

# Will I Get in? - Modeling the Graduate Admission Process for American Universities

Narender Gupta  
University of Illinois at  
Urbana-Champaign  
ngupta18@illinois.edu

Aman Sawhney  
University of New Mexico  
asawhney@unm.edu

Dan Roth  
University of Illinois at  
Urbana-Champaign  
danr@illinois.edu

## ABSTRACT

We study the graduate admission process in American universities. Our goal is to build a decision support model that provides candidates with pertinent information as well as the ability to assess their choices during the application process. This model is driven by extensive machine learning based analysis of large amounts of historic data available on the web. Our analysis considers factors including, but not limited to, standardized test scores and GPA as well as world knowledge regarding university *reputation* and *similarity* with other universities. The learning problem is modeled as a binary classification problem with latent variables that account for unavailable information, such as multiple graduate programs within the same institution.

An additional contribution of this paper is the creation of a new dataset of more than 25,000 student application data, covering hundreds of universities over several years. This dataset allows us to develop models that provide insight into student application behavior and university decision patterns. Our experimental study reveals some key factors in the decision process of programs and, consequently, allows us to propose a recommendation algorithm that provides applicants the ability to make an informed decision of programs to apply to given their profile, with high confidence of being accepted.

## CCS Concepts

•Information systems → Decision support systems; *Data mining*; •Computing methodologies → Machine learning;

## Keywords

Graduate Admissions, Decision Support, Learning Model, Latent Variables, Recommendations

## 1. INTRODUCTION

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

WOODSTOCK '97 El Paso, Texas USA

© 2016 ACM. ISBN 123-4567-24-567/08/06...\$15.00

DOI: 10.475/123.4

Every year, thousands of students apply to American graduate programs and in the process, discover that there is a dearth of reliable sources to aid them in making an informed decision. There are several sources that provide admission related statistics in aggregate form, but not tailored to individual profiles, thus, leaving the applicant with the only option of guessing and hoping for the best. Even though learning models have been used in variety of real-world applications, there is surprisingly little literature available on understanding admission dynamics and decision making.

A typical university application process involves the applicant submitting his transcripts, standardized test scores, a few letters of recommendation, a statement of purpose that expresses his aims, ambitions and research interests, and descriptive answers to a few additional questions. Test scores include GRE<sup>1</sup>, language test scores - such as TOEFL<sup>2</sup> or IELTS<sup>3</sup> etc. Since requirements, deadlines and the specific process to meet them is university specific, the applicant needs to first choose the universities he would apply to. Hence, an applicant is generally advised to start about an year in advance of the university deadlines to timely complete these requirements.

Given the uncertainty, a naive solution is to apply to a large number of universities. But, more applications imply a deeper investment of time and energy on the applicant's part. This also implies a greater monetary investment, which is a major concern for applicants from developing countries. An efficient strategy to circumvent this is to categorize the universities so that one can cover a wide spectrum by applying to a few representatives from each category. A popular scheme includes three categories: *Ambitious*: where the chances of admission are slim; *Reachable*: where the chances of admission are decent; *Safe*: where there is a fair certainty of being accepted. Multiple admission offers resulting from these decisions allow the applicant to choose, suboptimally, their best option.

The description of these categories is very subjective, and even more subjective is the applicant's ability to predict his probability of *Admit* i.e. chances of being admitted to a given program and, hence, the university categorization. This prediction is generally based on hearsay or semi-informed opinions, resulting in confusion and a waste of resources for both the applicant as well as the university.

In this paper, we address this problem by developing a machine learning approach which enables applicants to make

<sup>1</sup>Graduate Record Examination [www.ets.org/gre](http://www.ets.org/gre)

<sup>2</sup>Test Of English as Foreign Language [www.ets.org/toefl](http://www.ets.org/toefl)

<sup>3</sup>International English Language Testing System [www.ielts.org](http://www.ielts.org)

informed decisions by evaluating their chances of admission. We suggest a latent variable based generative modeling approach which is easily extensible. Another contribution of this paper is the dataset we provide which enables others to build upon the system we have created. Additionally, we analyze our system from student’s perspective but it can be easily extended to university’s perspective as well.

Some researchers have briefly reflected on the process of decision making, but only qualitatively [12]. The work done by Waters et al. models the problem from the university’s perspective [13]. They used a learning approach to aid the university admission committee by identifying the candidates that are unlikely to be offered admission. Their model is quite simplistic, considering the problem as a straightforward classification problem (via Logistic Regression) without attempting to reveal the diverse and rich patterns in the data. More importantly, their approach is university centric and does not provide any support to the decision process of applicants. The primary reason for the lack of such endeavours is the unavailability of a relevant dataset. One contribution of our work is the creation of such a dataset, which we will make available to the community. The dataset allowed us not only to determine the acceptability of an application but also suggest better choices.

Works such as Bruggink et al. and Moore et al. utilize domain knowledge to build statistical models [6, 10]. Bruggink et al. model undergraduate university admissions to a private liberal arts college [6]. The model treats application components as independent variables and assumes the decision to be dependent (*Admit* or *Reject*) on these. The independent variables include GPA<sup>4</sup>, SAT scores, other academic scores and extracurricular factors all of which have been quantized. Beyond strong statistical assumption, this approach assumes that the modelled university (and its application pool) provide a good representation of the whole distribution. Our study shows that this is not the case because different universities focus on different features, and hence produce decisions differently. We conclude that learning a decision model should be done separately for each university, if possible. Moore et al. model the problem with rule induction using ID3 algorithm [10]. Such an approach without care for bounded depth is prone to overfitting. Similar to Bruggink et al. the model is centered around one university and considers very small applicant sample size [6].

