# Using random forests to study physics graduate school admissions

## Nicholas T. Young

Center for Academic Innovation, University of Michigan
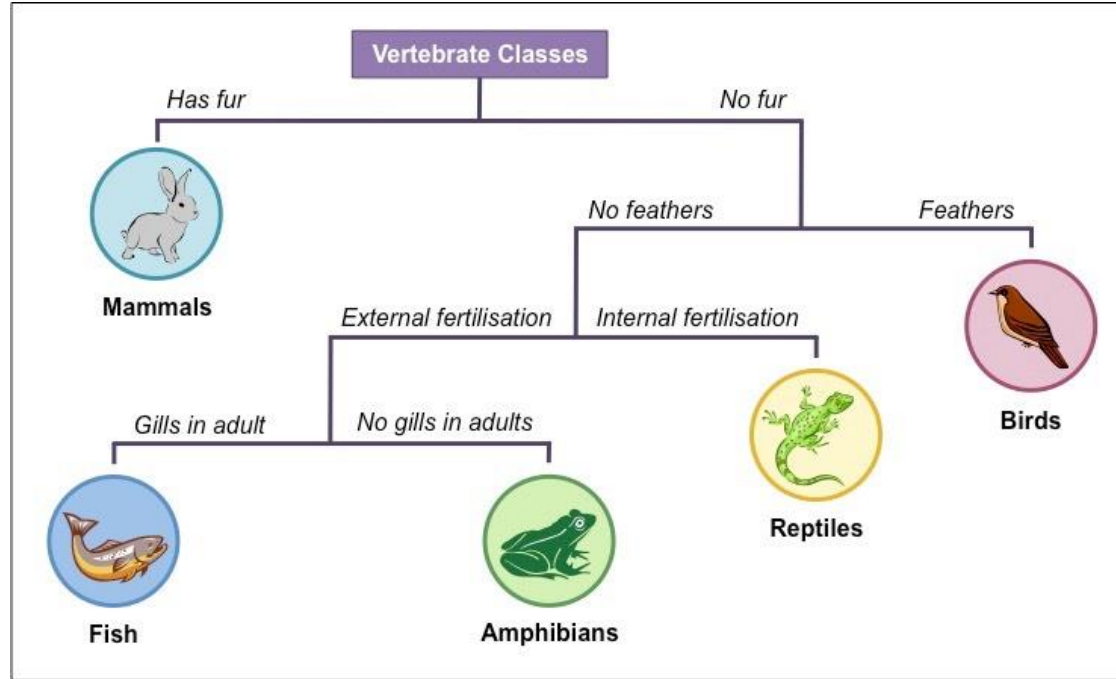
@NickYoungPER

ACADEMIC INNOVATION
UNIVERSITY OF MICHIGAN

# Dichotomous keys



Source: *BioNinja*

# Decision trees

# Decision tree example



**SHOULD I ATTEND THIS COMMITTEE MEETING?**

Am I on the committee?

YES — Attending

NO — Are there free donuts?

NO — Is the topic of interest?

YES — Attending!

YES — Attending

NO — Don't attend

ACADEMIC INNOVATION
UNIVERSITY OF MICHIGAN

# The random forest

| Run 1 | Run 2 | Run 3 | Run 4 | Run 5 |
|-------|-------|-------|-------|-------|
| V1 | V3 | V1 | V3 | V2 |
| V2 | V4 | V2 | V5 | V4 |
| V3 | V7 | V6 | V7 | V5 |
| V4 | V8 | V8 | V9 | V6 |
| V5 | V9 | V10 | V10 | V8 |

ACADEMIC INNOVATION
UNIVERSITY OF MICHIGAN

# The random forest

# Why use Random Forest

- No assumptions on the shape of the data

@NickYoungPER

ACADEMIC INNOVATION
UNIVERSITY OF MICHIGAN

# Why use Random Forest

- No assumptions on the shape of the data
- Scaling of continuous variables is irrelevant

