



Data Science Capstone

First-Gen @ Wellesley:

Comparisons of Post-Bac Outcomes Between First-Gen and non First-Gen Alums at Wellesley

Briana Vigil '22

Objectives & Research Question:

First-generation students make a small yet integral part of the Wellesley community, and their experiences in higher education are varied, but research shows that first-generation students “often face psychological, academic, financial, and social challenges, and about one in three leave college within the first three years” (mghclaycenter.org) This brings me to my question of: Having persevered against the obstacles at Wellesley and once graduated, do first-gen Wellesley alums have different outcomes than other peers in terms of pursuing higher degrees?

Data Collection & Cleaning

In order to analyze this question, I used data scraped from the WellesleyHive, an online networking platform created exclusively for Wellesley alums, current students, and faculty. I filtered the data for 'First-generation' and 'Alum' to ensure that everyone I have in my first-generation data are as accurate as possible. For the non first-gen alums, I simply filtered for “Alum” and scraped the data from those profiles. I cleaned and organized the data using Python's BeautifulSoup and Pandas libraries. Once cleaned, the data were converted into csv files and imported into R to begin with the analysis.

Variable Selection & Methodologies

To facilitate my analysis, I created several “dummy variables” that indicated whether an alum has obtained a Master's or PhD, or if they are a current Graduate Students, as well as whether the alum's major was in the STEM, Social Science, or Humanities field. I also extracted the year they graduated and approximated their age that would be accurate given the assumption that they graduated within 4 years and were a traditional student rather than a Davis Scholar. Although it may be inaccurate to assume so, this method facilitated my approach to come to my analysis. With these variables, I was able to begin my analysis for the pursuit of higher education for first-generation and non first-generation Wellesley alums.

Results:

Upon seeing the results of the boxplot, I made the decision to subset my original data to only look at the people who are 65 and younger, due to decreasing of data as well as overfitting that was skewing the results. The boxplot in Fig. 1 aided in my decision since it clearly depicts the sparseness of the data once the alum age is over 65. In Fig. 2, we can see that the scatterplot matrix depicts the strongest correlation between ApproxAge and Masters, which