## 2. PROBLEM MODELING

### 2.1 Data

There are several online resources where applicants share their admission experiences; Edulix is one commonly used resource [9]. It is an active resource which hosts applicant profiles from all over the world. GRE scores, undergraduate university name, GPA, TOEFL scores and other accomplishments such as work experience and research publications pertinent to the graduate admissions are reported in the profile. In addition, users mention the universities that they applied to and the result of each application (*Admit*, *Reject* or *Result Not Available*).

We collected the data present on this website. Since the data is self reported, it had some erroneous reports, which we identified and removed by completely deleting any such

Table 1: Features of Data

General Features	
Total number of users before sanitization	36,207
Total number of users after sanitization	26,148
Features for CS related dataset	
Number of users	10,788
Application year range	[2001 2015]
Median Application Year	2013
Most Frequent Application Term	Fall
Number of universities with reported data	313
Number of applications per student (Mean)	6
Number of applications per university (Mean)	51
Number of undergraduate universities	2353
Degrees sought	[MS, PhD]

record. We also excluded any application that was not classified as either *Admit* or *Reject*. We observed that GRE and TOEFL scores have undergone various changes in grading scale over the years. Also, undergraduate institutions all over the world follow different scales for reporting GPA. As a standardization measure, we mapped these fields linearly to a scale of 0-100. An undergraduate university might be referred to by the differing names due to reasons such as usage of a popular acronym or spelling errors. We mitigated this problem by mapping the university names to their unique website URL.

In this paper we focus on modeling admission to computer science graduate programs, which form the plurality of our data. Our experiments are conducted only on this subset of the data. A few of the features of the resulting dataset for the computer science applications are in Table 1.

### 2.2 Supervised Learning

Each university offers binary decisions to applicants (*Admit* or *Reject*) and this decision is unaffected by the decision of other universities. Hence, the overall problem can be modeled as a set of individual binary classification problems. Supervised learning algorithms can be trained using features extracted from a labeled dataset and evaluated on prediction accuracy. However, accuracy, in this case, is a biased metric since most universities have very low acceptance rates, some of them being as low as 5%. In such a scenario, even by blindly rejecting all the applications the classifier will still evaluate to over 90% accuracy, without the need to learn anything. Also, acceptance rates vary across universities and thus accuracy will not be a true representative of the learning model’s performance. Hence, we chose F1 as our evaluation metric which combines precision and recall. Since number of admitted students is less than that of the rejected ones, we report our results for *Admit* label which provides a stricter bound on evaluations. F1 score takes into account not just the correct number of predictions made for *Admit* (*Precision*), but also the ratio of *Admit* students out of true *Admit* count (*Recall*).

$$F1 = \frac{2 \times Precision \times Recall}{Precision + Recall}$$

$$Precision = \frac{TP}{TP + FP}$$

<sup>4</sup>Grade Points Average

**Table 2: Classification Context**

		True Condition	
		Admit	Reject
	Predicted Condition	Admit	Reject
		TP	FP
		FN	TN

$$Recall = \frac{TP}{TP + FN}$$

**Feature Engineering:** The dataset contains several fields such as standardized test scores and academic history records. We extract several numerical features from these such as GRE test scores (AWA<sup>5</sup>, Verbal & Quantitative), undergraduate GPA, language test scores (TOEFL), as well as categorical features such as program applied to (e.g. MS, PhD), term (e.g. Fall, Spring) etc. These features are used for learning Logistic Regression, Support Vector Machine and Random Forest. Since most work available in the literature is confined to approaches we have mentioned so far, we'll regard it as the baseline for any novelties that we propose.

**Ensemble Learning:** To improve the system performance, we use ensemble learning. Training decision trees without bounded depth is prone to overfitting. But we can hope to generalize better if we use several limited-depth decision trees using partial data, and then feeding each decision as a feature into another regularized classifier. Constrained by the variance-bias trade-off, the second classifier captures variance of data through multiple underlying decision trees while keeping a limit on its own variance by choosing simple models such as linear separators e.g. soft-margin support vector machine. We create  $d$  such limited-depth decision tree classifiers where each classifier is trained on a bootstrapped sample from the original dataset. Corresponding to decision of each such classifier, we get a feature for the next classifier. In section 3.1, we perform grid search to find optimum value of  $d$  for all universities.

### 2.3 Generative Modeling

For a few universities, we noticed that our discriminative classifiers do not perform as expected, despite relatively large amount of data. we attribute it to the fact that some universities offer multiple degree programs that might target different kinds of applicants and have different admission criteria, e.g. professional master's versus thesis master's program at UIUC<sup>6</sup>. Certain other universities, such as CMU<sup>7</sup>, offer specific programs at same degree level such as master's degree in Machine Learning versus HCI<sup>8</sup> which fall under the purview of Computer Science. These distinctions, however, are not captured in the dataset, and are thus *hidden* from our models. This distribution causes difficulties for linear separators which assume that data is coming from single source. In case the data is coming from multiple sources such as in Figure ??, it is impossible to linearly separate the data using a single classifier. To accommodate these phenomena, we use an *Expectation-Maximization (EM)* algorithm that allows us to learn a model with a latent variable, capturing the missing information. The rest of this section, describes

how we model the problem in this case.

