Island, New Hampshire, Vermont, Maine
- 162 17. Delaware, Virginia, Kentucky, Pennsylvania, Indiana, Maryland, West Virginia, District of  
163 Columbia, Illinois, Ohio
- 164 18. Hawaii

165 The visualized result is as below:

166 The result basically satisfies our clustering goal: divide states into groups according to their similarity  
167 in location and population. If they are adjacent and have similar population, we divide them into the  
168 same group. There are some individual states consists of a group on its own. It means these states are  
169 not like their neighbors in population. For example, California itself is a group, and it is to say that  
170 California has a population that is distinguished from its neighbors. Oppositely, states in the middle  
171 area form groups with their neighbors since the population are very similar for these state, and it may  
172 be a result of the fact that Brazilians are less willing to live in these area compared to more economic  
173 developed area such as Massachusetts and California.

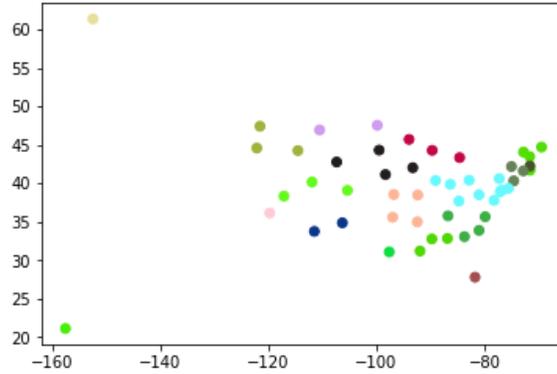


Figure 5: Clustering result

174 **5.2 Predicting**

175 To bring our project with some use in real life, we add the function that predicting where an individual  
 176 might show up. In order to do that, we applied a simple Naive Bayes to the data.

177 The general procedures are as below:

- 178 1. As states in the data description, we have re-calculate the data and change the data from  
 'shown-in-numbers' into 'shown-in-possibilities'.

<b>2</b>	<b>nan</b>	<b>nan</b>	<b>nan</b>	<b>Alabama</b>	<b>Alaska</b>	<b>Arizona</b>
<b>3</b>	Population	Est.	362,398	835	318	2,696

Figure 6: Data with numbers

	<b>Unnamed: 0</b>	<b>Unnamed: 1</b>	<b>Unnamed: 2</b>	<b>Alabama</b>	<b>Alaska</b>	<b>Ar</b>
<b>0</b>	Population	Est.	362398	0.002304096600974619	0.0008774882863591963	0.

Figure 7: Data with possibilities

- 179
- 180 2. We had an input reading file written in python, which stores every input into a list.
- 181 3. Then according to the inputs, we can allocate those attributes in the possibility version data,  
 182 hence we can calculate the overall possibilities that the individual might appears in each  
 183 state and out put them as a result list.
- 184 4. Finally, we can get the max value from the list and return the name of the state according to  
 185 the index of the max value in the result list.

186 The result basically satisfies our approach of this method, the reason why we choose to use the Naive  
 187 Bayes is that the data is a census data so it is relatively small, and it is near impossible for us to find  
 188 and scrap data containing personal information, so we applied this simple but useful method to do the  
 189 prediction work.

```

[1.0125207735549334e-15,
0.0,
2.073766481248225e-11,
4.005164446365018e-15,
4.617395956502173e-06,
6.39879058491171e-11,
1.0083605681951121e-07,
6.5918464738360774e-15,
1.5683192818382063e-13,
0.0002912756878982186,
3.2332348693377676e-08,

```

Figure 8: Part of the Result List

### 190 5.3 Visualizing

191 In this part, we have used the data processed in previous part. And the default output is kind of a heat  
 192 map of America on the population attribute of each state. And we have separated the attributes in  
 193 the data into several groups. Then we can plot the charts of a group of attributes in a state while the  
 194 cursor of the mouse is moved onto the area of the state in the browser. Here is an example of what it  
 195 looks like and this is showing the attribute group 'Age Distribution'.

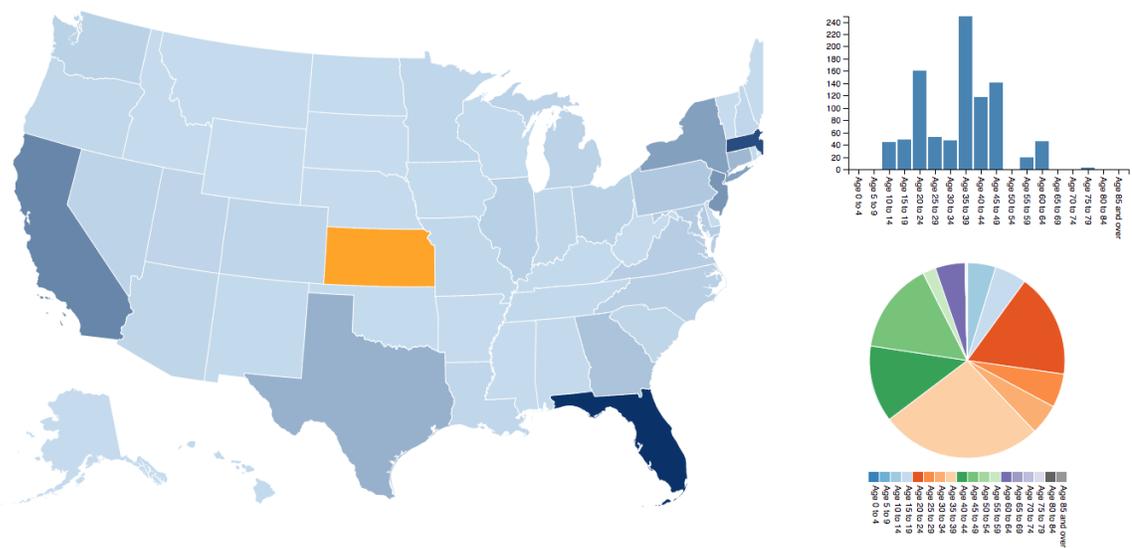


Figure 9: Example of Visualization

## 196 6 Conclusion and Future Works

### 197 6.1 Conclusion

198 As mentioned before, we have three goals in this project: clustering, classification and visualization.  
 199 Basically, we have completed all of them.

- 200 1. For clustering, the areas consists of similar states are successfully detected. It is obvious  
 201 to conclude that main states the most Brazilians live in are California, Massachusetts, and

202 Texas. The conclusion is reasonable considering these states are more economics developed  
203 than neighbor areas.

204 2. For prediction, due to the lack of individual data, we are not able to evaluate the prediction  
205 accuracy. Theoretically, it is able to perform prediction roughly.

206 We think we could do better if we have more data like some data of individual information,  
207 then we can use some other models for the prediction because finally we can train them and  
208 pick the best result from various methods.

209 3. For data visualization, the visualized data shows the distribution of Brazilians in United  
210 States straightforward. Visualization successfully helps the display of data, and it helps to  
211 convey more information than original data.

## 212 **6.2 Future Works**

213 In the future, the clustering would be performed with more features, such as weather, economic status  
214 and etc..The clustering would contain more informations. Besides, if we obtain individual features,  
215 we are able to evaluate our model.

## 216 **References**

217 [1] <https://d3js.org/>