@NickYoungPER

ACADEMIC INNOVATION
UNIVERSITY OF MICHIGAN

# Why use Random Forest

- No assumptions on the shape of the data
- Scaling of continuous variables is irrelevant
- Interested in predicting outcomes

ACADEMIC INNOVATION
UNIVERSITY OF MICHIGAN

# What Random Forest cannot do

- Be a magic solution to the problem

# Let's try it out!

ACADEMIC INNOVATION
UNIVERSITY OF MICHIGAN

# Data

GRE scores
GPA
Undergrad school
Research interest


N=512 applications
2014-2017

# The confusion matrix

|  |  | Model predicts | |
| --- | --- | --- | --- |
|  |  | True | False |
| What the true answer is | True | $N_{TT}$ | $N_{TF}$ |
|  | False | $N_{FT}$ | $N_{FF}$ |

ACADEMIC INNOVATION
UNIVERSITY OF MICHIGAN

# The confusion matrix

Model predicts

|  |  | True | False |
|---|---|---|---|
| What the true answer is | True | $N_{TT}$ | $N_{TF}$ |
|  | False | $N_{FT}$ | $N_{FF}$ |

$$Accuracy = \frac{N_{TT} + N_{FF}}{N_{TT} + N_{TF} + N_{FT} + N_{FF}}$$

ACADEMIC INNOVATION
UNIVERSITY OF MICHIGAN

@NickYoungPER

# Receiver Operating Characteristic (ROC) Curve

# Results

<u>All Variables</u>
Average Testing Accuracy:
75.6% ± 0.6%
Null accuracy: 52.7%

Average Testing
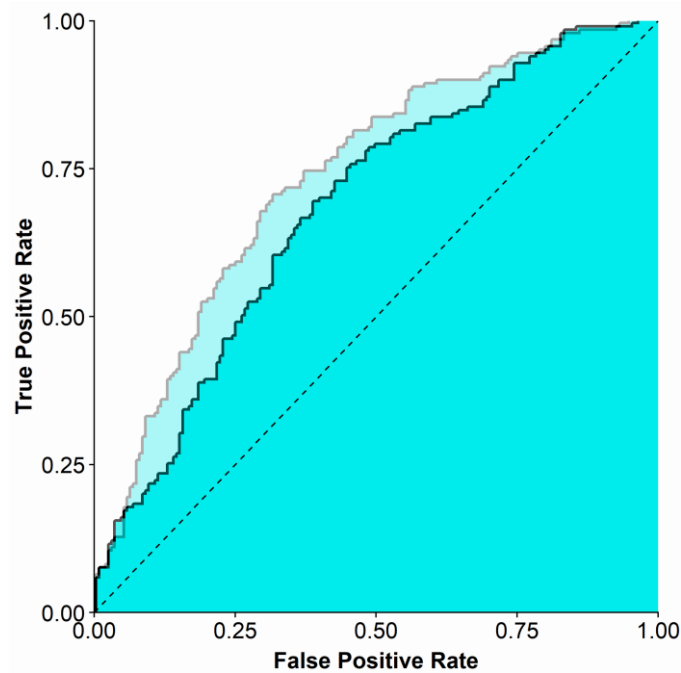Area Under the Curve (AUC):
0.756 ± 0.006

| Representative Run | | | |
| --- | --- | --- | --- |
| Is applicant admitted to the physics graduate program? | | Actual Decision | |
| | | Not Admitted | Admitted |
| Model Prediction | Not Admitted | **40.3%** | 14.9% |
| | Admitted | 9.1% | **35.7%** |

