



DONOR'S CHOICE

Capstone 2

Bill Murphy

8/12/18

WHO IS DONOR'S CHOICE?



- **Support a classroom. Build a future.**
- Teachers all over the U.S. need your help to bring their classroom dreams to life. Choose a project that inspires you and give any amount.
- “Our team of 80 has vetted and fulfilled over 600,000 classroom project requests that range from butterfly cocoons, to robotics kits, to Little House on the Prairie. Many of us are former teachers, so our operation feels like a cross between a startup and a schoolhouse.”

WHAT PROBLEM NEEDS TO BE SOLVED?

Small Group Learning for Big Success

My students need basic supplies for our classroom. Right now our classroom comes with a whiteboard and a bookshelf. I would greatly appreciate it if you are able to contribute to their learning!

My Students

I teach in an amazing elementary school in Lawrence Massachusetts. The school district was recently taken over by the state which means the funds are low and the stakes are high! Almost all my students come from low-socioeconomic homes, homes with trauma and many from single parent households. My goal is to provide a safe and loving learning environment where they know they can make mistakes.

My students did not grow up speaking English and all of them have a various learning disability, but they come to school ready to learn and ready to give it their all!

The goal of my project is to supply my students with everything they need in order to succeed in school, from pencils to books with



Ms. Lindholm

NEVER BEFORE FUNDED

Grades 3-5

South Lawrence East Elementary School
Lawrence, MA

Nearly all students from low-income households

[Remind me about this project](#)

- 186,000 proposals last year
- Each proposal is currently reviewed by a staff member → 2,325/staff member per year
- Number of proposals is growing rapidly, and the present system is unsustainable.
- Need proposals that help will reduce the load on staff

WHAT DATA IS AVAILABLE?

kaggle

581 total participants
Top roc_auc score:
0.82812

Trials top 3:
161, 112, 80

- Data source is a previously completed Kaggle competition
- 3 Types of Data Supplied:
 - Numerics: # previous projects, item costs, etc.
 - Categorical: subject, state, teacher title (Mrs., Ms., Mr., etc)
 - Text Fields: Descriptions of resource needs, essays
- Training Data:
 - 186,000 lines of project requests
 - 85% approval rate as baseline

