

GoodGrant Foundation: Next Generation Investing

A Model for Post-Secondary Educational Grants

February 1, 2016

Summary

These are the strategies of the consulting firm Team #54301. It's five-year mission: to explore new ranking systems; to seek out new return on investments and new student performance metrics; to boldly go where no charity has gone before.

Contents

1	Introduction	3
1.1	The Problem	3
1.2	Assumptions and Rationale/Justification	3
1.3	Summary of our Approach	3
2	Model Design and Approach	5
2.1	Notations and Definitions	5
2.2	Benchmark Selection	5
2.3	Categorizing Institutions into Classes	6
2.4	Investment Strategy	6
2.4.1	Number of students-Helped	6
2.4.2	Adjusted Poverty	6
2.4.3	Return on Investment (ROI)	7
2.4.4	Maximization of return	7
2.4.5	Simplifying the complexity	7
2.5	Limiting Schools	8
3	Model Application	8
3.1	Partitioning and Benchmark	8
3.2	Variable Selection	9
3.2.1	Filtering by Grant	9
3.2.2	Scorecard Variables	9
3.2.3	IPEDS Variables	10
3.3	kNN Selection	10
3.4	Maximizing the Return	11
4	Discussion	12
4.1	Sensitivity Analysis	12
4.2	Strengths	13
4.3	Weaknesses	13
5	Extensions	13
6	Conclusions	13
7	Letter to Mr. Alpha Chiang	14

1 Introduction

1.1 The Problem

The Goodgrant Foundation is a charitable organization with a mission to improve educational performance of students attending colleges and universities in the United States. The foundation intends to donate a total of 100,000,000 USD to a group of schools each year, for five years. They also wish to avoid investing in universities that have already received grants from major charitable organizations such as the Gates and Lumina Foundation.

We present a mathematical model to outline an ranking and investment strategy to provide the maximal return on investment(ROI) appropriate for an education charity. The model is split into two sections and is applicable in choosing $1 \cdots N$ institutions.

The first section outlines determining bench mark points and comparing the rest of the institutions to them via a k-Nearest Neighbours(kNN) Algorithm. In the second section, we derive a ROI measuring the effect the grant would have on improving student performance. In particular, our ROI reflects changes in terms of incoming poverty and outgoing poverty. In the Model Testing section, we input a subset of the U.S. National Center on Education Statistics data set and College Scorecard data set and output a list of $1, \cdots N$ schools, the investment amount, and number of students affected. In Sensitivity Analysis, we restrict the distance in the kNN algorithm to obtain a subset of our prioritized schools and assess the model's stability. Lastly, we summarize the findings of the paper and suggest future additions.

1.2 Assumptions and Rationale/Justification

- **The Gates Foundation and Lumina Foundation will continue to donate to the same school for at least two years:** Typically large charities often donate sums over the course of number years to allow highest impact. [4] We are able to remove those particular schools from our candidate list to avoid duplicating the two foundation efforts.
- **Students that attend GoodGrant Grant Recipients will continue to graduation.:** Lack of financial support is the leading cause of students dropping out of institutions. Students who receive grants have less financial stress and can focus on graduating.[1]
- **Attainment of higher education improves poverty rates** According to the US Census Bureau, individuals with an advance degree earn on average \$72,824 per year which is significantly higher then individuals with a high school degree at \$45,400.[2] Generally, attending post-secondary education lifts individuals above the poverty-line.

1.3 Summary of our Approach

- Create a benchmark setting: top 10 schools for each degree type school, that we consider to have good performance. Each school represent a class.
- Identify universities that have similar characteristics to the benchmark list.
- Prioritize universities to be invested in by the amount of grants they received from other sources.

- Determine the amount of investment and allocation to maximize the return on investment.

2 Model Design and Approach

2.1 Notations and Definitions

1. $U_j^{[i]}$ = university j of class i
2. $U_o^{[i]}$ = benchmark university for class i
3. $A_j^{[i]}$ = amount invested in university j for class i
4. $T_j^{[i]}$ = Tuition in university j for class i
5. $n_j^{[i]}$ = number of students receiving our grant in university j for class i
6. $IP_j^{[i]}$ = poverty rate of students entering university j for class i
7. $OP_j^{[i]}$ = poverty rate of students exiting university j for class i
8. $\phi_j^{[i]}$ = number of students in university j for class i
9. Φ = Total number of universities
10. $g_j^{[i]}$ = maximum amount of grants to be given to university j for class i
11. Group, refers to the Carnegie Classification of the university
12. Class, refers to which benchmark within the group, the university is closest to.

