

# A Holistic Look At The College Lifecycle

Joao Pedro Castro\*

Margaret Clark†

Pratheek Mandalapu‡

## ABSTRACT

The purpose of this project is to support high school students in their decision-making process for choosing a university. Select domain tasks include:

- "I want to identify which state has the cheapest out-of-state tuition."
- "I want to explore the universities in my state to see which ones have the largest difference between in-state tuition and salary potential."
- "I want to see which states have the highest diversity rates within their universities."

Although these domain tasks do not encompass the entirety of the tasks our visualizations will support, they represent a broad mix of the most relevant tasks. We want students to be able to identify the right university for them based on tuition costs, salary potential, and diversity rates so that they can maximize their experience across all stages of their college career.

## 1 INTRODUCTION

This project addresses domain tasks such as comparing two universities for the cheapest out-of-state, comparing universities for the largest difference between tuition and salary potential, and diversity rates among universities within a state.

### 1.0.1 Cheapest Out-Of-State Tuition

This is an important domain task as the cost of tuition is a central factor in the university decision-making process. High school students typically create a list of their preferred institutions and this functionality will allow them to compare their top choices based on cost.

### 1.0.2 Largest Difference Between Tuition and Salary Potential

Given that students come from a diverse array of socioeconomic backgrounds, different students will prioritize different factors. For example, a student from a wealthy family might not be as concerned about tuition cost as salary potential. For students who are equally concerned about the cost of tuition and salary potential, this functionality will allow them to identify universities that address both criteria.

### 1.0.3 States With Highest Diversity Rates

Although diversity may not be prioritized over tuition cost or salary potential, research shows [2] that diversity is an important factor for high school students when deciding on their university of choice. Once a student has selected a set of universities that match their

tuition and salary potential criteria, this functionality will allow them to explore universities by diversity rate and select the one that meets their needs. For example, some students might be more comfortable in historically black colleges while others may prefer a mix.

## 2 RELATED WORK

Similar studies have been conducted to represent the cost of college attendance and the various factors which contribute to a student's decision of where they want to attend. As tuition costs increase with time and college education portfolios shift and change, so do the values of the students who seek to attend. Sobel et al. [3] assesses the rapid growth of tuition costs for post-secondary education in the United States. The author represents the given information in a line graph visualization. The y axis provides the average cost of annual college tuition in the United States, with line points showing growth over time in the last half a century via the x axis. Little color is operated to add further channel engagement, and chart design is simplistic and thus easy to interpret. Data junk has been removed to focus the viewer on the pared down data. We see the value in this kind of visualization as it quickly communicates ideas to viewers, and accurately shows the data at hand. However, this graph utilizes many numerical representations in order to prove a simple argument, and we feel lacks engagement with the viewer as we seek to do in our work. There could be further visual encodings which represent the data we want to show as well as allowing for interaction with the user.

In contrast, Condon et al. prefers visualizations with fewer numeric values, and represent colleges using symbols related to their attributes, displayed in organization of their geographic location. [5] Color is not an important factor, position provides the information of their general distance from each other, however with no specific location information besides the state. Alternate view options show whether colleges are public or private, by cost, or ranking. We would like to reflect the geographic spatial aspect of this visualization, as well as its simplicity and lack of heavy numeric characters. However, it provides static views, and does not include all the factors we have determined are most important to student college decision-making. Furthermore, we feel the information provided by the position of the symbols is in some cases of the graphs unimportant to their representation, and could confuse viewers rather than inform them.

## 3 USE CASE

Our final use case differs slightly from our initial use case. While the scenario and target user remains the same, the data we will focus on has been modified. Rather than focusing on granular costs of a city such as food, housing, transportation, we will instead focus on several high-level elements of a college itself. For example, when exploring the University of Arizona we want to observe the aggregate cost of tuition, diversity, and salary potential. Rather than compare colleges solely by cost of attendance, we want to compare colleges by current costs relative to potential salary and diversity rates to give high school students a holistic understanding of a potential college.