BASELINE TO JUDGE SUCCESS

- Kaggle used ROC
- Better: Confusion Matrix

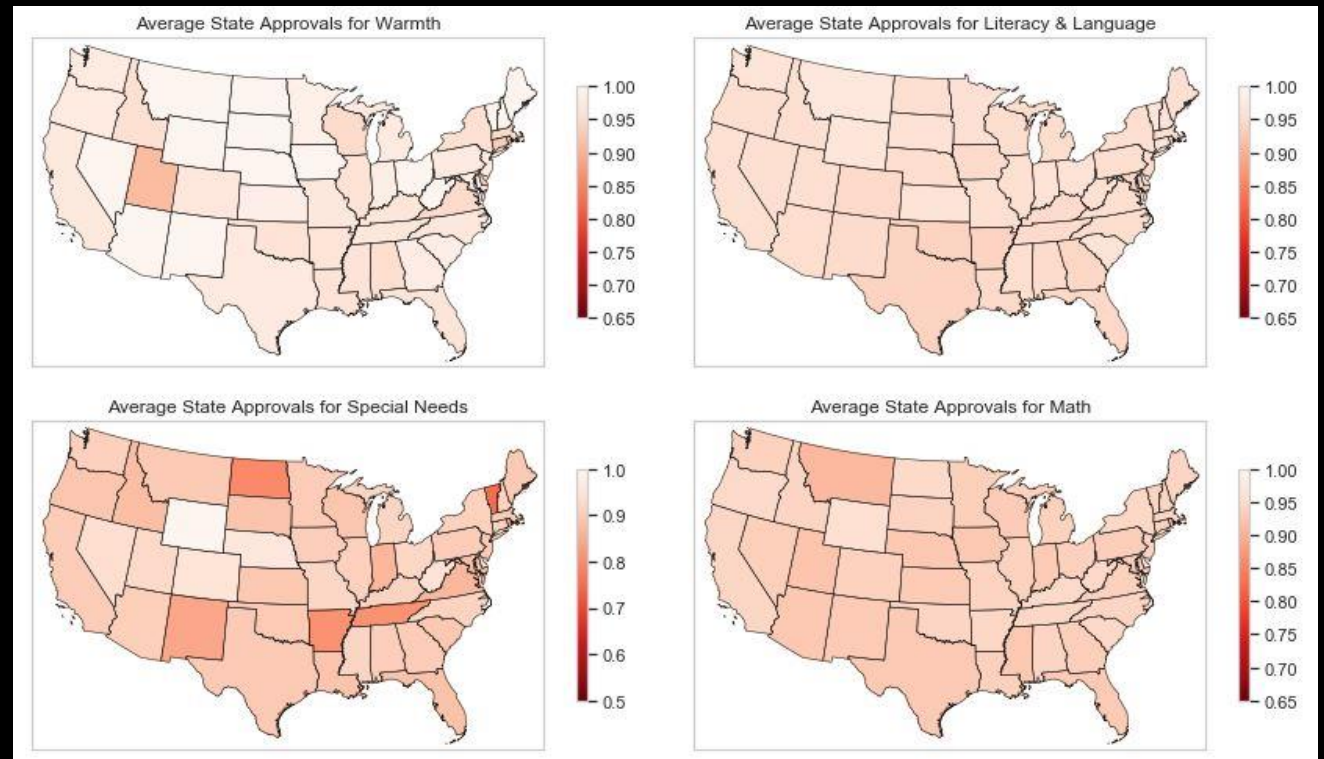
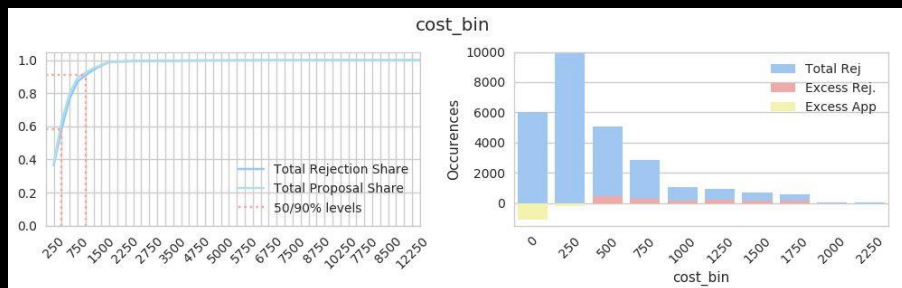
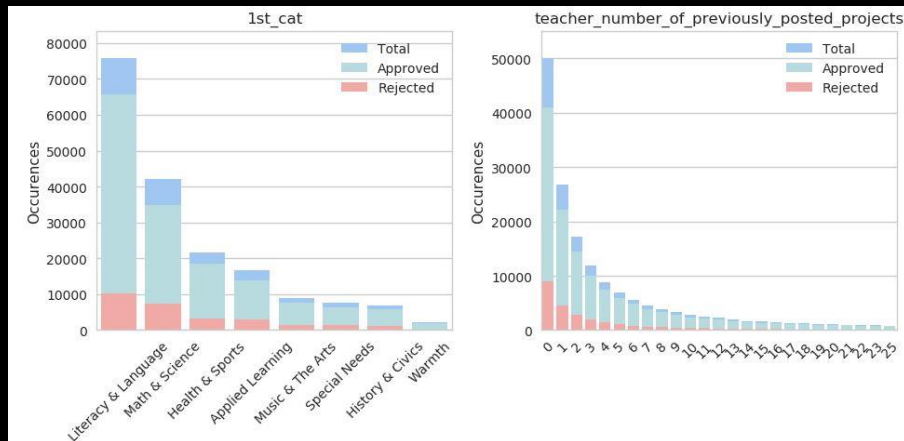
		Actual	
		0	1
Prediction	0	0	0
	1	27,909	158,151