EM is one of the very effective techniques for finding a *Maximum Likelihood Estimate* with *hidden variables*. For the sake of brevity, we are not going into details of *EM*. We will be using a setup similar to the one Grove et al. used in [7].

We assume a hidden variable called Program Type,  $z$ , which might carry assume multiple values. These values can hold different semantics based on university, e.g. different for UIUC and CMU, and, hence, we refer to them simply by the value assigned by the model to the hidden variable such as  $1, 2 \dots k$  etc. Once we learn the most likely model with the hidden variable  $z$  taking  $k$  values we can use it in two different ways. In a soft-boundary setting, each student can be from either Program Type with some probability. In a hard boundary setting, the most likely cluster (program) completely owns the student record, and we can learn individual linear classifiers for each cluster. A student belongs to cluster  $z = i$  if:

$$P(student|z = i) > P(student|z = j), \forall j \neq i \quad (1)$$

where

$$\sum_j P(student|z = j) = 1 \quad (2)$$

We assume that an applicant *belongs* to a program  $z \in \{1, 2, \dots, k\}$  with probability  $\alpha_r = p(z = r)$ . Given this program, feature  $x_i$  of the applicant  $x$  is generated independently by a Gaussian distribution with model parameters  $(\mu_i, \sigma_i)$ .

Let us define a hidden variable,  $z \in \{1, 2, \dots, k\}$ . A student record (sample),  $x$ , consists of  $(n+1)$  features,

$$x = (x_0, x_1, x_2, \dots, x_n).$$

The likelihood of sample is given by

$$P(x) = \sum_z p(x|z)p(z) \quad (3)$$

Incorporating each feature probability, the likelihood of the data sample can be expressed as:

$$P(x) = \sum_z \left( p(z) \prod_i p(x_i|z) \right) \quad (4)$$

Starting with an initial set of parameters  $\theta$ , the probability that a data point  $x^j = x^j = (x_0^j, x_1^j, \dots, x_n^j)$  comes from each of the  $k$  values of  $z$  is given by:

$$P_r^z = p(z = r|x^j) = \frac{p(z = r) \prod_i p(x_i^j|z = r)}{\sum_z (p(z) \prod_i p(x_i^j|z = r))} \quad (5)$$

Let  $p_i^z = p(x_i|z = r)$ , then we can compute the expected log-likelihood as follows:

<sup>5</sup>Analytical Writing Analysis

<sup>6</sup>University of Illinois at Urbana-Champaign

<sup>7</sup>Carnegie Mellon University

<sup>8</sup>Human Computer Interaction

$$\begin{aligned}
E(LL) &\equiv E \left( \sum_j \log P(x^j | \alpha_z, p_i^z) \right) \\
&\equiv \sum_j E \left( \log P(x^j | \alpha_z, p_i^z) \right) \\
&\equiv \sum_j \left( \sum_z P_z^j \cdot \log P(z, x^j | \alpha_z, p_i^z) \right) \\
&\equiv \sum_j \left( \sum_z P_z^j \cdot \log(\alpha_z \cdot \prod_i p(x_i^j | z)) \right) \quad (6)
\end{aligned}$$

Assuming numerical features to be generated from a Gaussian distribution with parameters  $(\mu, \sigma)$ , above equation can be expanded as:

$$\begin{aligned}
E &\equiv \sum_j \left( \sum_z P_z^j \cdot \log \left( \alpha_z \cdot \prod_i \frac{1}{\sigma_i^z \sqrt{2\pi}} \exp \left( -\frac{(x_i^j - \mu_i^z)^2}{2(\sigma_i^z)^2} \right) \right) \right) \\
&\equiv \sum_j \left( \sum_z P_z^j \cdot \left[ \log \alpha_z + \sum_i \left( -\log \sigma_i^z - \frac{(x_i^j - \mu_i^z)^2}{2(\sigma_i^z)^2} \right) \right] \right) \quad (7)
\end{aligned}$$

Differentiating with respect to all parameters, we can find new values of  $\alpha_z$ ,  $\sigma_i^z$ , and  $\mu_i^z$ , for which the expected (log) likelihood receives an extremal value. Since,  $\alpha$  can only assume  $k-1$  independent values because of  $k$  states of  $z$ , we have:

$$\alpha_k = 1 - \sum_{i \in [1 \dots k-1]} \alpha_i \quad (8)$$

Using this equation and differentiating Eq (7) partially with respect to model parameters, we get following update rules:

$$\alpha_z = \begin{cases} \alpha_k \frac{\sum_j P_z^j}{\sum_j P_k^j} & \text{if } z \neq k \\ 1 - \sum_{i \in [1 \dots k-1]} \alpha_i & \text{if } z = k \end{cases} \quad (9)$$

$$(\sigma_i^z)^2 = \frac{\sum_j \sum_z P_z^j (x_i^j - \mu_i^z)^2}{\sum_j \sum_z P_z^j} \quad (10)$$

$$\mu_i^z = \frac{\sum_j \sum_z P_z^j x_i^j}{\sum_j \sum_z P_z^j} \quad (11)$$

Using above update rules, and various initialization schemes, we performed multiple experiments with EM. Current results are reported for  $z=2$  in a hard-boundary setting.

### 3. EXPERIMENTS

Our experimental study is designed to investigate the following issues:

- The ability to make reliable prediction on whether a specific student can be admitted to a given program.
- Our ability to identify sub-programs in a given university, and its significance on the performance of our admission model.
- Understanding the factors that contribute to admission.

- Understanding the differences among universities in terms of their admission decisions.

