
Freshman Seminar Assignment Problem Final Report

Keenan Gao

Binghui Ouyang

Hanwen Zhang

Yiming Zong

Department of Mathematical Sciences
Carnegie Mellon University
Pittsburgh, PA 15213

Abstract

As the number of college students increases, an optimization algorithm that can automatically assign students to classes is becoming a pressing need in many universities. In attempts of solving the algorithm, various questions arised. What is the objective function to optimize? What are the aspects that need to be considered? In this final report, we will deliver solutions to the Freshman Seminar Assignment Problem, our team's final project for Operations Research II, Fall 2014. We will first briefly describe the real-world problem and create a basic mathematical model based on the data. Then, we will feed the pre-processed data into several different algorithms with various objective functions and run-time constraints, and compare their final results. Eventually, we will discuss some generalizations of our current problem and propose potential algorithms for solving them. Our algorithms improve the quality of assignment significantly in terms of student preference compared to the traditional manual approach.

1 Problem Overview

The assignment problem we are solving in this project was initially broached by Dietrich College of Humanities and Social Sciences at Carnegie Mellon University. The main goal is to assign freshmen to mandatory seminars in a way such that all seminars are filled, and that students are assigned to the seminars that they are interests in. In the real-world dataset that we are provided, there are 308 students and 22 seminars. Each seminar can have at most 16 students, and each student must enroll in exactly one seminar.

To facilitate the matching process, students are asked to rank one “first-choice” seminar and three “second-choice” seminars. We associate costs with assigning a student to a seminar based on the student’s preference for the seminar. To find the optimal assignment, we approached the problem in two ways – minimizing the total cost across all seminars and smoothing out enrolled students’ preference for each seminar.

2 Mathematical Model

Based on the problem description in *Section 1*, we can build a mathematical model for it in order to describe an algorithm for solving the problem with any valid input data.

2.1 General Input

- n : Number of students ($n > 0$);
- m : Number of seminars ($m > 0$);
- k : Max number of selections that a student can make ($1 \leq k \leq m$);
- $s_{i,j}$: The j^{th} selection of i^{th} student, where $1 \leq i \leq n$ and $1 \leq j \leq k$. $s_{i,j} = 0$ when the Student i makes no corresponding choice for Rank j ;
- q_k : The quota for k^{th} seminar, where $1 \leq k \leq m$.

2.2 Input Constraints

- Positivity: $n, m, k > 0, \forall k \in \{1, \dots, m\}, q_k > 0$;
- Number of selections for student is bounded by number of available seminars: $k \leq m$;
- (?) Student rankings are valid and unique: $\forall (i, j), 1 \leq s_{i,j} \leq m$. And, for each i , all non-zero entries $s_{i,j}$ ’s take unique values.

2.3 Decision Variables

- $Y_{i,j}$: Indicator variables for whether Student i is assigned to Seminar j , where $1 \leq i \leq n$ and $1 \leq j \leq k$;

2.4 Data Pre-Processing

In order to deal with cases when a student is only willing or allowed to rank $k' < m$ seminars, we automatically set all “unassigned” priorities to $(k + 1)$. Also, we change the representation of students’ preference from $(\text{student}, \text{ranking}) \mapsto \text{seminar}$ to $(\text{student}, \text{seminar}) \mapsto \text{ranking}$ to make further calculations easier, i.e.

$$X_{i,j} = \begin{cases} l & \text{If Student } i \text{ ranked } j \text{ as } l^{\text{th}} \text{ option, or } s_{i,l} = j \text{ for some } l \in \{1, \dots, k\} \\ M & \text{If Seminar } k \text{ is not on Student } i\text{'s list, or } s_{i,l} \neq j \text{ for all } l \in \{1, \dots, k\} \end{cases},$$

where M is an arbitrarily large value in order to discourage the algorithm from assigning a student to a seminar that s/he did not list.

2.5 General Constraints

- $Y_{i,j}$ ’s are indeed indicator variables: $\forall (i, j), Y_{i,j} \in \mathbb{Z}, Y_{i,j} \geq 0, Y_{i,j} \leq 1$;

- Each student is assigned precisely one seminar: $\forall i, \sum_{l=1}^m Y_{i,l} = 1$;
- Each seminar is within enrollment quota: $\forall j, \sum_{l=1}^n Y_{l,j} \leq q_j$;

3 Approach for Various Heuristic Functions

Due to the flexibility of the original problem, we are proposing different objective functions for optimization, including minimizing the total “rank” given by the students, maximizing the number of student getting their top λk choice (where $\lambda \in (0, 1)$), etc. In the following sub-sections we present our approach for each heuristic in mathematical terms.

3.1 Minimize Total Rank of Students

In this case, our goal is to minimize the sum of all student rankings for their assigned seminars. To do so, our objective is to minimize $W = \sum_{i=1}^n \sum_{j=1}^m X_{i,j} Y_{i,j}$.

