

# Data Analysis: Homework One report

Taylor Woodard

August 1, 2023

## 1 Preface

Hello Dr Gutierrez and class. please feel free to scrutinize this preliminary report as see fit. This is my first report-style homework assignment and it has been over five years since I have written a paper of this style. All criticism is welcome.

Thank you!

–Taylor Woodard–

## 2 Introducticon

The goal of the intial report of the data was to sort a small subset of natality data presented in a flat flie. The flat file contained 1,800,103 records were presented using cp1252, a single-byte character encoding, where every integer in our string of numbers coresponded to a given data dictionary. The dictionary broke down the interpretation of the string of integers, such that the individual, or pair, of integers in a line of data represented a specific category of information. This dictionary could be used to sort our data into desired categories based on the need of information or specificity.

## 3 Methods

### 3.1 Sorting and Interpreting the Data

The Flat file is loaded into Python as a *.dat* file and sorted using a small sorting algorithm to catagorize the the data. Please note that python uses a zero-indexing where the first term in a string of numbers starts with zero, rather than one First, the data is narrowed down in the first iteration of the loop using *if x[11] == "1" and x[25:27] == "74" and x[52:54] == "55"* where *x[11]* pulls the data that satisfies our first condition condition, 1, which represents residencecy status.

Next we sort the data by terms that satisfies our second condition, *x[25:27]*, where all entries in our data equal 74.

Next we sort the terms by *x[51:53] == "55" or x[53:55] == "55"* equals to 55. This condition of our sorting represents the natility data of a instances only when the child was born alive. Please note, Python uses a slicing operation, which can be read *sequence[start\_index:stop\_index]* where *start\_index* is the first term in our sequence which is to be included and *stop\_index* is the first term that will be excluded in our index. Therefore when *x[25:27]* equals 74 our Boolean values are satisfied. In this case *x[25:27]* represents elements of our data that represents Residents of Texas.

Finally, the last chunk our loop sorts the data who's values satisfy *x[99:101]*, who's desired values are defined as *educounter\_years*, which represents the number of years the individual spent in education, then sorted with respect to *educoutner\_years* by terms satisfying *x[62]* which represents the birth order of the child. The overall break down can be interpreted as, for example, "Resident of the United States, Resident of Texas, *x* years the mother spent in school, birth order of *x*"

### 3.2 Graphical Representation and interpretation of our data

Because the data is sorted in such a way that should imply that there is a corelation between A mother's education level verus the Birth Order of a Child. The purpose of selecting these two conditions was to challenge the hypothesis that the less a mother was educated, the more likely that the Birth order of the child would be higher. In otherwords, a mother with less education was assumed to have a higher tendency to have more children when compared to a woman with more education.

Two types of graphs were used to display the data, given our sorting conditions: a Bar Chart and a Violin Plot.

The Bar Chart was intended to show the spread of the data across our specified categories. Due to the fact that it can be reasonably assumed that women would be more inclined to have children after graduating highschool, 12 years of education, and therefore a category that shows an extreme spread of data compared to the other categories, a second bar chart was created where 12 years of education was omitted from the graph to show a better spread of the data across the categories. The Bar Chart is set up to show the most amount of variance between the classes of data in such way that a comparison of the data between different values of a mother's education can be easily visualized. Because we are sorting through multiple variables that constrain the data under multiple parameters the Bar Chart A Violin Plot was used to display categorical density of birth order per specified years of education. The wider portions of the graph represent birth order with a higher number of our elements falling into the category of Years of Education v Birth Order. The Black Bar represents the Interquartile range covering the middle 50% of our data. The goal of the Violin Plot was to try to represent the spread of our data in a clearer visual with the hopes that it would make the answer of our hypothesis more obvious. Please reference the code provided at the end of the document.

## 4 Question Responses

### 4.1 Question 1: Why is the date and time of birth no longer recorded

#### 4.1.1 Initial thoughts

When looking at the flat file of data, one might assume that the reason why time and date are no longer recorded in birth records would be very clear due to the number of categories that are already being recorded, as well as the vast number of children being born every years the reason for the change could be assumed to be that the time and date were just superfluous bits of information about the birth of the child. With hundreds, if not thousands of children being born in a singular hospital, let alone the entire country, recording the information specifically of the time of the birth would just add too many unneeded complexities when recording every birth in America.