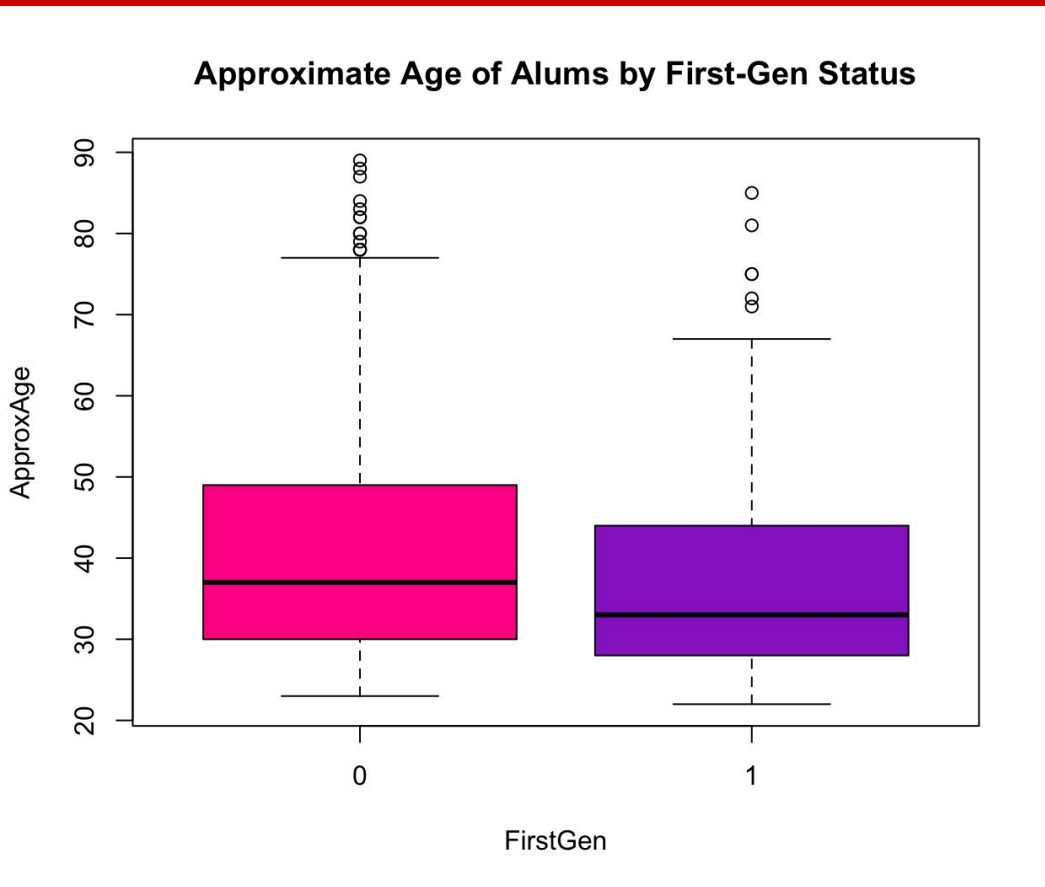


Fig. 1: Boxplot of Approximate Age and First Gen status of Wellesley Alums

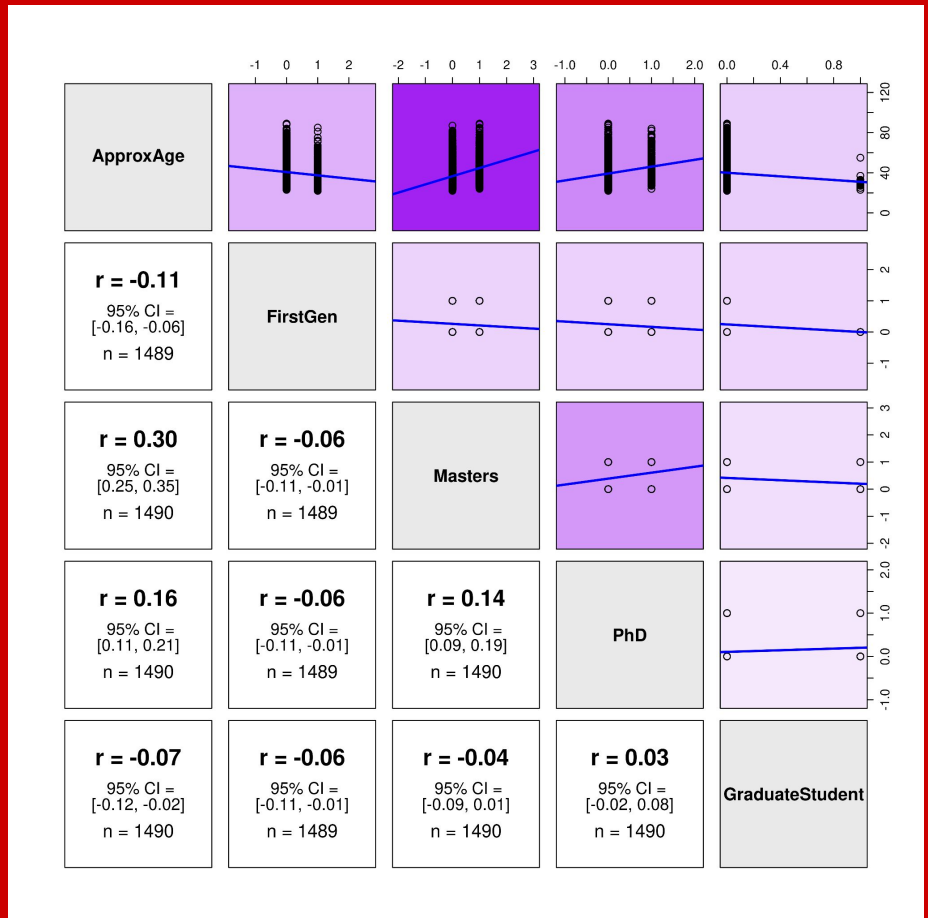


Fig. 2: Scatterplot Matrix depicting all correlations of all variables in dataset. Darker colors signify stronger relationships between variables