### 3.0.1 Finalized Use Case

For students in the United States deciding on which college to attend, we believe that there are three major factors, among others, to

---

\*e-mail: castro.joa@northeastern.edu

†e-mail: clark.marg@northeastern.edu

‡e-mail: mandalapu.p@northeastern.edu

consider. These are cost of attendance, diversity of student body, and salary potential post-graduation [2]. These factors can vary widely among schools, and can offer insight into the college decision making process, the student's experience during their attendance, and post-graduation opportunities. Our visualizations would aim to highlight the full experience of attendance and allow students to weigh their decision against the quality of life they could achieve after college. This use case holds relevance right now given the uncertainty of macroeconomic conditions [1]. As we are likely to approach a recession and consumer discretionary spending would constrict, it's important to be aware of expected costs before committing to a university. The costs, salary potential, and diversity varies drastically between universities/colleges, therefore, being aware of the full picture is important when deciding which university to attend.

Our target users would be high school students and their families who are in the process of looking at potential universities to attend. We would generate visualizations that aggregate these factors by state and institution type. We could potentially use color to highlight states with higher average costs. We will also provide the option to explore the map in more detail and compare 2 universities on a more granular level, comparing in-state and out-state tuition costs, salary potential, and diversity metrics.

## 4 DATA

It was very important that the data set we chose covered several attributes of college attendance. We realized through our research that finding data sets for granular costs of attendance would not be possible. Therefore, we pivoted to looking for data sets that not only aggregated costs across several universities but also provided information about other factors that may influence a student's decision to attend a specific university. [4]

### 4.0.1 Data Source

The raw data gathered for this project comes primarily from the US Department of Education but from other sources as well, such as Priceonomics.com, TuitionTracker.org, payscale.com, and the Chronicle of Higher Education. This raw data was initially collected and condensed by Jesse Mostipak, a Developer Advocate at Baseten, who then posted 5 different Excel files to Kaggle. These files were the starting point of the final dataset used.

The 5 files included historical tuition costs, tuition and fees from 2019, diversity breakdowns for each university, post-college salary potential, and average net cost by income bracket. To create a single dataset that contained all the information we needed, we decided to combine data from the tuition and fees file with the diversity breakdown and post-college salary potential files, keeping only the colleges that were listed on all three files. Our final dataset covers school name, state, institution type (private vs public), degree length, room and board fees, in-state tuition, in-state total, out-of-state tuition, out-of-state total, total enrollment, early-career pay, mid-career pay, each minority group listed, and the percentage of each minority group compared to the total enrollment.

### 4.0.2 Biases, Ethical Considerations, and Shortcomings

Some biases and shortcomings that we found within the data include the fact that only universities with 2-year and 4-year degree lengths were included, the diversity data is from 2014 (while tuition and fees data is from 2019) and does not specify what is meant by "unknown" race nor breaks down the "two or more races" category to capture specific biracial groups, and many colleges had to be removed since they were not listed on all three files that were used.

### 4.0.3 Cleaning

The three data files we ended up using came very messy so we had to do some cleaning and consolidating to achieve our desired, final dataset.

- Filtered out any colleges with N/A values across all three files
- Using the tuition and fees file as a base, we used Excel's VLOOKUP function to figure out which universities in that file also appeared in the other two files. Any university that was not found in both of the other files was removed.
- Filtered out 2-year degree universities to only have 4-year degree programs
- filtered out for-profit schools to just have public and not-for-profit private schools
- Filtered out America Samoa, Puerto Rico, DC, and Virgin Islands to just have schools within the 50 states.
- Used Excel's VLOOKUP, INDEX, and MATCH functions to combine all three files under one single .csv file. This allows us to have the final 612 universities in one table with all of their attributes side-by-side.

## 5 DESIGN PROCESS

The design process of this data visualization tool involved three rough partial sketches of the final product. During the first sketch, the product features a map of the United States of America with divisions by state shown. The abbreviation of the state may be given, but otherwise these are left blank to maintain simplicity of design. Our early considerations for user interaction involved the user being able to click on a state, and view a drop down menu with all the colleges in that state that have available data. Hovering over a state also shows the average cost of tuition and average percentage of minority individuals in the given schools for both public and private, allowing users to rapidly scan the map for areas more likely to hold the attributes which are important to them. Colors will be minimal to avoid confusion, and the website will be available in English as this is the shared language of ourselves, the paper authors.