The rest of this section describes the details of our experimental study, and its results. For practical purposes, we selected a university for classification only if it has more than a certain number of applicant records. We chose the threshold of at least 10 admitted and at least 10 rejected students. In all our experiments we are reporting average F1 over 5-fold cross-validation.

#### 3.1 Discriminative Classifiers

We ran multiple experiments with various classifiers such as - SVM with linear kernel or Radial Basis Function kernel, Logistic Regression, Adaboost with decision trees, and Random Forest. We also experimented with all of the above classifiers by adjusting class weights according to sample frequency. Each of these setups is used in two different settings:

- Using simple features extracted from student records
- Training  $d$  decision tree classifiers with bounded-depth (depth=3) and then using predictions of these classifiers as features. Each decision tree is trained on 50% data, selected randomly.

Using grid search, we found that  $d=60$  yields maximum average F1 over all universities. These classifiers can provide the probability of the label as well which we utilize in Section 5. We observed that SVM with RBF kernel and Random Forest (with adjusted class weights) perform the best, hence we will report only their results for further experiments. These experiments were conducted using software implementation of Scikit-learn [11].

#### 3.2 Feature Ablation

These experiments were aimed at understanding the value in each feature. We trained multiple classifiers using single features and evaluated their performances. Then we iteratively added more features to each of the classifiers and evaluated gain in performance. Each feature, when considered individually, is the only classifying parameter. We call its corresponding F1 result the Discriminative Power of feature. It was observed that undergraduate GPA has the highest discriminative power and has an average F1 over all universities close to 0.65. Figure 4 shows how overall F1 increases if we sort the features into descending order of discriminative power and keep on adding them to the classifiers.

#### 3.3 EM

The EM model formulated in Section 2.3 was used to cluster students into different groups, representing potential programs. Subsequently, individual classifiers were learned for each of the clusters. Results for few of the universities for which significant growth was observed are listed in Table 5. As per the model assumption, EM bifurcates the data into two clusters ( $z=2$ ), each of which can be separated in a better way than the earlier cumulative cluster thereby increasing the performance of the models significantly. Improvement in F1 due to EM clustering is reported in Fig 5.

#### 3.4 Understanding Institution Rankings

This experiment explored the role of *reputation* or *rank* of the applicant’s undergraduate institution in graduate admissions. Ranks could be categorical or numerically ordered. Each university may have a different view of how they perceive an undergraduate institution and, hence, rank it differently. But these ideas do not have a firm ground and the following experiments analyzed this hypothesis.

First, we investigated if the notion of an undergraduate institution’s rank or category even exists. If it does, providing this extra knowledge should help improve the classifier’s performance. Our dataset has applicant records from a variety of departments belonging to thousands of undergraduate institutions across many countries. Since it is not practical to rank all of them in numeric order, we used information from other proxies to group them into categories. These proxies included, but were not limited to, rankings provided by US News [4], QS (Top Universities) rankings [3], Shanghai rankings [1], and other lists provided by various government or non-government agencies such as ‘List of Institute of National Importance in India’ [2].

Using these proxies, we created four categories:

- A- Institutions ranked as top tier and widely recognized.
- B- Institutions ranked in the middle tier or recognized regionally.
- C- Institutions ranked in the low tier.
- D- Institutions that are neither recognized nor ranked.

We want to emphasize that our goal is not to come up with an ‘*ideal*’ university rank-list but to evaluate the existence of such a list. This can be answered by answering two sub-questions:

1. Does a list provide any gain in discriminative power of classifier?
2. Does any other list, similar or random, provide any gain in discriminative power of classifier?

Our hypothesis was that if a largely agreed upon rank-list existed, and if our proxies are representative of such a list, then this rank-list should provide gain to the classifier. At the same time, any other list which deviates drastically from such a list should not provide comparable gain during classification.

This category distribution was referred to as ‘**Original Rank List**’ (ORL). ORL had following category distribution: A=47, B=217, C=363, D=2354. Next, we **consciously shuffled** this list using following rules:

1. A university can have either the same category as it originally had, or it can move to its closest category, e.g. B can move to either A or C. The probability of an institution moving to neighbor category is linearly proportional to the target size.
2. Each category still has the same number of institutions as it originally had.

Since, we shuffled the institution categories based on precise rules, we called the result as ‘**Consciously Shuffled Rank List**’ (CSRL). We were taking into account that ranking of a university varies with ranking agencies or regions. In addition, we maintained the original category distribution (size of category) of the institutions. ‘**Randomly**

**Shuffled Same Distribution Rank List**’ (RSSDRL) was created by assigning a randomly chosen category to each institution but by maintaining original category distribution. Finally, we created a ‘**Randomly Shuffled Uniform Distribution Rank List**’ (RSUDRL) by assigning a random category to each institution, without the constraint of maintaining original distribution. In RSUDRL, each category has uniform probability of occurring within the rank list. Results of these are shown in Fig 6.

### 3.5 Impact of Change in Application Year

In this experiment, we asked the question - Do universities change taste of students over time? Hence, we explored the change in decision to an application in a different application year. Some assume that since there is an increment in the number of applications every year, admissions become more competitive over time. We performed a carefully controlled experiment to test the validity of this hypothesis. In this setting, for every university:

1. Choose a training set (80%) and test set (20%), by random selection.
2. For each record in test set, record the application year, admission decision and prediction of classifier.
3. For each record in test set, change the application year (choose randomly between 2001 and 2015), and record the new prediction, using the classifier used in the previous step.
4. Perform this experiment for  $n(=100)$  iterations.