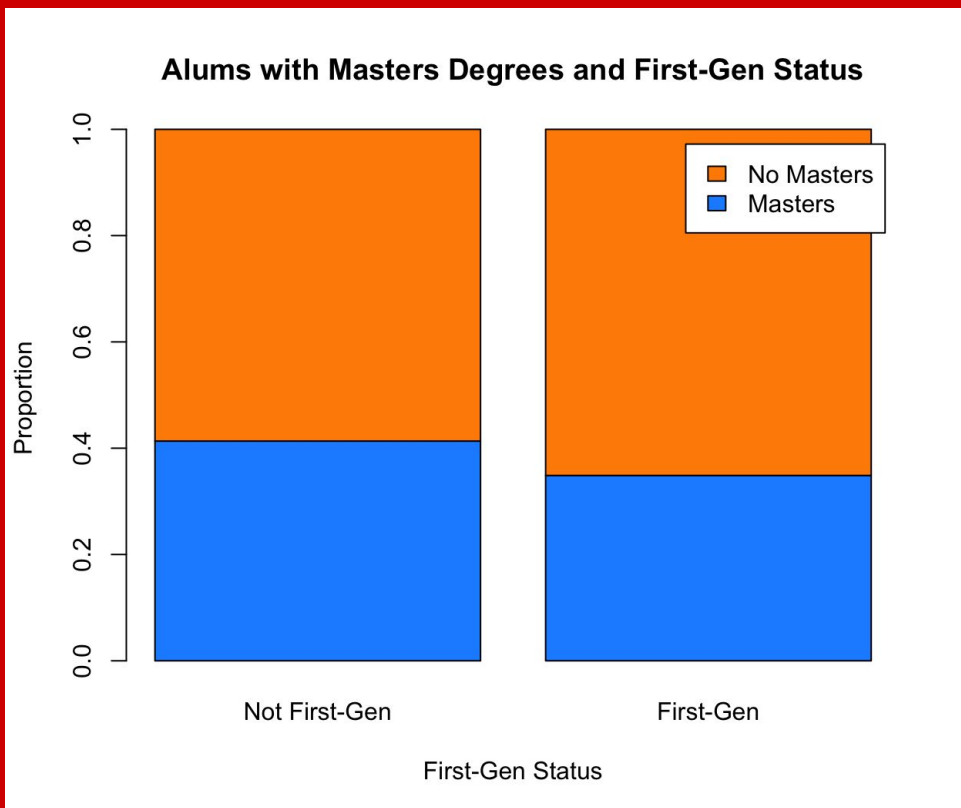


Fig. 3. Barchart that depicts the proportion of Wellesley alums (both First Generation and non) and the amount of Masters degrees that they received