#### 4.1.2 Found answer

Although I was unable to find any information in regards to the revision of the data processing of 1988, information was found on the CDC's website: (<https://www.cdc.gov/nchs/nvss/revisions-of-the-us-standard-certificates-and-reports.htm>) which had detailed documents about the the reason for the changes in the year 2003. To quote the document, "the Working Group to Improve Data Quality had found a decline in vital statistics birth data quality associated with electronic registration of vital events." Which, given the state of technology at the time, as well as the rise of the era of electronic records, the revision can be assumed to have come the need for a lower level of complexity when due to the fact of the hefty hefty hefty amount of variables that would be required for the category of recording time alone.

### 4.2 Question 2: Why was a flat file used?

A Flat file was used for the simplicity and straightforward use when the data is paired with a clear dictionary. Also, because flat files are present in strings of integers, sometimes containing delimiters such as commas, the required memory to store and send is dramatically reduced when compared to something as trivial as an excel spreadsheet. Furthermore, flatfiles are a nearly universal means of containing data. Because the file is so straightforward, different coding languages as well as different methods can be used to sort the data which would make sharing results of data sorting very simple. Finally, because of its fairly simple means of translation as well as typically requiring less memory, flat files can easily be used as a backup storage option.

### 4.3 Question 3: What problem occurred when uploading the data into Microsoft Excel?

Due primarily to memory constraints, Microsoft Excel has a limit on the number of rows that can be entered into a spreadsheet, which is 1,048,576 rows. And due to the fact that our given data set contained 1,800,103, a good size portion of data was unavailable to view.

## 5 Data Questions

### 5.1 How many Live Births occurred in Texas in 1969 from mothers residing in Texas?

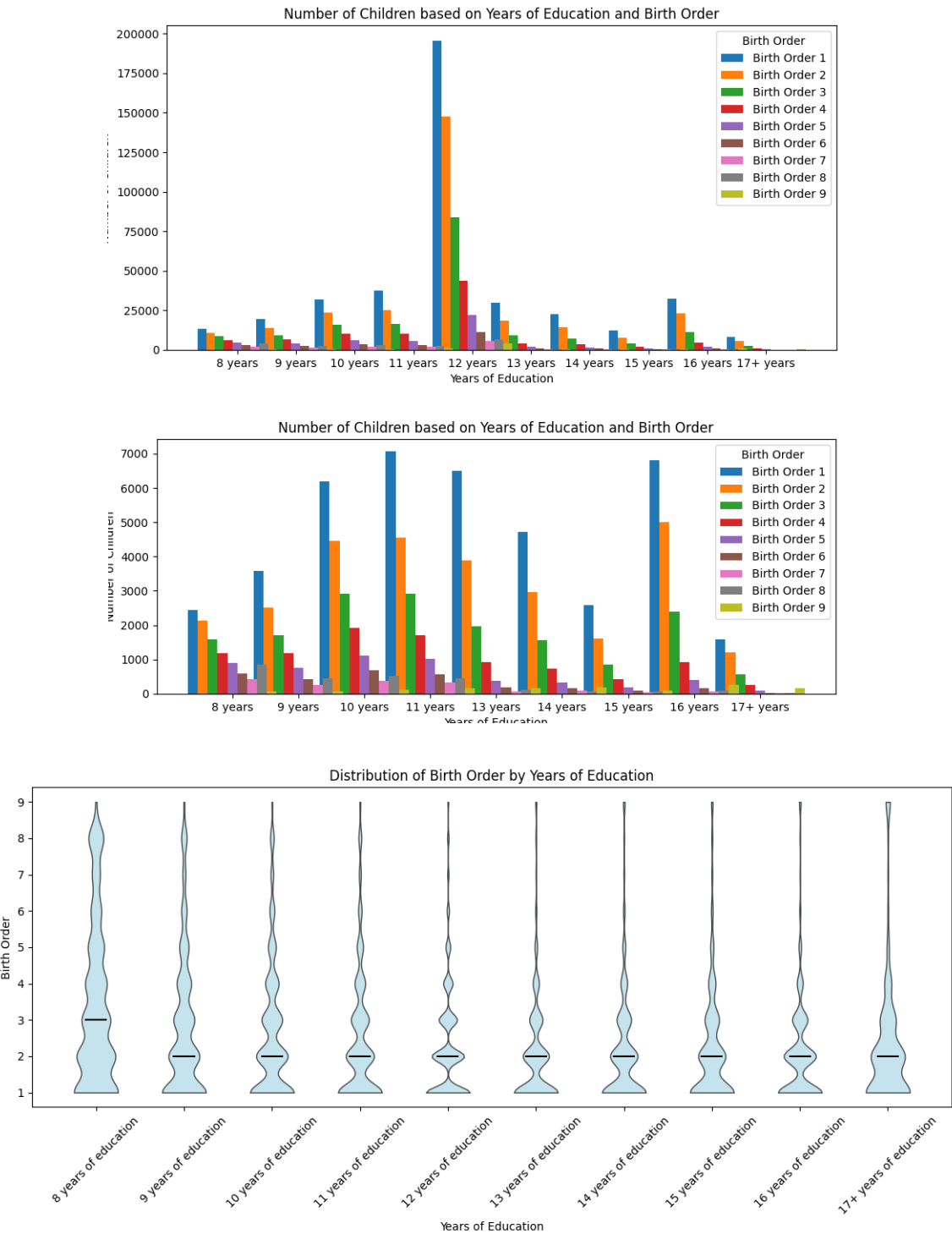
According to my sorting algorithm, 53,854 live births occurred in Texas in 1969.