# Use variable importance to determine what is useful in making a prediction

ACADEMIC INNOVATION
UNIVERSITY OF MICHIGAN

Shuffle a variable ⟶ observe change in some metric (e.g. AUC) ⟶ Order by the change in the metric
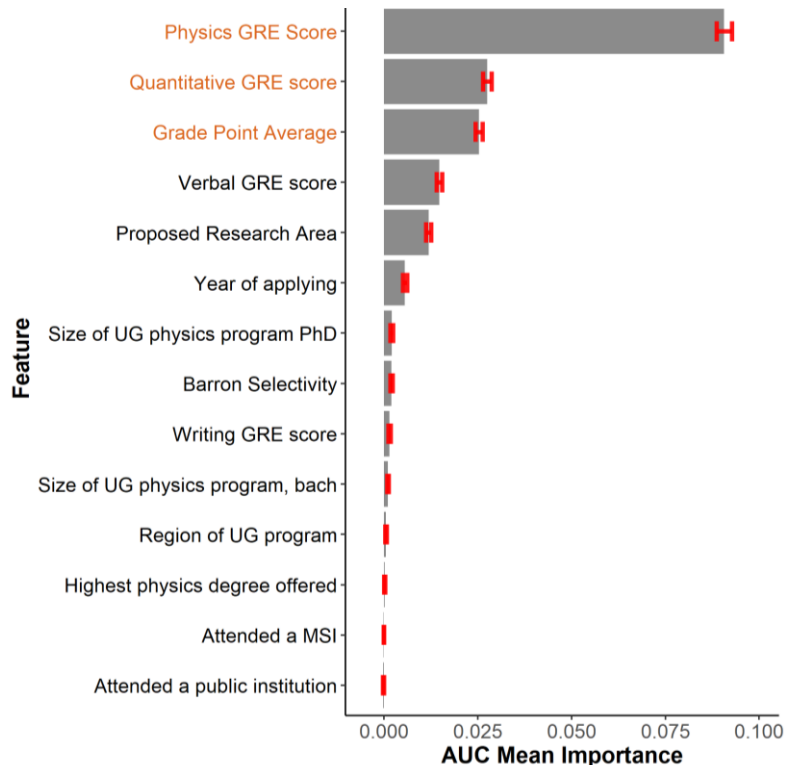
# Results

All Variables
Average Testing Accuracy:
75.6% ± 0.6%
Null accuracy: 52.7%

Average Testing
Area Under the Curve (AUC):
0.756 ± 0.006

ACADEMIC INNOVATION
UNIVERSITY OF MICHIGAN

# How do we know what matters?

ACADEMIC INNOVATION
UNIVERSITY OF MICHIGAN

# Results

All Variables
Average Testing Accuracy:
75.6% ± 0.6%
Null accuracy: 52.7%
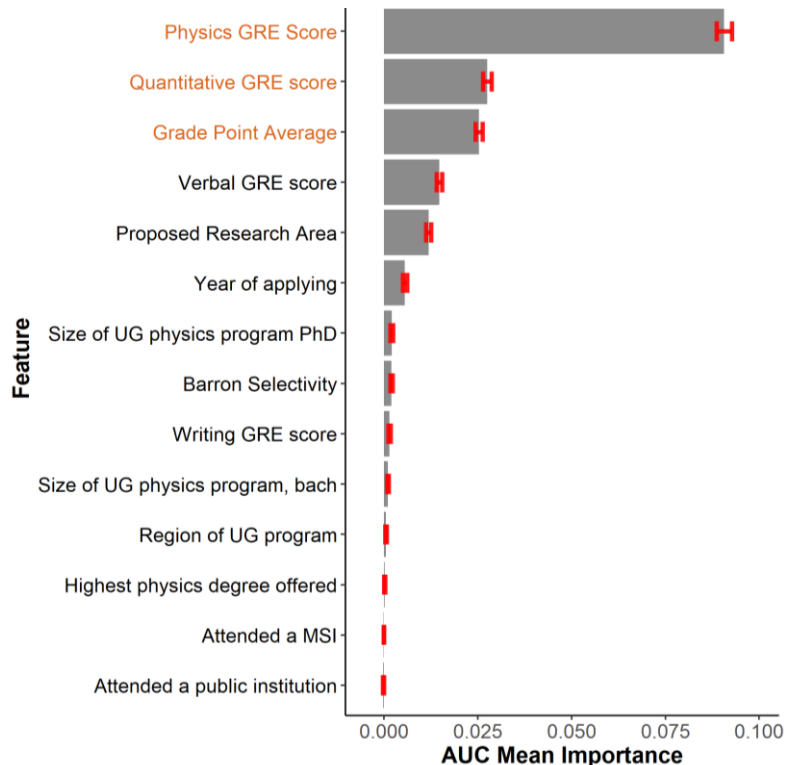
Average Testing
Area Under the Curve (AUC):
0.756 ± 0.006