2.2 Benchmark Selection

First we omitted any below medium size universities or communal universities such as racial or religious schools. Then we partition all universities by number of years required for graduation. Lastly, we rank the universities by the amount of their poverty elevation standards (PES), this is done based on the result and recommendation of Trombitas [1] and College-Score Card [11].

$$R_j^{[i]} = \frac{IP_j^{[i]} - OP_j^{[i]}}{IP_j^{[i]}} \quad (1)$$

Next we choose the top, highest ranking universities per ‘years of graduation’-group. We set these chosen universities as our benchmarks for each class in each group.
i.e.

$$o^{[i]} = \underset{j}{\operatorname{argmax}} R_j^{[i]}$$

Where o , is to represent the benchmark index of the class i .

2.3 Categorizing Institutions into Classes

We aim to classify each university into some class i . This is done via the kNN algorithm: a non-parametric method used in classification supervised learning problems. The output is determined by a majority vote, where the object is assigned to the class most common with its neighbors. [7] We set up the following model: Assuming $j \in [1, n]$, where n = amount of variables. For a detailed selection list please see Appendix 1. Let $y_j^{[i]}$ be the j -th variable for the i -th class benchmark. Let x_j be the j -th variable of a particular university. A university u , is considered a member of the k -th class if the variables x_i of that university is ‘close’ to a particular variable of a given class in euclidean-space.

That is, if u belongs to the k -th class then

$$k = \underset{i}{\operatorname{argmin}} \sqrt{\sum_{j=1}^n (x_j - y_j^{[i]})^2}$$

2.4 Investment Strategy

2.4.1 Number of students-Helped

Let $A_j^{[i]}$ be the amount we invest in university j of the class i , and $T_j^{[i]}$ be the tuition of low-income students. Then we can compute the number of students who we help as:

$$n_j^{[i]} = \lfloor \frac{A_j^{[i]}}{T_j^{[i]}} \rfloor \quad (2)$$

2.4.2 Adjusted Poverty

Next we computed the adjusted poverty rate. We make a crucial assumption here about the correlation between grant-money and retention rates. We assume that every student who receive the grant will stay in the program and university (this assumption can be adjusted by a future study on the probabilities between grant and retention rates).

Let $\phi_j^{[i]}$ be the total number of students in university j for class i .

The adjusted rate of existing poor students based on our assumptions is:

$$OP_j^{[i]*} = \frac{OP_j^{[i]} \phi_j^{[i]} - n_j^{[i]}}{\phi_j^{[i]}} \quad (3)$$

We can then provide a rank to this university, in the similar fashion we computed the rank for the benchmark universities (recall eq. (1)):

$$R_j^{[i]*} = \frac{IP_j^{[i]} - OP_j^{[i]*}}{IP_j^{[i]}}$$

We use this to compute the increase to the social well-being of students going to this university.

2.4.3 Return on Investment (ROI)

The return on investment for our model will be the change in the level of alleviation of poverty provided by our investment.

We do this as follows:

$$R_j^{[i]}(1 + r_j^{[i]}) = R_j^{[i]*} \quad (4)$$

where $r_j^{[i]}$ is the return on investment for a particular school.

We also note that (4) can be rewritten using (1) and (3):

$$r_j^{[i]} = \frac{OP_j^{[i]} - OP_j^{[i]*}}{IP_j^{[i]} - OP_j^{[i]}} \quad (5)$$

This will later be used for simplification purpose of the maximization problem.

We then average these rates, to get the approximated return:

$$\bar{r} = \frac{\sum_{j=1}^{\infty} \sum_{i=1}^{\infty} r_j^{[i]}}{\Phi} \quad (6)$$

Where Φ is the total number universities.

2.4.4 Maximization of return

Now that we have a measure of return on investment, we wish to maximize this measure based on the physical restriction that are given to us.

Since we only have a limit of \$100M, and would not want to overshoot our investment by giving any particular university more than what the top-tier schools are given. We put the following restriction forth.