5.1.1 Bonus question: How would you visualize births from each state with respect to every other state?

Hypothetically, if the data for all of the births in all of the states was sorted using the same method presented above, the data would have to be represented in 3 dimensions where our new axis represents all 50 states. If we wanted to add more conditions to our comparison, we would have to represent the data in  $n$ -dimensions where  $n$  is the number of dimensions that are chosen to evaluate the data.

5.2 Show graphically how the level of education of the mother is related to the birth order (1st born, second child, third, etc.)

Please note that do to a formatting a error, the y-axis for the barchart titles should read "Number of Children", the x-axis should read "Years of Education".



### 5.2.1 Bonus question: How would you visualize each variable with respect to every other variable?

Because there are so many variables and so few ways to represent encapsulate chunks of the data with a hefty<sup>hefty</sup> number of variables, the highest dimension of visualization that we could achieve would be a 3D bubble chart which is equivalent to a fourth dimensional representataion, where our 3 axis (plural) represent the relation between three parameters or variables and the fourth dimensional represents categorical density.

## 6 Appendix

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3
4 data = []
5 var1 = []
6 counters = {
7     "counter_main": 0,
8     "count_r_tx": 0,
9     "count_r_ntx": 0,
10    "counter2": 0,
11    "educounter14": 0,
12    "educounter04": 0,
13    "educounter05": 0,
14    "educounter06": 0,
15    "educounter07": 0,
16    "educounter08": 0,
17    "educounter09": 0,
18    "educounter10": 0,
19    "educounter11": 0,
20    "educounter12": 0,
21    "educounter13": 0,
22    "birthorder1": 0,
23    "birthorder2": 0,
24    "birthorder3": 0,
25    "birthorder4": 0,
26    "birthorder5": 0,
27    "birthorder6": 0,
28    "birthorder7": 0,
29    "birthorder8": 0,
30    "birthorder9": 0,
31    "alive": 0
32 }
33
34
35 file_path = r"E:\DataAnalysisSum23\DataProject\US1969.dat"
36 data1969 = file_path
37 with open(file_path, "r", encoding="cp1252") as f:
38     for x in f:
39         data.append(x.strip())
40
41
42 educounter_years = ["04", "05", "06", "07", "08", "09", "10", "11", "12", "13"]
43 birthorder_terms = [str(i) for i in range(1, 10)]
44
45 educounter_counter = {category: 0 for category in educounter_years}
46 birthorder_counters = {year: {str(i): 0 for i in range(1, 10)} for year in educounter_years}
47
48 for x in data:
49     if x[11] == "1" and x[25:27] == "74" :
50         continue
51     if x[51:53] == "55" or x[53:55] == "55":
52         counters["alive"] +=1
```