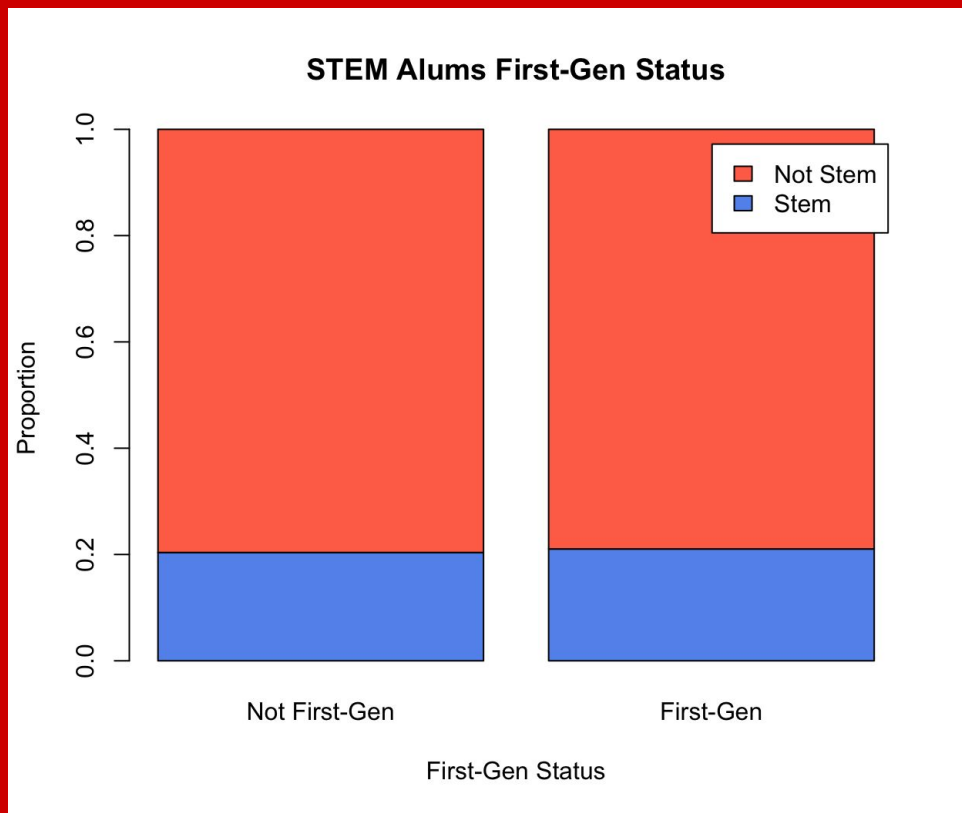


Fig. 5. Barchart that depicts the proportion of Wellesley alums (both First Generation and non) and the amount of alums who are primarily STEM oriented in their academic history