Where $g_j^{[i]} = \max(\frac{G_o^{[i]}}{\phi_o^{[i]}} - \frac{G_j^{[i]}}{\phi_j^{[i]}}, 0)$, is the difference between per student investment with the benchmark school. We formulate the problem as follows:

$$\begin{aligned} & \underset{A}{\text{maximize}} \quad \bar{r} \\ & \text{subject to} \quad \sum \sum A_j^{[i]} = 100,000,000 \\ & \quad \forall i, j \quad 0 \leq A_j^{[i]} \leq g_j^{[i]} \phi_j^{[i]} \end{aligned} \quad (7)$$

2.4.5 Simplifying the complexity

In order to simplify the complexity of the previous problem we outlined before. We reformulated the problem as follows:

First, we rewrote (5) using (2) as,

$$\begin{aligned} r_j^{[i]} &= \frac{1}{\phi_j^{[i]}(IP_j^{[i]} - OP_j^{[i]})} n_j^{[i]} \\ r_j^{[i]} &= \Delta_j^{[i]} n_j^{[i]} \end{aligned} \quad (8)$$

We note that $\Delta_j^{[i]}$ from (4) is the information given by our data. It must be positive as

$$\Delta_j^{[i]} \implies IP_j^{[i]} < OP_j^{[i]}$$

Which means that the considered university is getting more in funding than the benchmark one. If this happens, we omit that university from our dataset.

Whereas $n_j^{[i]}$ is the number of students we invest in, which is a variable that will be controlled by us.

Hence the problem as outlined in (6) return can be written as:

$$\bar{r} = \sum_{j=1}^{\infty} \sum_{i=1}^{\infty} \frac{1}{\Phi} r_j^{[i]} = \frac{1}{\Phi} \sum_{j=1}^{\infty} \sum_{i=1}^{\infty} \Delta_j^{[i]} n_j^{[i]} \quad (9)$$

So our new maximization problem can be reformulated as follows:

$$\begin{aligned} & \underset{\forall n_j^{[i]}}{\text{maximize}} && \frac{1}{\Phi} \sum_{j=1}^{\infty} \sum_{i=1}^{\infty} \Delta_j^{[i]} n_j^{[i]} \\ & \text{subject to} && \sum \sum n_j^{[i]} T_j^{[i]} = 100,000,000 \\ & && \forall i, j \quad 0 \leq n_j^{[i]} \leq \frac{g_j^{[i]} \phi_j^{[i]}}{T_j^{[i]}} \end{aligned} \quad (10)$$

2.5 Limiting Schools

The model outlined in the previous section, is given the ability to reject schools. This can be seen in the condition:

$$\forall i, j \quad 0 \leq n_j^{[i]} \leq \frac{g_j^{[i]} \phi_j^{[i]}}{T_j^{[i]}}$$

However, we wanted to eliminate this ability, in order to provide the capacity to maximize over any number of schools. So we rewrote the condition as

$$\forall i, j \quad 1 \leq n_j^{[i]} T_j^{[i]} \leq g_j^{[i]} \phi_j^{[i]} \quad (11)$$

As a side note, this would mean that $g_j^{[i]} \phi_j^{[i]} \geq 1$. Meaning, that some schools may be omitted before the optimization process starts.

3 Model Application

3.1 Partitioning and Benchmark

As discussed in the previous section, the first goal in the analysis is to find a list of benchmark schools. However, before doing so, we partitioned the schools by the type of degrees that they predominantly provide. This was due to the fact that otherwise, many pieces of data become unusable, and the kNN algorithm becomes biased towards schools that have less-years-to-graduation.

	INSTNM	PES
1	Washington County Community College	0.15391917
2	Carteret Community College	0.13717617
3	Independence Community College	0.12846153
4	Western Iowa Tech Community College	0.12329040
5	Sauk Valley Community College	0.11671234
6	Cloud County Community College	0.10996537
7	Central Community College	0.10869516
8	Elgin Community College	0.08656496
9	Otero Junior College	0.08413926
10	Southwestern Illinois College	0.07424521

(a) 2 Year Program

	INSTNM	PES
1	Vermilion Community College	0.2574341
2	Eastern Wyoming College	0.2301173
3	Dawson Community College	0.1869565
4	Lewis and Clark Community College	0.1839621
5	University of Wisconsin Colleges	0.1838409
6	Dakota College at Bottineau	0.1815356
7	Highland Community College	0.1643508
8	Southwestern Michigan College	0.1630952
9	Montcalm Community College	0.1614305
10	Southwestern Oregon Community College	0.1603809

(b) Bachelor Program

	INSTNM	PES
1	University of North Carolina at Asheville	0.2689571
2	Ohio University-Main Campus	0.2332701
3	Valley City State University	0.2301947
4	The University of Tennessee-Chattanooga	0.2211343
5	SUNY at Purchase College	0.2190117
6	Lake Superior State University	0.2181373
7	Central Michigan University	0.1996328
8	University of Science and Arts of Oklahoma	0.1979167
9	Georgia College and State University	0.1976155
10	University of Minnesota-Morris	0.1938063

(c) 6 Year Program

Once we acquired this data, we chose the top 10 schools from each school type as our benchmark schools. It is also important to note the schools here do not only provide the designated degree, but rather have a predominate focus on such degrees.