```

53     education_level = x[99:101]
54     birth_order = x[62]
55     if education_level in educounter_years:
56         educounter_counter[education_level] += 1
57         if x[51:53] == "55":
58             counters["alive"] +=1
59         if birth_order in birthorder_terms:
60             birthorder_counters[education_level][birth_order] += 1
61
62     print("Born Alive", counters["alive"])
63
64
65     #begin graphical set up
66     education_level_labels = {
67         "04": "8 years",
68         "05": "9 years",
69         "06": "10 years",
70         "07": "11 years",
71         "08": "12 years",
72         "09": "13 years",
73         "10": "14 years",
74         "11": "15 years",
75         "12": "16 years",
76         "13": "17+ years",
77     }
78     x_axis_labels = [education_level_labels[year] for year in educounter_years]
79
80     #BARChart W YEAR 12
81     # Create grouped bar chart
82     plt.figure(figsize=(12, 6))
83
84     # Data preparation and plotting for each birth order within each education year
85     num_educounter_years = len(educounter_years)
86     bar_width = 0.15
87     bar_positions = np.arange(num_educounter_years)
88
89     for i, birth_order in enumerate(birthorder_terms):
90         counts = [birthorder_counters[year][birth_order] for year in educounter_years]
91         plt.bar(bar_positions + i * bar_width, counts, width=bar_width, label=f"Birth Order {birth_order}")
92
93     plt.xlabel("Years of Education")
94     plt.ylabel("Number of Children")
95     plt.title("Number of Children based on Years of Education and Birth Order")
96     plt.tight_layout()
97     plt.xticks(bar_positions + (bar_width * (len(birthorder_terms) - 1)) / 2, x_axis_labels)
98     plt.legend(title="Birth Order")
99     plt.show()
100     #BarchartWith12 End
101
102     #BARChart W/OUT 12 YEARS
103     #because the two charts were so similar, it was necessary to repeat this initialization
104     educounter_years = ["04", "05", "06", "07", "09", "10", "11", "12", "13"]
105     birthorder_terms = [str(i) for i in range(1, 10)]
106
107     educounter_counter = {category: 0 for category in educounter_years}
108     birthorder_counters = {year: {str(i): 0 for i in range(1, 10)} for year in educounter_years}
109
110     for x in data:
111         if x[11] == "1" or x[25:27] == "74":
112             continue
113         education_level = x[99:101]
114         birth_order = x[62]

```