Our hypothesis was that if yearly factors do not have an effect then changing application year should not change the decision.

### 3.6 Which Universities Go Together

One of the unique features of our dataset construction is that applicant records capture various university combinations that the users apply to along with their results. This allowed us to find patterns, and formulate similarities among universities. Apriori algorithm [5] produces interesting results that are reported in table 8. But Apriori favors heavily populated universities over the less frequent ones. Hence, we expanded our experiments to include null-invariant measures. We computed similarity of two universities based on candidate acceptance using several null-invariant measures such as: AllConf, Jaccard, Cosine Similarity, Kulczynski coefficient, MaxConf defined in [8]. Results of the experiments are reported in table 7 and 8.

## 4. RESULTS

### 4.1 Discriminative Classifiers

To garner insights out of classification results for different universities, we needed to sort them in some order. To explore the trends i.e. the variation in prediction performance on universities of varied reputation, an ordering on the basis of reputation was required. We chose this order to be US News Graduate School Ranking, primarily, because

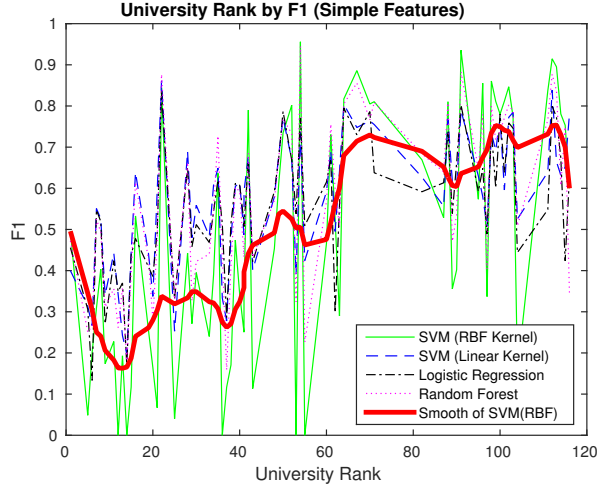


Figure 1: F1 scores for classification using supervised learning on simple features.

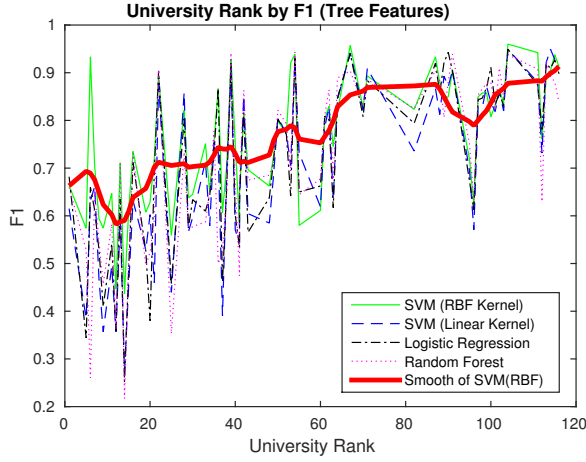


Figure 2: F1 scores for classification using supervised learning on bounded-depth decision tree features.

of its popularity among applicants. Whenever a rank was not available in this resource, a similar resource was consulted. In US News ranks, sometimes multiple adjacently ranked universities were stacked up on a single rank, and the resulting emptied out slots were left vacant. We flattened each stack to its nearest available slots. Fig 1 shows F1 on y-axis and university ranking on x-axis.

These results are for 69 universities having more than a thresholded number of applicants. For simplicity purposes, we treated these discrete values as a continuous function while plotting the graph. It can be observed that SVM with RBF kernel yields best results, hence, a smooth curve, using moving averages of span 5, of it is also plotted for easy visualization. Similar graph is plotted for ensemble learning with Decision Tree features in Fig 2. Fig 3 shows that Decision Tree feature classifier is a huge improvement over simple features.

As we can see from Fig 2 that we can get very good pre-

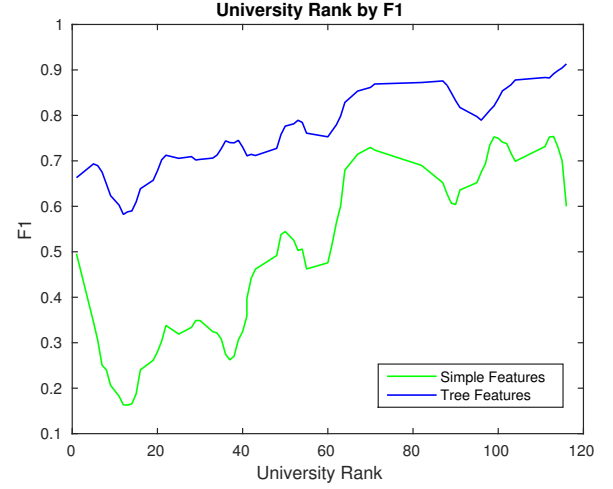


Figure 3: Comparison of Decision Tree feature classifier with baseline of simple feature classifier.