We used the schools noted in Figures (a),(b) and (c) as our benchmark schools for the kNN method. These school already perform well based on our ranking, and we wanted to find the next best schools that would benefit from an investment by our charity.

3.2 Variable Selection

3.2.1 Filtering by Grant

Since we did not want to overlap with grants given by other institutions, we filtered the next list using the data given by the Gates[5] and Lumina [6] Foundations.

3.2.2 Scorecard Variables

First we created a dictionary for all the variables that were given to us by Scorecard. We had to categorize these variables as follows:

- **Irrelevant** [0], values such as location of website, or miscellaneous ids.
- **Informative** [1], flags for women-only, historically black colleges and others.
- **Selected** [2], the values we chose to focus on, such as SAT, ACT, retention rates etc.
- **Percentile** [3] values that refer to the composition of the student body (such as racial, sexual or income composition).

We did make use of [1], and [3], for cleaning and producing other variables (such as the PES rank), however we did not run the kNN algorithm on them.

3.2.3 IPEDS Variables

Since many of the variables provided by scorecard, overlapped with the ones from IPEDS, we only withdrew the following variables:

- **Number of students**, which we used for computation of kNN and for the maximization algorithm.
- **Average Tuition of Low Income Students**, which we used mainly for the maximization.
- **Total Grants Received**, we summed the various grants that each university received to get this number. We then used it for both kNN and maximization.

3.3 kNN Selection

After this point, we used the kNN algorithm as the method of selection for the nearest best schools. As we wish to avoid investing in the top tier schools, and would rather invest in schools with the potential to become the top-tier.

Prior to the application of the kNN algorithm we normalized the data by transforming each variable to a scale in $[0, 1]$, in order to avoid bias for larger scaled-values.

This identification provided us with the following schools as can be seen in fig.1.

	Name	Class	Distance
1	[University of South Carolina-Union]	0	0.0115237864
2	[Ohio State University Agricultural Technical Institute]	0	0.0523607966
3	[Rainy River Community College]	0	0.1576588469
4	[Warren County Community College]	1	0.0028690275
5	[Mitchell Technical Institute]	1	0.0066463260
6	[Williston State College]	1	0.0077858030
7	[Coconino Community College]	1	0.0083305728
8	[Lamar Institute of Technology]	1	0.0084850685
9	[Texas State Technical College-West Texas]	1	0.0114658642
10	[Pennsylvania Highlands Community College]	1	0.0129051735
11	[Dabney S Lancaster Community College]	1	0.0141197056
12	[Northwest Iowa Community College]	1	0.0151541146
13	[Manhattan Area Technical College]	1	0.0172878854
14	[Lake Region State College]	1	0.0177886428

Figure 1: Sample of kNN results

(please note that this is not the entire dataset, the full dataset will be provided separately).

3.4 Maximizing the Return

Since our ultimate goal is to invest in schools that would offset poverty the most. We focused on the change in PES. The amount invested in each school should therefore maximize the increase in PES. This increase is our measurement for ROI.

The set up for this problem is outlined in sections 2.4-2.5, using equations (10) and (11) and the aggregation of all the school's variables into a single list we can get the following:

$$\begin{aligned}
 &\underset{\forall n_j}{\text{maximize}} && \frac{1}{\Phi} \sum_{j=1}^{\Phi} \Delta_j n_j \\
 &\text{subject to} && \sum_{j=1}^{\Phi} n_j \cdot T_j = 100,000,000 \\
 &&& \forall i, j \quad 1 \leq n_j \cdot T_j \leq g_j \phi_j
 \end{aligned}$$

We used the linear-optimization method described in Winston [10], in order to solve this problem. Once done, the results for the top 20 school as an output of our method is seen in Figure 2.

The overall return on investment that we received for this model was 1.01%. It is also important to note that we applied the same process to the unpartitioned version of this model. When doing that we saw that many variables were dropped in the process (due to being overwhelmingly null) and the return on investment dropped to 0.75%.