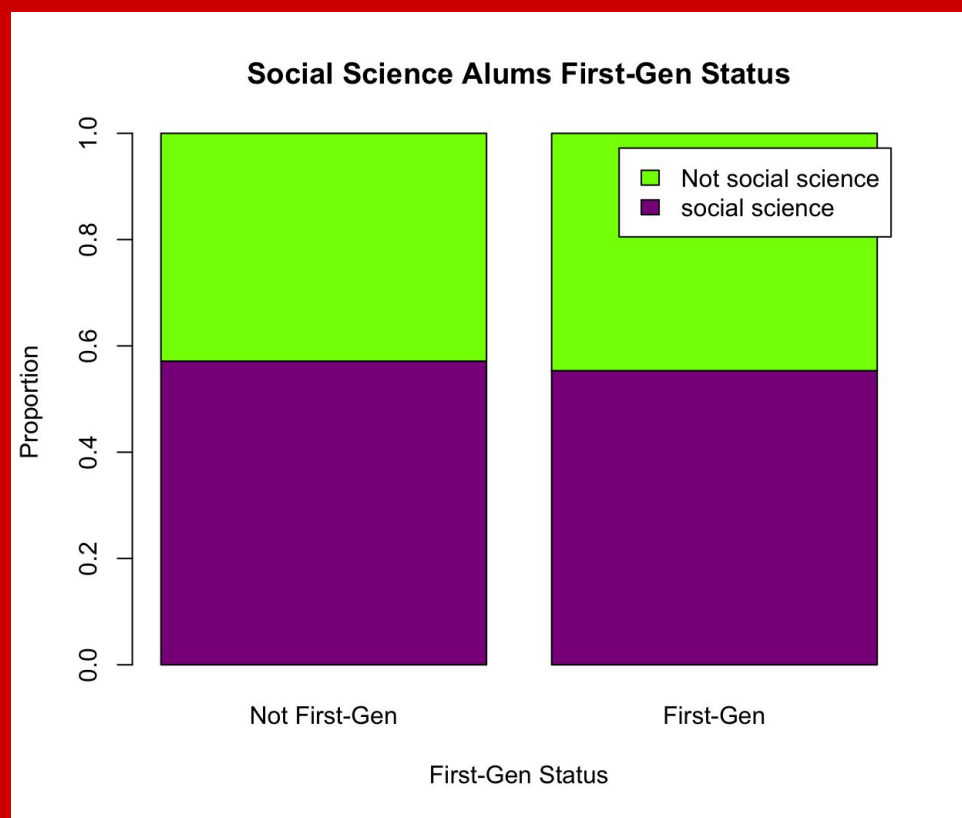


Fig. 7. Barchart that depicts the proportion of Wellesley alums (both first gen and non) and the amount of alums who are primarily Social Science oriented in their academic history

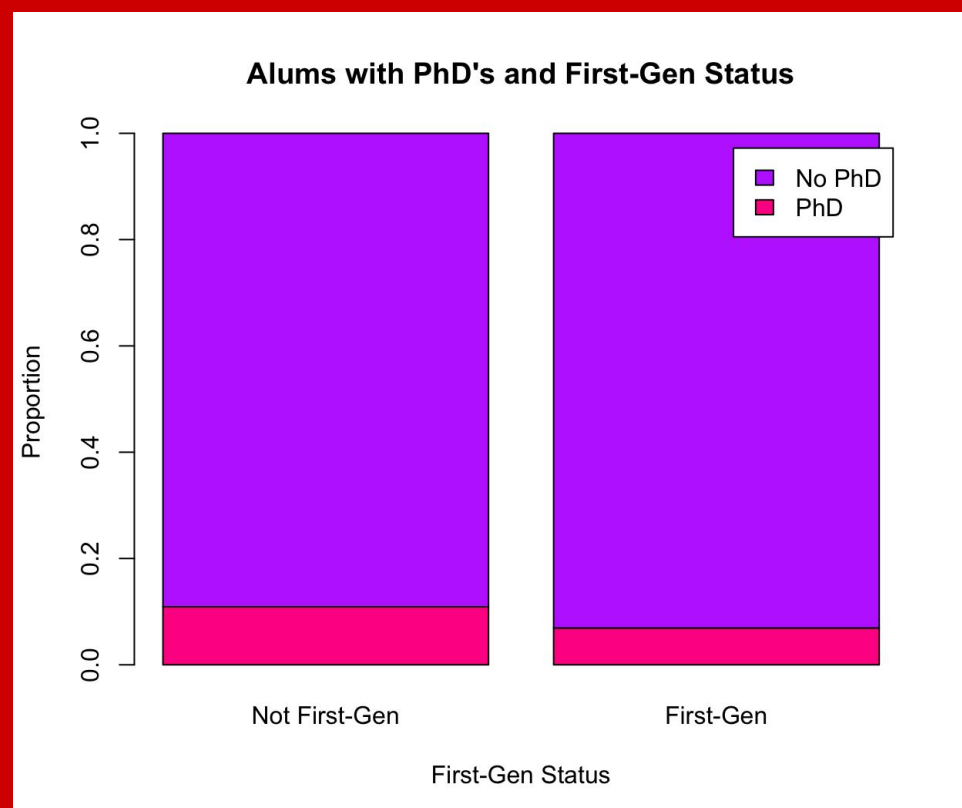


Fig. 4. Barchart that depicts the proportion of Wellesley alums (both First Generation and non) and the amount of PhD's that they received