Only Selected Variables
Average Accuracy:
75.4% ± 0.6%
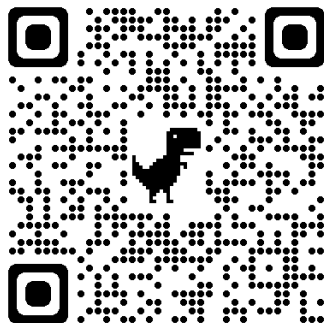Average Area Under the Curve:
.754 ± 0.006

# Learn more

Rubric-based holistic review represents a change from traditional graduate admissions approaches in physics

Nicholas T. Young,[1, 2, *] N. Verboncoeur,[1] Dao Chi Lam,[3] and Marcos D. Caballero[1, 2, 4, 5, †]

[1]Department of Physics and Astronomy, Michigan State University, East Lansing, Michigan 48824
[2]Department of Computational Mathematics, Science, and Engineering, Michigan State University, East Lansing, Michigan 48824
[3]Department of Statistics, Michigan State University, East Lansing, Michigan 48824
[4]Center for Computing in Science Education & Department of Physics, University of Oslo, N-0316 Oslo, Norway
[5]CREATE for STEM Institute, Michigan State University, East Lansing, Michigan 48824
(Dated: December 14, 2021)

Rubric-based admissions are claimed to help make the graduate admissions process more equitable, possibly helping to address the historical and ongoing inequities in the U.S. physics graduate school admissions process that have often excluded applicants from minoritized races, ethnicities, genders, and backgrounds. Yet, no studies have examined whether rubric-based admissions methods represent a fundamental change of the admissions process or simply represent a new tool that achieves the same outcome. To address that, we developed supervised machine learning models of graduate admissions data collected from our department over a seven-year period. During the first four years, our department used a traditional admission process and switched to a rubric-based process for the following three years, allowing us to compare which parts of the applications were used to drive admissions decisions. We find that faculty focused on applicants' physics GRE scores and grade point averages when making admissions decisions before the implementation of the rubric. While we were able to develop a sufficiently good model whose results we could trust for the data before the implementation of the rubric, we were unable to do so for the data collected after the implementation of the rubric, despite multiple modifications to the algorithms and data such as implementing Tomek Links. Our inability to model the second data set despite being able to model the first combined with model comparison analyses suggests that rubric-based admissions does change the underlying process. These results suggest that rubric-based holistic review is a method that could make the graduate admissions process in physics more equitable.

I. INTRODUCTION

While graduate school has historically been seen as a route for students to begin careers in academia, graduates are increasingly pursuing careers across industry, government, and academia. The National Science Foundation's Survey of Doctorate Recipients finds that less than half of all PhDs work at an educational institution while only 2 out of 5 physics PhDs do [1]. As such, universities have a duty to ensure that their students are able to achieve their chosen career trajectories.

Yet, the data suggests that isn't always the case. Only 3 out of 5 physics students who enroll in a PhD program will successfully complete their program [2, 3]. As undertaking graduate study involves a significant time and financial investment from both the student and institution, failing to ensure students graduate leads to a waste of resources. Solutions must consider both the admission and retention sides to this problem. In this paper, we will focus on the former.