	Name	PES Rank	Investment
1	Georgia Institute of Technology-Main Campus (1397...	1.9459e-03	16036000
2	University of Colorado Denver (126562)	1.3750e-03	14735000
3	North Central Kansas Technical College (155593)	9.3614e-04	494500
4	Nashua Community College (183141)	7.8803e-04	3336700
5	Bates Technical College (235671)	7.8307e-04	7523200
6	Anoka Technical College (172954)	6.4288e-04	1313500
7	Clayton State University (139311)	6.1314e-04	2999600
8	Quincy College (167525)	5.4067e-04	3566800
9	Washburn Institute of Technology (155423)	5.1475e-04	833090
10	Parkland College (147916)	4.6194e-04	9539200
11	Indiana University-Purdue University-Fort Wayne (15...	4.0705e-04	12738000
12	Blinn College (223427)	2.2880e-04	5871400
13	Southeast Technical Institute (219426)	1.4047e-04	461120
14	Community College of Aurora (126863)	1.1598e-04	976160
15	Lone Star College System (227182)	1.1406e-04	11166000
16	Clackamas Community College (208406)	8.9894e-05	2764100
17	Eastern Oklahoma State College (207050)	6.0287e-05	213570
18	Mt Hood Community College (209250)	1.3525e-05	219370
19	Salt Lake Community College (230746)	2.0437e-06	237870
20	Tarrant County College District (228547)	1.5481e-06	234310

Figure 2: Top 20 Schools

4 Discussion

4.1 Sensitivity Analysis

In order to test model stability, we impose a radius condition to the kNN selection algorithm. The model testing outlined in Section 3 used an unrestricted subset of the data. We compute the the mean and standard deviation for each class i . Then each value exceeding k from $k = 1, \dots, 3$ standard deviations away was omitted from the data passed to the investment algorithm. The goal is to partition the dataset such that the subset data will have 68%, 95%, and 99.7% of the data.

With all other variables the same the ROI for each standard deviation(STD) is:

1. STD 1 ROI: 0.0353
2. STD 2 ROI: 0.0499
3. STD 3 ROI: 0.0474

All ROI are relatively the same suggesting that the model does well with any size of data input.

4.2 Strengths

- **Wide Variety of Variables** The investment strategy looks at all variables provided in the COMAP College Score Card data set. This approach provides a large number of factors to assess an institution's potential for investment rather than the usual ones used by different education foundations.

4.3 Weaknesses

- **Institutions with Missing Data** During the early stage process of the data cleaning, schools that did not have data for the two bench mark variables "Incoming Poverty" and "Outgoing Poverty" were removed from the list. It is a possibility that those particular institutions could have skewed the results.
- **Specialized Institutions** Similar to the point above, the data was filtered to remove any institutions that identified as religious-affiliated, race-affiliated, and gender-affiliated. In our first initial tests, many of the specialized institutions were at the top of the ranking systems due to the special nature of their student body. However the trade-off is smaller data set for the investment strategy.

5 Extensions

- Determine the correlation between tuition-grants and students' retention rates
- Refine the number of observed variables used for the kNN method, run a regression model that would determine the most relevant variables to use as predictors
- Run the model on major types
- Determine the true poverty rate by adjusting for major type as

6 Conclusions

7 Letter to Mr. Alpha Chiang

Dear Mr. Alpha Chiang,

References

- [1] Trombitas, Kate. *Financial Stress: An Everyday Reality for College Students*. Inceptia, July 2012. PDF.
- [2] US Census Bureau *The Big Payoff: Educational Attainment and Synthetic Estimates of Work-Life Earnings*. Report, 2006 . PDF.
- [3] Carnevale, Anthony, Ban Cheah, and Martin Van Der Werf. *Ranking Your College*. Georgetown University Center on Education and the Workforce, Dec. 2015. PDF.
- [4] Conkey, Christopher 2006. "Big donations to charity often include spending advice *The Wall Street Journal Asia* 05:
- [5] "Bill & Melinda Gates Foundation." *Bill & Melinda Gates Foundation*. Web. 29 Jan. 2016.
- [6] "Grants Database." *Grants Database*. Web. 29 Jan. 2016. <<https://www.luminafoundation.org/grants-database/strategy/student-financial-supports>>.
- [7] Hastie, Trevor, Robert Tibshirani, and Jerome Friedman. "Prototype Methods and Nearest-Neighbors". *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. 2nd ed. New York: Springer, 2009. Print.
- [8] "IPEDS Data Center." *IPEDS Data Center*. Web. 30 Jan. 2016. <<https://nces.ed.gov/ipeds/datacenter/DataFiles.aspx>>.
- [9] Hanushek, Eric and Margaret Raymond. 2005. "Does School Accountability Lead to Improved School Performance? *Journal of Policy Analysis and Management* 24(2): 297-329.
- [10] Winston, Wayne L. "Linear Programming." *Operations Research: Applications and Algorithms*. 4th ed. Belmont, Calif.: Duxbury, 2003. Print.
- [11] "Using Federal Data To Measure And Improve The Performance Of U.S. Institutions of Higher Education" Sept 2015. <<https://collegescorecard.ed.gov/assets/UsingFederalDataToMeasureAnd-ImprovePerformance.pdf>>.