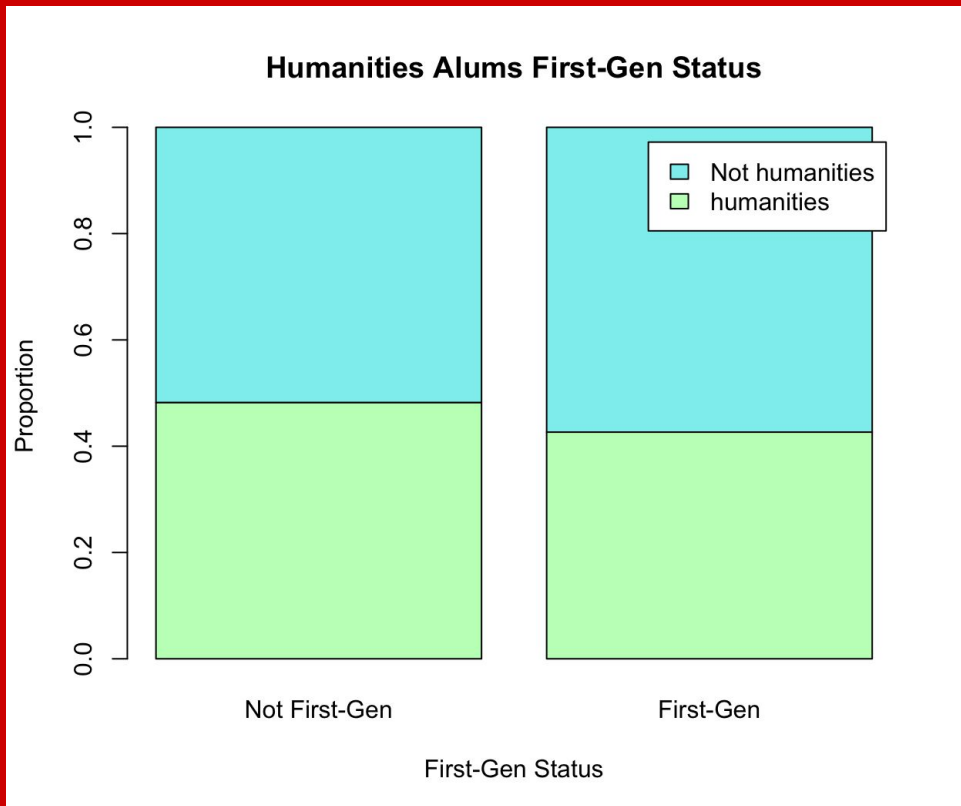


Fig. 6. Barchart that depicts the proportion of Wellesley alums (both First Generation and non) and the amount of alums who are primarily Humanities oriented in their academic history

