

Previous research and conference experience influences students' motivations for pursuing a Ph.D. (Guerin, C. et al., 2014), but successful completion of a Ph.D. relies on perception of support from advisors, faculty, and other students (Litalien, D., & Guay, F., 2015). Conversely, social isolation can motivate Ph.D. candidates to drop out of the program (Young, S. N., et al., 2019), suggesting that prospective Ph.D. candidates' success depends on their fit in an institution's social atmosphere. Not every Ph.D. hopeful lives near an institution that offers a program within their interests, and travel to assess atmosphere can be costly. Without firsthand experience, Ph.D. applicants must somehow answer the question:

This project presents a personalized method to analyze fit between PhD applicants and Ph.D. programs. To analyze programs, faculty, and aspects of the application process, this example includes a set of inclusion criteria personalized to the author's research interests.

CollegeScorecard:
The CollegeScorecard dataset includes data on 6,543 post-secondary institutions with 233,979 different fields of study at various academic levels. This analysis uses the following variables:
INSTNM. The institution's name
INSTURL. The URL to the institution's webpage
CITY. The city in which an institution is located
STABBR. The state in which an institution is located
CREDLEV. The credential level of the program
CIPDESC. The field of study the program offers

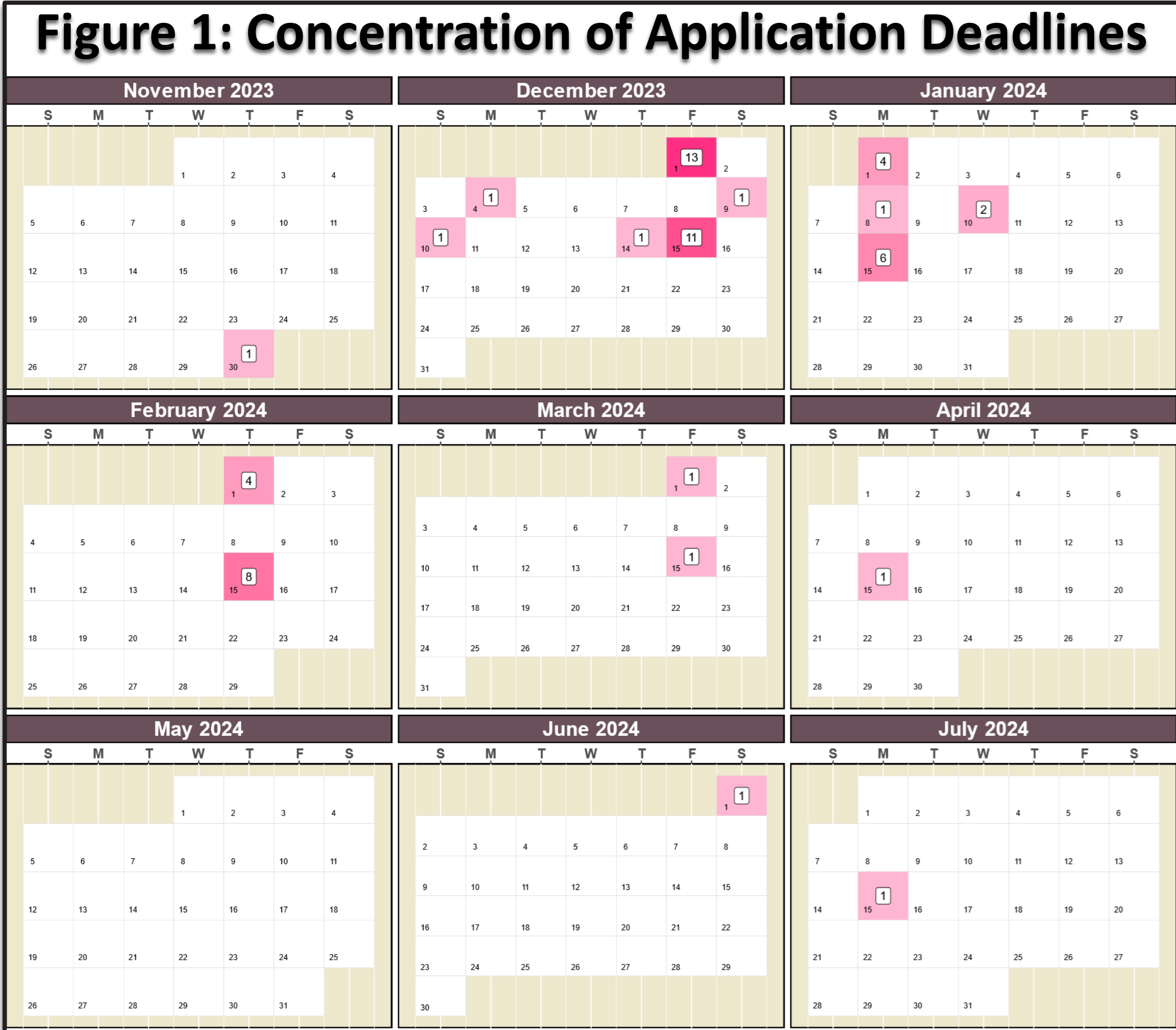
Inclusion Criteria.
Institutions included in the final dataset must:

- Be accredited & offer at least one graduate degree (11,411 programs from 1,828 institutions)
- Offer a doctoral degree program in statistics, data science, computer science, math, or neuroscience (8,012 programs from 1,219 institutions)
- Foster interdisciplinary research opportunities & offer courses in graph theory, probability, machine learning, or neuroscience (220 programs from 60 institutions)

Web Scraping Program & Faculty Data.

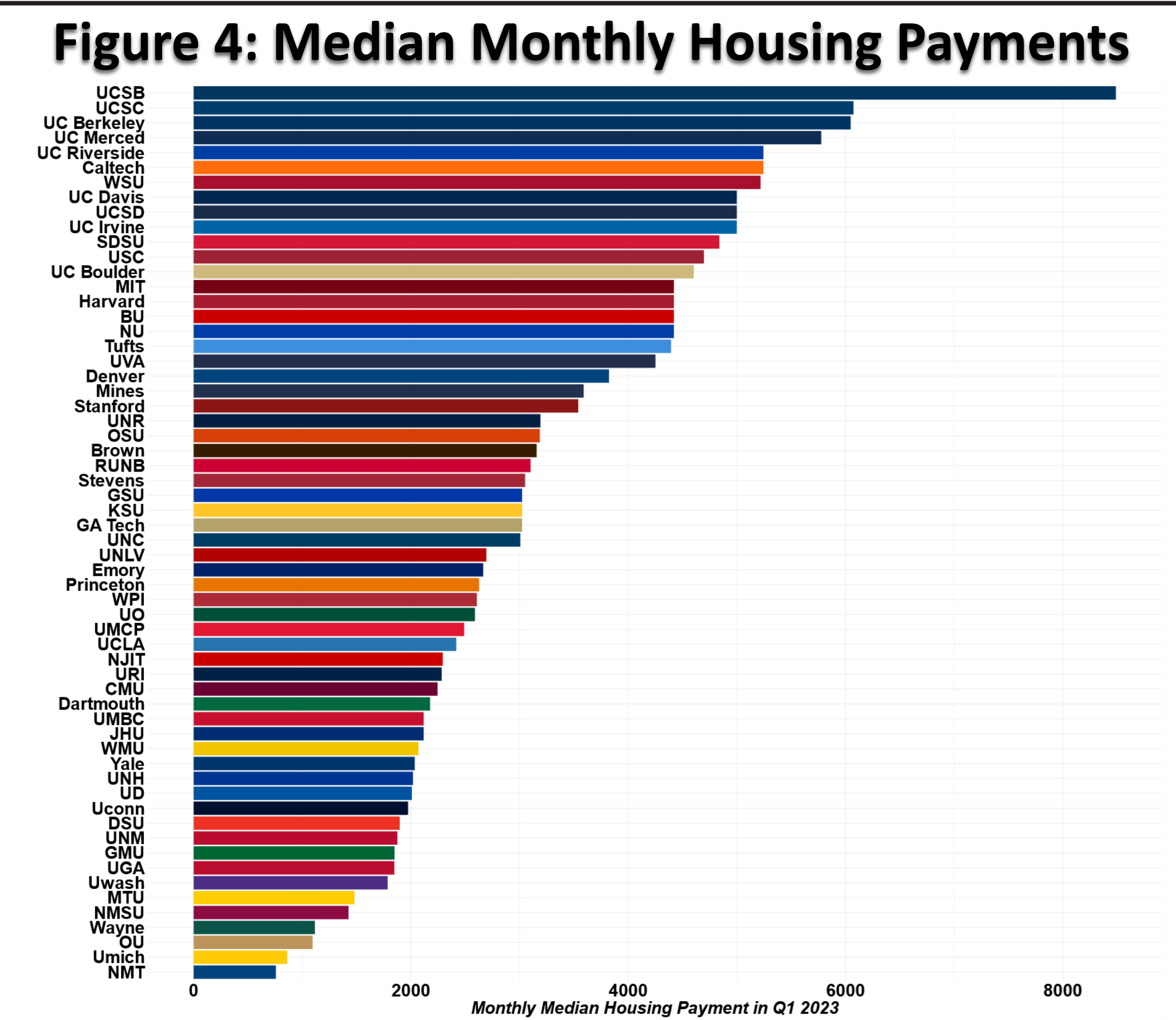
The data included faculty publication data scraped from Google Scholar and institutional webpages and directories using the following methods:

- *Beautifulsoup*. A Python library that calls/extracts webpage HTML. Used to collect program and application information (e.g., degree outlines, statement prompts, faculty member information, application deadlines) using official institution URLs from the CollegeScorecard. This project used this library to scrape program webpages and institutional directories to collect potential faculty members of interest
- *SerpAPI*. A Python API service that automatically scrapes Google webpages (e.g., Google Scholar). Used to collect publication information (i.e., title, abstract, year, co-author information)
- The final dataset included 31,109 publications by 383 faculty members from 8 schools



Nathalie Jones & Alexis Lowery

Figure 2: Frequency of Application Fees by Cost



University	Number of Faculty Members
Stanford	95
UC Berkeley	60
Cornell	60
UC San Diego	52
Brown	40
MIT	30
CalTech	28
KSU	20

A bar chart with 'Frequency' on the y-axis (0 to 40) and two categories on the x-axis: 'Fee Waiver Available' and 'GRE Required'. For each category, there are two bars: a dark pink bar for 'No' and a light pink bar for 'Yes'. The 'No' bar for 'Fee Waiver Available' is approximately 18, and the 'Yes' bar is approximately 42. The 'No' bar for 'GRE Required' is approximately 45, and the 'Yes' bar is approximately 14.

Category	No	Yes
Fee Waiver Available	18	42
GRE Required	45	14

University	Number of Programs Offered
UC Irvine	7
UC Davis	6
GMU	6
RUNB	6
Umich	6
UMCP	6
GA Tech	6
UCLA	6
UC Berkeley	6
Yale	5
Harvard	5
BU	5
UCSC	5
Brown	5
JHU	5
Princeton	4
UNH	4
GSU	4
UNLV	4
UD	4
UC Boulder	4
UVA	4
UGA	4
UCSD	4
MIT	4
USC	4
UC Riverside	4
Catech	4
Stanford	4
WPI	4
URI	3
Wayne	3
NMT	3
UNM	3
MTU	3
UNR	3
Stevens	3
WSU	3
Uconn	3
NJIT	3
UCSB	3
NU	3
CMU	2
NMSU	2
DSU	2
WVU	2
OU	2
KSU	2
UO	2
Mines	2
Emory	2
Tufts	2
Dartmouth	2
Denver	2
GSU	2
UNC	1
UMBC	1
UC Merced	1
SDSU	1

Horizontal stacked bar chart showing the percentage of publications by topic for eight universities. The x-axis represents the percentage of publications (0% to 100%). The y-axis lists the universities. The legend identifies five topics: Harmonics & Hierarchical Modeling (pink), ML/AI & Language Modeling (light green), Probability, Graphs, & Networks (yellow), Marketing, Economics, & Business (light blue), and Social Statistics (grey).

University	Harmonics & Hierarchical Modeling	ML/AI & Language Modeling	Probability, Graphs, & Networks	Marketing, Economics, & Business	Social Statistics
UC San Diego	35%	15%	40%	5%	5%
UC Berkeley	40%	20%	35%	0%	5%
Stanford	45%	30%	15%	5%	5%
MIT	40%	20%	25%	2%	13%
KSU	35%	10%	15%	25%	15%
Cornell	40%	15%	25%	2%	18%
CalTech	45%	10%	35%	0%	10%
Brown	30%	20%	25%	5%	20%