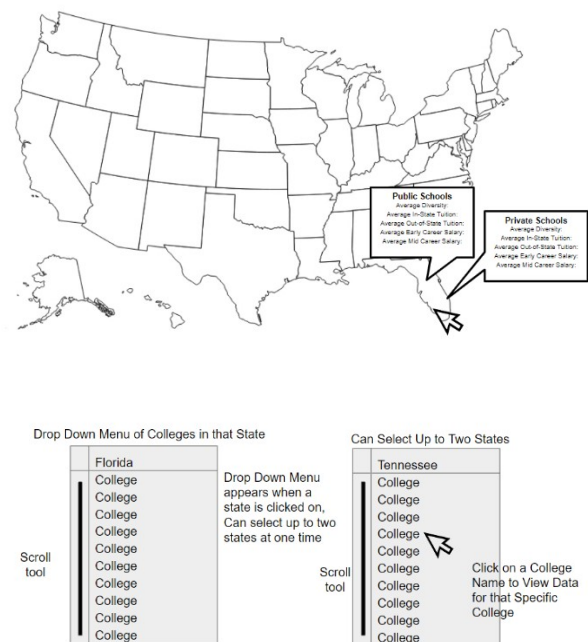


Figure 1: Initial sketch with map and interactive drop down

Upon designing these first two sketches, we assessed that not enough data was available to the user. We also wanted users to be

able to compare two colleges against one another. We added a data visualization for the average diversity with divisions of racial groups in a pie chart, if students are seeking a specific community. They can also view a bar chart comparing mid and early salary expectations against the cost of attendance per year for in and out of state tuition. This allows students to assess specific colleges rather than general information about a state. It also provides the data easily visually available. The final subject will be minimally colored with color representations provided in an immediately available legend. Users can interact with the visualizations by hovering over the bar chart rectangles, sparking a highlight on the same measurement of the other school. For example, highlighting the out of school tuition bar for College A will highlight the out of school tuition bar for College B.

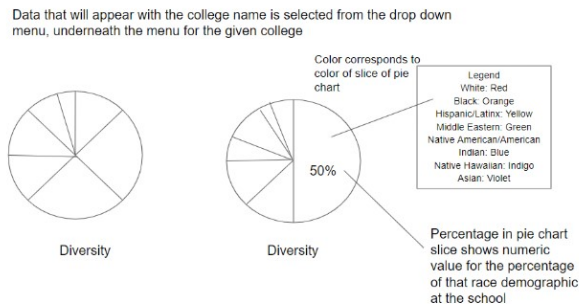


Figure 2: Initial sketch of diversity pie chart

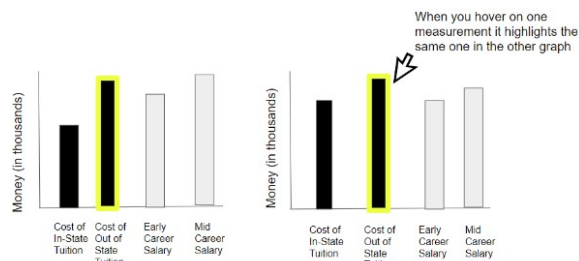


Figure 3: Initial sketch of interactive bar chart

## 6 FINAL DESIGN

Our final design involves marks and channels to communicate the data provided through visual encoding techniques. First, the tool provides an interactive map with state borders and different colleges. Area as a mark shows the area of different states in the United States. Position as a channel, both vertical and horizontal, is important to showing where colleges appear on the map in terms of the state they reside in. It does not communicate their specific position, but the general state only. The mouse hover function can be used to displayed averages of the statistics of each state in the map, and the click function can be used to display a drop down menu. The drop down menu will allow users to select a state a show further data about its specific diversity, tuition, and career prospects. They can select up to two states to show further details, and compare the same statistics on either side. Ideally, this would provide an easy tool for students to weigh different factors in their college decision.