As the Council of Graduate Schools notes in one of its reports, "Better selection [of graduate students] can result in higher completion rates" [4]. Historically and continuing to today, graduate school admissions in the US have tended to be an exclusionary process that favors certain groups over others. Previous research into the graduate admissions process in physics has found that the process relies heavily on the quantitative metrics such as grade point average (GPA) and General and Physics GRE scores [5–10]. These metrics have been found to benefit groups already overrepresented in higher education. For example, prior work has shown students from groups underrepresented in higher education (e.g., first generation, low income, Black, Latinx, Native) suffered a grade penalty relative to their more privileged peers with students from minoritized racial groups suffering the largest penalties [11]. Other work has shown that the General and Physics GREs are biased against women and students from minoritized racial and ethnic groups [2, 12] as well as against students from smaller or less prestigious universities [13]. Furthermore, the high costs associated with these often-required tests, despite limited evidence that these tests serve a predictive purpose [2, 14, 15], prevent some students from applying [16, 17].

The inequities in the admissions process and the fact that traditional admissions methods "miss many talented students" [18] have led various programs and organizations to consider alternative admission approaches such as holistic admissions, which considers a "broad range of candidate qualities including 'noncognitive' or personal attributes" [19]. These efforts are often supported by rubrics to ensure that all applicants are assessed on the

* Current email: ntyoung@umich.edu
† Corresponding Author: caball14@msu.edu

arXiv:2112.06886v1 [physics.ed-ph] 13 Dec 2021

arxiv:2112.06886

ACADEMIC INNOVATION
UNIVERSITY OF MICHIGAN

# Recap

- Random forest is a good technique if you know the outcome of your data, the data has a complex relationship to outcome (non-linear), and you are interested in predicting the outcome

ACADEMIC INNOVATION
UNIVERSITY OF MICHIGAN

# Recap

- Random forest is a good technique if you know the outcome of your data, the data has a complex relationship to outcome (non-linear), and you are interested in predicting the outcome
- Can also determine what features are most predictive of the outcome

# Recap

- Random forest is a good technique if you know the outcome of your data, the data has a complex relationship to outcome (non-linear), and you are interested in predicting the outcome
- Can also determine what features are most predictive of the outcome
- Will always get an answer; want to make sure it is a reasonable answer

ACADEMIC INNOVATION
UNIVERSITY OF MICHIGAN

# Recap

- Random forest is a good technique if you know the outcome of your data, the data has a complex relationship to outcome (non-linear), and you are interested in predicting the outcome
- Can also determine what features are most predictive of the outcome
- Will always get an answer; want to make sure it is a reasonable answer

Get in touch: ntyoung@umich.edu

ACADEMIC INNOVATION
UNIVERSITY OF MICHIGAN

# Resources & Recommended readings

- Breiman, L. (2001). Random forests. *Machine learning*, *45*(1), 5-32. https://doi.org/10.1023/A:1010933404324
- Breiman, L. (2001). Statistical modeling: The two cultures (with comments and a rejoinder by the author). *Statistical science*, *16*(3), 199-231. https://doi.org/10.1214/ss/1009213726
- Janitza, S., Strobl, C., & Boulesteix, A. L. (2013). An AUC-based permutation variable importance measure for random forests. *BMC bioinformatics*, *14*(1), 1-11. https://doi.org/10.1186/1471-2105-14-119
- Strobl, C., Boulesteix, A. L., Zeileis, A., & Hothorn, T. (2007). Bias in random forest variable importance measures: Illustrations, sources and a solution. *BMC bioinformatics*, *8*(1), 1-21. https://doi.org/10.1186/1471-2105-8-25
- Young, N. T., Allen, G., Aiken, J. M., Henderson, R., & Caballero, M. D. (2019). Identifying features predictive of faculty integrating computation into physics courses. *Physical Review Physics Education Research*, *15*(1), 010114. https://doi.org/10.1103/PhysRevPhysEducRes.15.010114
- Zabriskie, C., Yang, J., DeVore, S., & Stewart, J. (2019). Using machine learning to predict physics course outcomes. *Physical Review Physics Education Research*, *15*(2), 020120. https://doi.org/10.1103/PhysRevPhysEducRes.15.020120

ACADEMIC INNOVATION
UNIVERSITY OF MICHIGAN