3.2 Minimize Variance of Students Preference in Each Seminar

In addition to minimizing total rank of student, it would also be helpful if we could “balance out” students’ preference of their assigned seminar for each seminar. For example, we do not want to have an assignment where some seminar has all students listing it as their first choice, yet some other seminar has none of the students listing it in their choices at all. With this idea in mind, for each seminar we wish to enforce a hard limit on the number of students enrolled that put it as first tier, second tier, etc.

Meanwhile, the hard limit can also affect the optimal solution that minimizes total rank of students in the final assignment. Therefore, care needs to be taken while picking the hard limits.

4 Exact Algorithm

The problem can also be reduced to an *Assignment Problem* when we include “dummy seminars” and “dummy students”. When we apply *Hungarian Algorithm*, we can obtain an absolute optimal solution that minimizes the total “cost” of students. Details about the algorithm are as follows:

4.1 Data Pre-Processing

We start with the matrix $X_{i,j}$ as obtained in *Section 2.4*. For each column that represents Seminar j , we create extra $(j - 1)$ dummy seminars by duplicating the same column q_j times. After that, we make our cost matrix square by adding zero rows at the bottom of the cost matrix. This gives us a matrix that we may feed into *Hungarian Algorithm*.

4.2 Hungarian Algorithm

Given the square cost matrix from previous section, we may simply apply *Hungarian Algorithm*, which returns a *student-seminar* assignment with minimal total cost. The algorithm completes in polynomial time [reference], and given the result we may simply assign each student to the actual seminar that the assignment corresponds to.

5 Approximation Algorithm

While the previous exact algorithm minimizes the sum of ranks of students across all seminars, it does *not* balance out the preference of students in different seminars. In order to do so, a heuristic for approximation is to place an artificial hard quota, Q , on the number of students in any seminar that places it as its first option. While the optimal solution for different values of Q can be obtained by running Hungarian algorithm with different Q ’s by creating different combinations of dummy seminars, its runtime (more than one hour) makes it undesirable.

In this case, our solution is to use randomized greedy algorithm for arbitrarily number of times, and return the best solution to the user. Following is the detail of the algorithm:

5.1 Ranking-Based Greedy Algorithm

For the greedy algorithm, we first satisfy (a portion of) all students' first choices, then second choices, and so on. Depending on the "popularity" of each seminar, we may limit the number of students allowed to be added to a seminar at each ranking. The algorithm (as *Algorithm 1* on next page) is outlined as follows, and it can be run multiple times in order to select an assignment with least amount of students that are not assigned to their ranked list.

Algorithm 1 Ranking-Based Greedy Algorithm

Output: $\text{asn}_i \leftarrow$ seminar assignment for Student i based on greedy algorithm
for $r = 1$ to k **do**
 for $i = 1$ to m **do**
 $\text{pool}[i] \leftarrow \{\text{unassigned student } s \mid s \text{ listed seminar } i \text{ as } r^{\text{th}} \text{ choice}\}$
 $\text{pool}[i] \leftarrow$ random subset of itself with certain size limit (e.g. seminar quota)
 Merge each $\text{pool}[i]$ into asn
Fill in still unassigned students

6 Summary of Results

So far we have tried two algorithms to solve the seminar assignment problem, namely the Exact Algorithm which applies the Hungarian Algorithm, and the Approximation Algorithm which utilizes randomization. Both algorithms are tested on the real data for the incoming class of Dietrich College for Year 2013 ($n = 308, m = 22$). The cost for assigning a student to a seminar is $g(x) = 2x^2$, where x is the tier that student ranks the seminar; if the student does not rank the seminar, the cost would be $g(x) = 10\,000$, which is a huge constant inspired by Big-M method.

6.1 Exact Algorithm

We use the Hungarian Algorithm to find the student-seminar assignment with the minimal total cost. The algorithm gives a very satisfactory assignment within five minutes. To be specific, among all the 308 students, 207 of them are assigned to their first choice of seminars, consisting of 67.2% of the population; and 82 are assigned to their second choice of classes, consisting of 26.6% of the population. Therefore, a total of 93.8% of the students are assigned to either their first or second choice of seminars. The remaining 19 students do not provide any preference of seminars, so they are randomly assigned to seminars which are not full. We also get a satisfying balance among seminar enrollments. Among all the 22 existing seminars, 18 of them get full enrollment of 16 students (16 was the highest enrollment allowed in 2013); one of the seminar get 15 students; one get 5 students; and the other two left seminars get 0 students enrolled. For the seminars that get 0 student, we would simply cancel it. The total cost of the optional solution is 1 900 164.