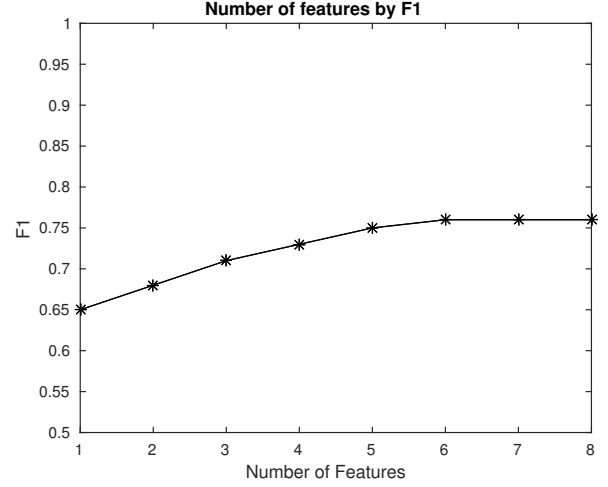


Figure 4: Cumulative F1 as we keep on increasing features on top of GPA

dictions for certain universities ( $F1 \geq 0.8$ ), while it is decent but not as good for certain others ( $0.6 \leq F1 < 0.8$ ), and low for certain other universities. But it is interesting to see the growth trend in F1 as we go from top-tier to lower-tier rank universities. The trend represents that we can classify with very high confidence for lower-tier universities using available features, in spite of lack of qualitative features such as Statement of Purpose or Letter of Recommendation, Research Publications etc. At the same time, top-tier universities might put more focus on these qualitative features.

## 4.2 Feature Ablation

We observed from experiments that GPA results in a higher F1 score compared to any other single feature for any university. Also, if we calculate F1 by excluding individual features, exclusion of GPA causes maximum loss in F1. Both of these observations lead to the conclusion that GPA has the

**Table 3: Discriminative Power of each feature**

Index	Feature Name	Cumulative Individual F1	F1
1	GPA	0.65	0.65
2	GRE Quant	0.68	0.53
3	GRE Verbal	0.71	0.58
4	GRE AWA	0.73	0.45
5	TOEFL	0.75	0.58
6	Program	0.76	0.31
7	Term	0.76	0.35
8	Previous Department	0.76	0.42

**Table 4: F1 without each feature. Less F1 due to missing feature indicates more discriminative power of that feature.**

Ignored Feature Name	F1
GPA	0.7234
GRE Quant	0.7415
GRE Verbal	0.7422
GRE AWA	0.7491
TOEFL	0.7455
Program	0.7654
Term	0.7647
Previous Department	0.7518

highest discriminative power among all available features. The result is intuitive as it validates the expectation that, broadly speaking, GPA is the prime factor in the admission process. It can also be seen that as the number of features that are considered is increased, performance goes up and is the highest when all the features are considered. Table 3 lists individual discriminative powers of each feature, as well as cumulative power for  $i$  features i.e. when features 1, ...,  $i$  are used for classification. Table 4 lists loss in F1 due to exclusion of each feature during classification. Fig 4 plots the cumulative discriminative power when we keep on adding features.

### 4.3 EM

Our model before the use of EM relied on the fact that data for each university is coming from a single source. The improvements in F1 as result of splitting the data according to EM formulation indicates that our model is able to capture underlying different distributions of source data. Also, since we know that UIUC offers different degree programs (Professional, Thesis), and CMU offers different specifications (Machine Learning, HCI etc) for the data reported as CS, it is probable that several other universities have more than one underlying distribution because of other factors.

<sup>9</sup>University of California Santa Cruz

<sup>10</sup>San Jose State University

<sup>11</sup>University of California Los Angeles

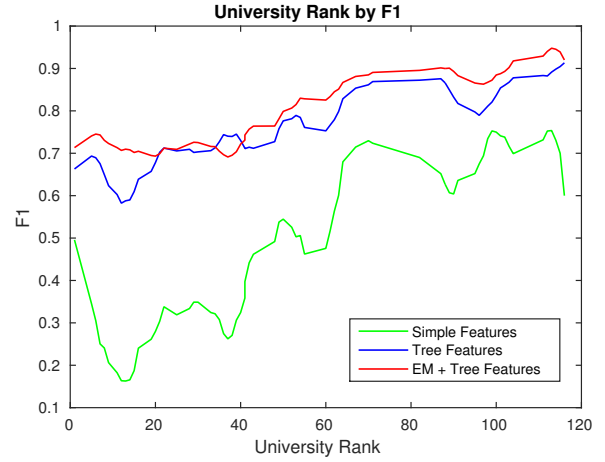
<sup>12</sup>University of Maryland College Park

<sup>13</sup>State University of New York Binghamton

<sup>14</sup>University of Texas Austin

<sup>15</sup>University of Colorado Boulder

<sup>16</sup>Texas A and M University College Station

**Figure 5: Improvement in F1 scores due to EM clustering****Table 5: Gain in F1 score due to EM clustering**

University	Tree	EM + Tree	EM Gain
UCSC <sup>9</sup>	0.58	0.84	0.26
SJSU <sup>10</sup>	0.62	0.86	0.24
UCLA <sup>11</sup>	0.45	0.68	0.23
UMD <sup>12</sup>	0.43	0.66	0.22
SUNY Binghamton <sup>13</sup>	0.76	0.94	0.17
UT Austin <sup>14</sup>	0.57	0.74	0.17
UC Boulder <sup>15</sup>	0.80	0.94	0.14
TAMU <sup>16</sup>	0.75	0.86	0.11
UIUC	0.57	0.65	0.08
CMU	0.66	0.71	0.05

Fig 5 shows overall increase in F1 over all of the universities using EM splitting in two clusters ( $z=2$ ). Table 5 reports significant improvement for some of the universities in our dataset by splitting data using EM.