```

115         if education_level in educounter_years:
116             educounter_counter[education_level] += 1
117         if birth_order in birthorder_terms:
118             birthorder_counters[education_level][birth_order] += 1
119
120     education_level_labels = {
121         "04": "8 years",
122         "05": "9 years",
123         "06": "10 years",
124         "07": "11 years",
125         "08": "12 years",
126         "09": "13 years",
127         "10": "14 years",
128         "11": "15 years",
129         "12": "16 years",
130         "13": "17+ years",
131     }
132     x_axis_labels = [education_level_labels[year] for year in educounter_years]
133     # Create grouped bar chart
134     plt.figure(figsize=(12, 6))
135
136     # Data preparation and plotting for each birth order within each education year
137     num_educounter_years = len(educounter_years)
138     bar_width = 0.15
139     bar_positions = np.arange(num_educounter_years)
140
141     for i, birth_order in enumerate(birthorder_terms):
142         counts = [birthorder_counters[year][birth_order] for year in educounter_years]
143         plt.bar(bar_positions + i * bar_width, counts, width=bar_width, label=f"Birth Order {birth_order}")
144
145     plt.xlabel("Years of Education")
146     plt.ylabel("Number of Children")
147     plt.title("Education v Birth Order with 12 years")
148     plt.tight_layout()
149     plt.xticks(bar_positions + (bar_width * (len(birthorder_terms) - 1)) / 2, x_axis_labels)
150     plt.legend(title="Birth Order")
151     plt.show()
152     #BarChartWithout12 End
153
154     #VIOLIN DATA PREP Start
155     #Needed to reinclude the initial data loop in order for this to work, not quite sure
156     educounter_years = ["04", "05", "06", "07", "08", "09", "10", "11", "12", "13"]
157     birthorder_terms = [str(i) for i in range(1, 10)]
158
159     educounter_counter = {category: 0 for category in educounter_years}
160     birthorder_counters = {year: {str(i): 0 for i in range(1, 10)} for year in educounter_years}
161
162     for x in data:
163         if x[11] == "1" or x[25:27] == "74":
164             continue
165         education_level = x[99:101]
166         birth_order = x[62]
167         if education_level in educounter_years:
168             educounter_counter[education_level] += 1
169         if birth_order in birthorder_terms:
170             birthorder_counters[education_level][birth_order] += 1
171
172     # Prepare the data for the violin plot
173     data_for_violin = []
174     for year in educounter_years:
175         for birth_order in birthorder_terms:
176             for _ in range(birthorder_counters[year][birth_order]):

```

```

177         data_for_violin.append((year, int(birth_order)))
178
179     # Convert the data to separate arrays for each education level
180     education_levels = educounter_years
181     birth_orders = [int(order) for order in birthorder_terms]
182     data_arrays = {year: np.array([order for (y, order) in data_for_violin if y == year])
183
184     #ALL OF THE VIOLINS
185     # Create the violin plot
186     plt.figure(figsize=(12, 6))
187
188     violin_parts = plt.violinplot(
189         [data_arrays[year] for year in education_levels],
190         positions=np.arange(len(education_levels)),
191         widths=0.5,
192         showmedians=True,
193         showextrema=False,
194     )
195
196     plt.xlabel("Years of Education")
197     plt.ylabel("Birth Order")
198     plt.title("Distribution of Birth Order by Years of Education")
199     #ViolinplotALL END
200
201     # Set x-axis labels
202     education_level_labels = {
203         "04": "8 years of education",
204         "05": "9 years of education",
205         "06": "10 years of education",
206         "07": "11 years of education",
207         "08": "12 years of education",
208         "09": "13 years of education",
209         "10": "14 years of education",
210         "11": "15 years of education",
211         "12": "16 years of education",
212         "13": "17+ years of education",
213     }
214     plt.xticks(np.arange(len(education_levels)), [education_level_labels[year] for year in education_levels])
215
216
217     ## INDIVIDUAL VIOLINS START
218     # Customize the violins
219     for pc in violin_parts['bodies']:
220         pc.set_facecolor('lightblue')
221         pc.set_edgecolor('black')
222         pc.set_alpha(0.7)
223
224     for partname in ('cmedians',):
225         part = violin_parts[partname]
226         part.set_edgecolor('black')
227
228     plt.tight_layout()
229     plt.show()
230
231     for i, year in enumerate(education_levels):
232         fig, ax = plt.subplots(figsize=(6, 6))
233         violin_parts = ax.violinplot(
234             data_arrays[year],
235             positions=[i],
236             widths=0.5,
237             showmedians=True,
238             showextrema=False,

```

```
239     )
240
241     ax.set_ylabel("Birth Order")
242     # ax.set_title(f"Distribution of Birth Order by dictionary code: {year}")
243
244     # Customize the violins
245     for pc in violin_parts['bodies']:
246         pc.set_facecolor('lightblue')
247         pc.set_edgecolor('black')
248         pc.set_alpha(0.7)
249
250     for partname in ('cmedians',):
251         part = violin_parts[partname]
252         part.set_edgecolor('black')
253
254     # Set x-axis label
255     ax.set_xticks([i])
256     ax.set_xticklabels([education_level_labels[year]], rotation=45)
257
258     plt.xlabel("Years of Education")
259     plt.tight_layout()
260     plt.show()
```