6.2 Approximation Algorithms

The Approximation Algorithm returns a quite satisfying result with 10000 iterations under two minutes and a half. Among the 308 freshmen students, 151 of them are assigned to their first choices, which consists of 49.0% of the total population of students; 132 are assigned to their second choices, which consists of 42.9% of the population; 5 students are assigned to their later choices; and 19 of them do not show any preference of seminars and so are assigned to a random seminar which is still not full. Therefore, a total of 91.9% of the students who indicated preferences are assigned to either their first or second choice. Regarding the balance of enrollments among the seminars, 15 of the seminars get full enrollment of 16 students; 3 of the seminars get pretty good enrollments of 15, 14 and 13 students respectively; and 3 of the seminars are poorly enrolled with 7, 5 and 2 students respectively. The total cost of the optional solution is 2 200 270.

6.3 Comparison & Analysis

According to the results of the Exact Algorithm and the Approximation Algorithm, both are able to provide us with quite satisfying assignments. More than 90% of the students are able to get in their first or second choices of seminars. Also, the class enrollments are pretty balanced.

However, both algorithms have their own pros and cons. The Approximation Algorithm has a faster and scalable performance compared to the Exact Algorithm because it uses randomization instead of calculating the best result. Therefore, when the dataset gets really large, it's better to choose the Approximation Algorithm. However, the relative performance of the Approximation Algorithm cannot be guaranteed compared to the Exact Algorithm, and it is absolutely not as good as the Exact Algorithm which always gives the minimum cost. Therefore, when the dataset is not huge and the time of using the Exact Algorithm is realistic, the Exact Algorithm should be a better choice as it provides us with the best solution.

6.4 Limitations

Despite the satisfactory performance of the exact algorithm, its run-time complexity $\mathcal{O}(n^3)$ makes the algorithm undesirable for $n > 1000$. Therefore, further work should be focused on more effective heuristics for the approximation algorithm in order to improve the quality of result under a reasonable time constraint.

7 Further Work & Enhancements

7.1 Allowing flexible input parameters

For our implementation, the main test case is for seminars with the same enrollment quota, and that each student makes one “first-choice” and three “second-choices”. However, our implementation of algorithm in *Sections 4-5* also supports seminars with various quotas and also students that rank arbitrary number of seminars in arbitrary number of tiers. This makes our implementation applicable to much more real-world cases because in general the course sizes need not be the same, and that student should be given the opportunity to make flexible selections.

7.2 Supporting bi-directional preference/cost parameters with Stable Marriage Algorithm

One of the features we could implement in the future might be that the final assignment does not only depend on the students’ rankings on seminars but also the seminars’ rankings on students. For example if a certain seminar is very major-orientated and prefers history major students mostly, then we will want to enhance our algorithm by accommodating for the bi-directional preferences. Algorithms like the Stable Marriage Algorithm would be useful to base on, or we may simply use the sum of cost in both students’ and seminars’ perspectives as parameters in the algorithms specified in *Sections 4-5*.

7.3 Stable Assignment Optimization

Similar to the principle of *Stable Marriage Problem*, in the final seminar assignment we do not want to have two students A and B , such that A prefers B ’s section, and also vice versa (we call those two students *rogue pair*). This can be done by scanning each pair of students and fixing every *rogue pair*. The algorithm is outlined as follows:

7.4 Parallelize Approximation Algorithm

Since the approximation algorithm in *Section Five* depends on a user-defined number of independent trials, we can actually parallelize the algorithm by using multiple threads, such that each worker is able to make trials and aggregate the result to the *Master* node. For a modern machine with multiple CPU cores, this is able to reduce the algorithm runtime by at least 50%.

Algorithm 2 Rogue-Pair Fixing Algorithm

Input: $\text{asgn}_i \leftarrow$ current seminar assignment for Student i

Output: (i,j) if we found a rogue pair, otherwise null

```
function FINDROGUEPAIR( $\text{asgn}$ )  
  for  $i=1$  to  $n$  do  
    for  $j=i+1$  to  $n$  do  
      if Student  $i$  and  $j$  prefer each other's seminar then return  $(i,j)$   
  return null
```

Input: $\text{asgn}_i \leftarrow$ seminar assignment for Student i based on greedy algorithm

Output: asgn_i : rogue pair-free assignment for Student i

```
 $p \leftarrow \text{FindRoguePair}(\text{asgn})$   
while  $p$  not null do  
   $(i,j) \leftarrow p$   
   $\text{asgn}_i \leftrightarrow \text{asgn}_j$   
   $p \leftarrow \text{FindRoguePair}(\text{asgn})$ 
```

Acknowledgements

The authors would like to thank Professor Alan Frieze for proposing the Hungarian algorithm in *Section 5* and for holding weekly meetings to track our progress. Also, we thank Professor Brian W. Junker, Professor Joseph E. Devine, and Gloria P. Hill for providing us with the real seminar assignment data for Year 2013. And, the thank goes to Brian Clapper for his Python implementation of Hungarian algorithm, which made our final results and analysis possible.

References

- [1] Martello, Silvano, and Paolo Toth. Knapsack Problems: Algorithms and Computer Implementations. Chichester: J. Wiley & Sons, 1990. Print.