### 4.4 Undergraduate Institution Rankings

Fig 6 indicates that the addition of rank of the undergraduate institution feature led to the gain in performance. We evaluated the gain in terms of statistical significance over 100 iterations and it was significant with  $p\text{-value} < 0.0001$ . This leads us to the conclusion that institution rank does play a role in the admission decision. Yet another interest-

<sup>17</sup>Arizona State University

**Table 6: Gain in F1 due to various rank-lists**

University	ORL	CSRL	RSSDRL	RSUDRL
ASU <sup>17</sup>	-0.0047	-0.1093	-0.1307	-0.2356
CMU	-0.0207	-0.7284	-0.0679	-0.2669
Brown University	0.3856	-0.8298	4.0007	-0.0333
Purdue	-2.4818	-1.4320	-1.2061	0.5790
UT Austin	5.7343	4.1302	3.5169	5.2184
Virginia Tech	4.0004	5.4014	5.7771	4.5617

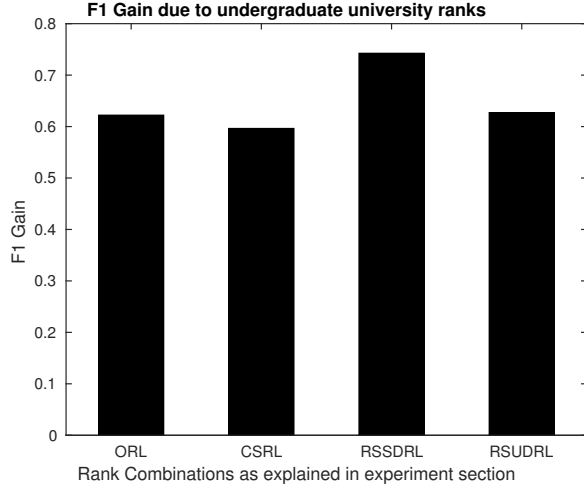


Figure 6: Gain in F1 due to various rank-lists

ing observation is that there is a comparable gain in all of the ways that we shuffled the rank-lists, consciously as well as randomly. This indicates that a popular rank-list provides only as much value as is provided by a random one.

Table 6 provides gain in F1 corresponding to each rank-list for some of the universities. We see that the universities show three types of behavior:

1. There is loss in F1 for each of the rank-lists e.g. ASU or CMU
2. Some rank-lists provide gain while others cause loss in F1 e.g. Purdue or Brown universities.
3. There is gain in F1 for each of the rank-lists e.g. UT Austin or Virginia Tech

A probable reason for this behavior (1) is that either the universities such as ASU or CMU don't use ranking system at all, and hence providing rank-lists causes the classifier to learn on irrelevant features causing a net loss, or none of the rank-lists is even close to the rank-lists used by these universities. For behavior (2), some of the rank-lists are close to the ones used by these universities while other rank-lists deviate drastically. For observation (3), all of the rank-lists are partially matching to the rank-lists used by these universities. Also, from 3<sup>rd</sup> and 4<sup>th</sup> rows in Table 6, it can be seen that one rank-lists provides gain to a specific university, while the other does the same for some other university.

Hence, it can be claimed that there is no single rank-list that every university uses, and there is no consensus over these ranks. So, two similar universities might respond differently to the same candidate even though all the other factors are consistent in his applications.

## 4.5 Impact of Change in Application Year

The experiment is aimed at testing the hypothesis that a change in the application year leads to a change in the decision. If an application was correctly classified initially then if there isn't any change in the competition of the application pool then it should still be correctly classified. Over 100 iterations, approximately 60K records were tested, out of which 55K were classified correctly. It was observed that

out of these, a vast majority (>98%) retained the same label even after the change of application year.

We define competition per application to increase if increase in application year changed the decision from *Admit* to *Reject*, and vice versa. Out of all of the decision changes, a record was assigned '+1' if it showed that the competition increased, and '-1' if the competition decreased. Overall sum of these scores for most of the universities was very close to 0. For all of the universities, the overall sum was -56, which means 56/55K, approximately 0.1% decrease in competition. Hence, it can be said that decision for an application depends solely on the university and the application, and not the application year.

## 4.6 Which Universities Go Together

As described earlier, applicants have the tendency to apply to universities in buckets (*Ambitious*, *Reachable*, *Safe*). This leads to the hypothesis that since applicants apply in buckets, it should be apparent through the similarity scores of the universities. As the results show this is infact the case. Kulczynski coefficient represents the average of conditional probability conditioned on each of the variables. Table 7 lists a few interesting associations based on the kulczynski coefficient and Table 8 lists such results for Apriori algorithm.

Table 7: Interesting similar universities based on Kulczynski score

University 1	University 2	Kulc
UChicago <sup>18</sup>	CSU <sup>19</sup>	0.286
UNCC <sup>20</sup>	UNLV <sup>21</sup>	0.303
CalTech <sup>22</sup>	UCR <sup>23</sup>	0.521
URI <sup>24</sup>	UWisc <sup>25</sup>	0.508

Table 8: Universities that go together based on Apriori algorithm

University 1	University 2	Support
UM Twin <sup>26</sup>	SUNY Stony <sup>27</sup>	218
UIC <sup>28</sup>	Indiana <sup>29</sup>	112
Cornell University	SUNY Stony	97
SUNY Buffalo <sup>30</sup>	GMU <sup>31</sup>	75

<sup>18</sup>University of Chicago

<sup>19</sup>Chicago State University

<sup>20</sup>University of North Carolina Charlotte

<sup>21</sup>University of Nevada Las Vegas

<sup>22</sup>California Institute of Technology

<sup>23</sup>University of California Riverside

<sup>24</sup>University of Rhode Island

<sup>25</sup>University of Wisconsin Madison



## 5. RECOMMENDATIONS

Since, now we have a system that can predict application decisions for a university, we can utilize it to aid students in making informed choices. In Section 4.6, we provide evidence that students apply to universities based on their notion of ‘Ambitious’, ‘Reachable’ and ‘Safe’ buckets. We include this notion into our algorithm to generate recommendations. Since the classifier system we have is not 100% accurate, it can generate erroneous recommendations if we simply classify for each university and return the results. Fig 5 shows that although the trend in university ranks is not strictly monotonous, it becomes very smooth if we cluster neighboring universities and then plot it. Hence, we cluster universities and employ multi-level classification to produce robust results while generating recommendations.