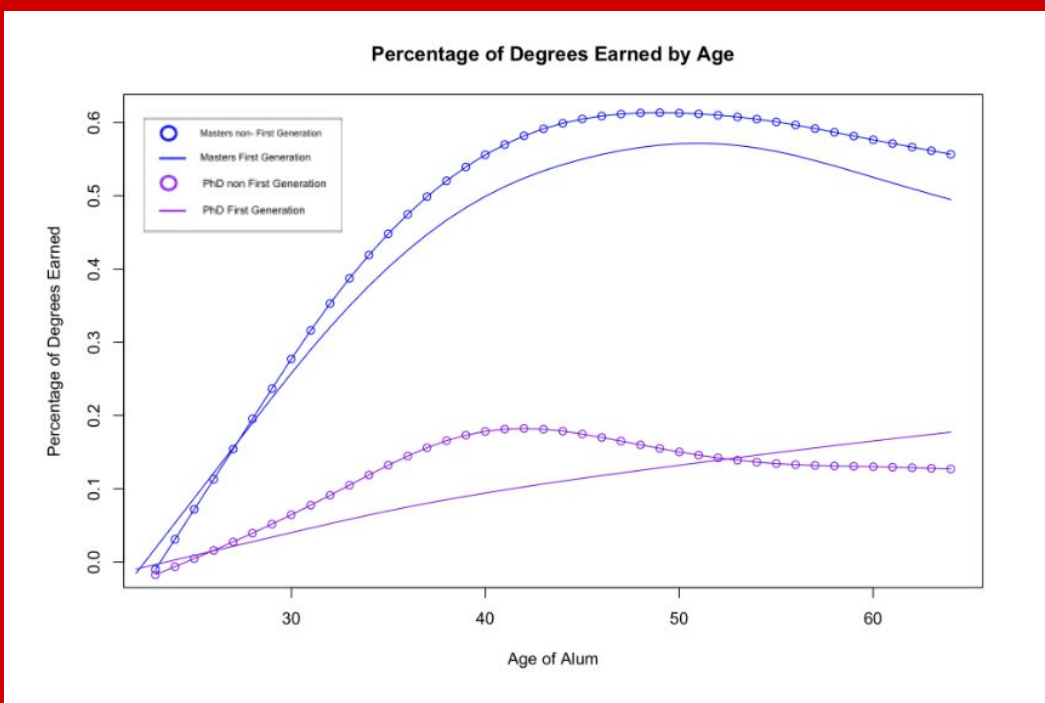


Fig. 10. Spline that depicts the percentage of Masters and PhD degrees received by both First Generation and non through the different ages of alums.

	Masters Degree: all predictors	Masters Degree: no disciplines	Masters Degree: interactions for First Gen and age variable	Masters Degree: interactions between age and STEM	Masters Degree: interactions between age and Social Science	Masters Degree: interactions between age and Humanities
Intercept	-3.33*** (0.268)	-3.79*** (0.225)	-3.20*** (0.281)	-3.23*** (0.281)	-2.94*** (0.305)	-3.66*** (0.342)
FirstGen	-0.387 (0.191)	-0.383 (0.194)	-0.340 (0.148)	-0.345 (0.148)	-0.276 (0.196)	-0.279 (0.196)
STEM	0.131 (0.160)	0.132 (0.160)	-0.240 (0.071)	0.221 (0.139)	0.133 (0.160)	0.133 (0.160)
SocialScience	0.007*** (0.005)	0.007*** (0.005)	0.007*** (0.005)	0.006*** (0.005)	0.006*** (0.005)	0.006*** (0.005)
Humanities	0.398*** (0.120)	0.398*** (0.120)	0.398*** (0.120)	0.398*** (0.120)	0.398*** (0.120)	0.398*** (0.120)
ApproxAge	0.002*** (0.000)	0.002*** (0.000)	0.002*** (0.000)	0.002*** (0.000)	0.002*** (0.000)	0.002*** (0.000)
FirstGen * ApproxAge			-0.003 (0.003)			
STEM * ApproxAge				0.003 (0.003)		
SocialScience * ApproxAge					0.003 (0.003)	
Humanities * ApproxAge						-0.003 (0.003)
Num.Obs.	1392	1392	1392	1392	1392	1392
AIC	1708.4	1708.4	1708.3	1708.9	1707.9	1707.7
BIC	1718.8	1718.3	1715.1	1716.5	1714.3	1714.1
Log Lik.	-886.151	-886.901	-886.180	-887.054	-886.941	-886.820
R-sq	0.11	0.12	0.11	0.11	0.11	0.11
p < 0.001, ** p < 0.01, * p < 0.05, *** p < 0.001						

Fig. 8. Model Summary of all logistic regression models conducted for Masters Degrees. This shows the different kinds of models that were run to see which one was optimal for the purposes of this research question