The second portion with further data includes a pie chart and a bar graph. The bar graph uses the mark lines, where their length corresponds to specific data. Their position, in terms of vertical, communicates the size of the wealth for that category. The horizontal

data does not a range of data, but informs the type of measurement being shown. The visualization allows for interactive functionality in highlighting the same measurement in both graphs when one is hovered over, to aid in ease of comparison. Difference in color may also denote whether the data concern tuition or future salary predictions, however does not communicate numeric data. Color is also an important channel for the pie chart, showing the different categories of race displayed. Area as a mark is important to the pie chart slices as well, as it denotes the percentage size of that race at the given school. Tilt is also a channel involved in the pie chart. These various graphical elements allow us to communicate various aspects of the given data in a clear way to the viewer. We specifically attempted to minimize overusing marks and channels that were unnecessary to the data we hoped to represent, and focus on what was essential. This allows the viewer to have clarity on the information and speeds up their comprehension at first viewing. From the rough sketches to the final sketch, we gradually added more functionality to improve user experience. The larger swath of data that is able to be visualized provides the user with more information to influence their decision about specific schools, rather than general information for different states. Our final design thus combines all of these goals.

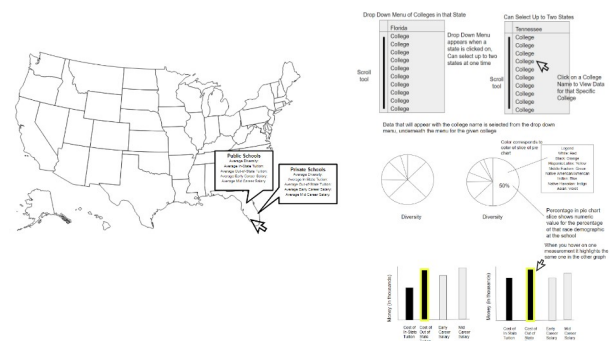


Figure 4: Final sketch of data visualization tool

## 7 DISCUSSION

Will Add Information.

## 8 CONCLUSION

Will Add Information.

## ACKNOWLEDGMENTS

Will Add Information.

## 9 APPENDIX

### 9.1 Data Abstractions

#### 9.1.1 Rows

- Each item in the dataset represents a U.S university

#### 9.1.2 Columns

- Column 1: University Name
  - Attribute Type: Categorical
- Column 2: State Name
  - Attribute Type: Categorical
- Column 3: State Abbreviation
  - Attribute Type: Categorical

- Column 4: Institution Type
  - Attribute Type: Categorical
- Column 5: Degree Length
  - Attribute Type: Categorical
- Columns 6-10: Room and Board Fees, In-State/Out-of-State Tuition, Total In-State/Out-of-State Cost
  - Attribute Type: Quantitative
- Columns 11-12: Early and Mid-Career Pay
  - Attribute Type: Quantitative
- Columns 13-19: Count of each Race Group's Enrollment
  - Attribute Type: Quantitative
- Column 20: Total Enrollment
  - Attribute Type: Quantitative

## 9.2 Task Abstractions

### 9.2.1 Cheapest Out-Of-State Tuition

- **High-Level: Discover**
  - The user is learning new information about out-of-state tuition costs.
- **Mid-Level: Lookup**
  - The user will be able to identify the spatial location of states but will need to lookup further information to understand tuition costs by clicking into the state to see the list of universities in that state.
- **Low-Level: Compare**
  - The user will compare 2 universities to see which one has the cheaper tuition costs.

### 9.2.2 Largest Difference Between Tuition and Salary Potential

- **High-Level: Discover**
  - The user is learning new information about tuition costs relative to salary potential.
- **Mid-Level: Locate**
  - The user doesn't know which universities within which states have the largest differential between tuition and potential salary and must locate those universities themselves.
- **Low-Level: Compare**
  - The user will compare 2 universities to see which one has the largest difference between tuition cost and salary potential.

### 9.2.3 States With Highest Diversity Rates

- **High-Level: Discover**
  - The user is learning new information about diversity rates.
- **Mid-Level: Locate**
  - The user knows they want to find states with high diversity rates, but they don't know where they are located without exploring the map.
- **Low-Level: Compare**
  - The user is looking to identify a single target (state with high diversity).

## REFERENCES

- [1] Four scenarios for 2023: Navigating uncertainty.
- [2] How students decide which college to attend.
- [3] C. et al. A visualization model based on adjacency data. *Decision Support Systems*, 33(4):349–362. doi: 10.1016/S0167-9236(02)00003-9
- [4] J. Mostipak. College tuition, diversity, and pay, Mar 2020.
- [5] A. E. Sobel. The escalating cost of college. *Computer*, 46(12):85–87, 2013. doi: 10.1109/MC.2013.438