The first level of decision is coarse and the next level result is fine-grained. While classifying coarsely over a range of universities, we mix the records of all universities inside a cluster and train a single classifier on all of them. If the universities in the cluster are similar, the classifier learns the common patterns of admission versus rejection and provides a more general decision than any of the component universities. While in Fine-grained classification an individual classifier is trained for each university.

This algorithm consists of 5 steps:

### 1. University clustering

- Cluster similar universities together based on US News rankings e.g. Universities in rank [1,10] fall into cluster 1, universities from [11,20] fall into cluster 2 and so on.

### 2. Coarse classification

- Using coarse decision, we find the cluster that offers *Admit* and is closest to the top-tier universities. We call this cluster as ‘Reachable’ because it is the best ranked university cluster that can offer Admission.

### 3. Reachable Universities

- Perform fine-grained classification on each of the universities in ‘Reachable’, and return those which produce an *Admit* with highest probability.

### 4. Safe Universities

- We call the cluster next to ‘Reachable’ as ‘Safe’ because it also offers admission and does so with higher probability. Then fine-grained classification is applied on ‘Safe’ to report Safe universities.

### 5. Ambitious Universities

- For ‘Ambitious’ universities, we find those universities which are similar to the ones produced by ‘Reachable’ and ‘Safe’ but are towards the top

tier universities and hence do not offer admission. These similar universities are based on the higher similarity score based on common admissions.

In step 1, the benefit of using US News rankings is that as the universities get closer to top rank, probability of admission of any candidate decreases. We also verified the same observation from data. As a future work, there many clustering schemes can be employed here, including university similarity scores reported in section 4.6, as long as proximity of various clusters is known.

## 6. CONCLUSIONS

This paper studies the graduate admission process in American universities using a machine learning approach. Our goal is to build a decision support model that allows candidates to make informed decisions on which schools to apply to, what are their chances of admission, and a slew of other decision-related issues. We modeled the decision process as a learning problem and presented a system that can achieve high accuracy and can be generalized across multiple universities. By employing many approaches towards solving this problem, such as supervised learning and generative modeling, we prove that a mixture of approaches can provide better results than any of the individual approaches. We also provide a dataset that provides avenues for further research. This work can be extended in multiple ways such as towards improving accuracy or validating common notions. Some of the additions of this work may include expanding the EM formulation by modeling further variables such as undergraduate institution ranking mechanism. We proved that every university has a custom ranking mechanism. This mechanism can be modeled as a distribution which can, then, be assigned to one of the hidden variables in the EM model. Theoretically, such a model has more expressive power and can, thus, learn better regarding application decisions. We believe this is but a brisk start to the research that can be performed on the topic. Many more enhancements are possible by expanding the dataset and extracting richer data features from Letters of Recommendation or Statement of Purpose. By asking these questions and providing this dataset, we hope to initiate a discussion that can lead to better understanding of how academia accepts its new members.

## 7. REFERENCES

- [1] Academic Ranking of World Universities, 2015.
- [2] List of Institutes of National Importance in India, 2015.
- [3] University Rankings | Top Universities, 2015.
- [4] US News Best Graduate Schools, 2015.
- [5] R. Agrawal, R. Srikant, et al. Fast algorithms for mining association rules. *Proc. 20th int. conf. very large data bases, VLDB*, 1215:487–499, 1994.
- [6] T. H. Bruggink and V. Gambhir. Statistical models for college admission and enrollment: A case study for a selective liberal arts college. *Research in Higher Education*, 37(2):221–240, 1996.
- [7] A. Grove and D. Roth. Linear concepts and hidden variables. *Machine Learning*, 42(1/2):123–141, 2001.
- [8] J. Han, M. Kamber, and J. Pei. 6 - mining frequent patterns, associations, and correlations: Basic

<sup>26</sup>University of Minnesota twin cities

<sup>27</sup>State University of New York Stony Brook

<sup>28</sup>University of Illinois Chicago

<sup>29</sup>Indiana University-Bloomington

<sup>30</sup>State University of New York Buffalo

<sup>31</sup>George Mason University

concepts and methods. In J. H. Kamber and J. Pei, editors, *Data Mining (Third Edition)*, The Morgan Kaufmann Series in Data Management Systems, pages 243 – 278. Morgan Kaufmann, Boston, third edition edition, 2012.

- [9] S. Kassegne. Edulix - premier site for scholars - 'education crowdsourced'.
- [10] J. S. Moore. An expert system approach to graduate school admission decisions and academic performance prediction. *Omega*, 26(5):659 – 670, 1998.
- [11] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [12] K. Raghunathan. Demystifying the american graduate admissions process, 2010.
- [13] A. Waters and R. Miikkulainen. Grade: Machine learning support for graduate admissions. In *Proceedings of the 25th Conference on Innovative Applications of Artificial Intelligence*, 2013.