	PhD Degree: all predictors	PhD Degree: no disciplines	PhD Degree: interactions for First Gen and age variable	PhD Degree: interactions between age and STEM	PhD Degree: interactions between age and Social Science	PhD Degree: interactions between age and Humanities
Intercept	-22.054 (611.131)	-3.681*** (0.340)	-22.763 (610.960)	-22.981 (604.620)	-23.330 (610.602)	-21.794 (614.240)
FirstGen	-0.554* (0.201)	-0.433* (0.230)	-0.497* (0.206)	-0.553* (0.276)	-0.551* (0.261)	-0.554* (0.287)
STEM	0.302*** (0.077)	0.348*** (0.070)	0.260 (0.147)	0.260 (0.147)	0.302*** (0.077)	0.302*** (0.077)
SocialScience	0.006 (0.010)	0.006 (0.010)	0.006 (0.010)	0.006 (0.010)	0.006 (0.010)	0.006 (0.010)
Humanities	0.398 (0.121)	0.398 (0.121)	0.398 (0.121)	0.398 (0.121)	0.398 (0.121)	0.398 (0.121)
ApproxAge	0.004*** (0.001)	0.004*** (0.001)	0.004*** (0.001)	0.004*** (0.001)	0.004*** (0.001)	0.004*** (0.001)
FirstGen * ApproxAge			0.045* (0.010)			
STEM * ApproxAge				0.045 (0.010)		
SocialScience * ApproxAge					-0.003 (0.010)	
Humanities * ApproxAge						0.003 (0.010)
Num.Obs.	1392	1392	1392	1392	1392	1392
AIC	618.8	619.6	617.5	618.5	620.1	620.8
BIC	630.2	630.3	624.2	630.2	637.4	637.4
Log Lik.	-303.278	-303.313	-303.750	-303.275	-303.302	-303.376
R-sq	0.06	0.06	0.06	0.06	0.06	0.06
p < 0.001, ** p < 0.01, * p < 0.05, *** p < 0.001						

Fig. 9. Model Summary of all logistic regression models conducted for PhDs. This shows the different kinds of models that were run to see which one was optimal for the purposes of this research question

aided in my final decision of subsetting of the data. In addition to this, I also wanted to visualize the proportion of First Generation alums to the different predictors that are in my models. Figures 3-7 depict this, where we can see the proportion of alums with: Masters degrees, PhDs, in STEM, Humanities, and Social Science fields. In all of these visualizations, we see that First Generation alums are slightly lower in proportion to non First Generation counterparts.

To view the variables that were most important, I also conducted a series of logistic regressions for both Masters and PhD with all other predictors that were relevant to the research question. Figure 8 depicts the series of logistic regression that were run for the Masters outcome, where I ultimately concluded that the model with all of the predictors plus the interaction between Humanities and Age was the best one because of the lower AIC of 1707.65. Figure 9 depicts the series of logistic regression that were run for the PhD outcome, where I ultimately concluded that the model with all of the predictors plus the interaction between First Generation and Age was the best one because of the lower AIC of 617.5. In addition to this, what we can also see from this summary is that the most significant predictor is age, but there are also negative correlations with the outcome and being First Generation.

Discussion:

Although the results were not exactly what I had hoped (ideally there would have been no differences), these findings reveal that the gap between First Generation and non First Generation alums is relatively small and although not as significant as the age variable, an alum identifying as First Generation does negatively correlate with having a Masters and/or a PhD. It is also important to note that there may be missing data in terms of current graduate students since alums may have not updated their education history yet.

To close the gap, there is a need for further support and programming for First-Generation upperclassmen at Wellesley. There are already amazing structures in place such as WellesleyPlus and the FLI Network, and these are especially helpful for the transition to Wellesley. Wellesley also needs structures for the transition out of Wellesley, so that first-generation sibs are fully equipped with the resources and confidence needed for further pursue of education.

Further Research:

Prof. Jeremy Wilmer provided direction to the Standard Occupation Classification (SOC) by the CDC that I hope to include in further iterations of this project in order to observe any differences or lack thereof between occupations pursued by First-Gen or non First-Gen alums.

Special thanks to Rebecca Garcia, former dean at Wellesley College, for her unconditional support of myself and all other First-Generation students that had the honor to work with her.