		Actual	
		0	1
Prediction	0	0	0
	1	15%	85%

ANALYSIS STRATEGY:

- Clean missing values, Unicode entries, non-English punctuation
 - Basic visualizations. analyze baseline
 - Added Cost information, re-analyzed
 - Added simple text features (spelling, capitalization, etc.) re-analyzed
 - Added varying numbers of TF-IDF features, re-analyzed
-
- Multiple classifiers but limit to more “traditional” classifiers to test differences against winning solution. Used Logistic Regression, Gradient Boost, and Random Forest from sklearn. Added LGB which has a different api.

THE DATA: IT ALL LOOKS LIKE IT HAS THE SAME SHAPE, UNTIL IT DOESN'T...



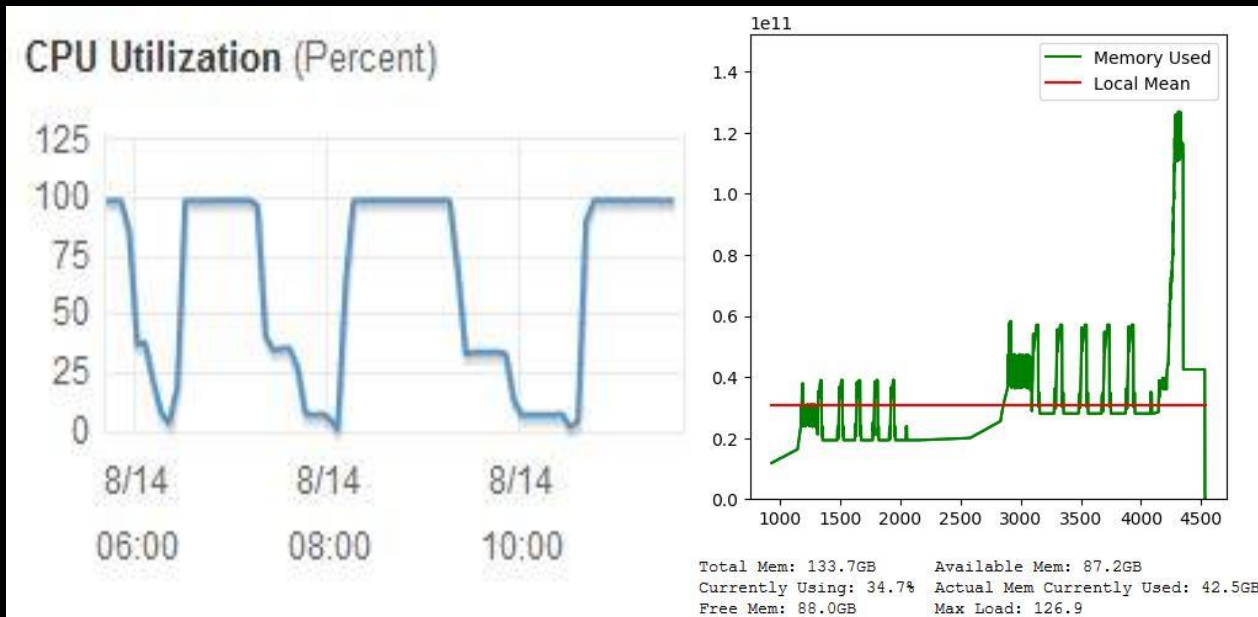
SOME ANALYSIS DETAILS

- Rounded most derived variables on the assumption that humans don't typically think in 64 bit math space.
- Limited textual analysis to ascii and English. There is some room for improvement here if there were more time.
- No feature reduction was included, again due to time constraints.
- Simple word scoring caused leakage, but offers room for improvement with some development. Classifier choice?

THE ARMS RACE IS OVER

AWS EC2 Instance: R54xl
(16 processors, 128gb memory)
Still some failures

My Laptop:
(8 processors, 16BG memory)



JUDGING SUCCESS

vs.

- Kaggle used ROC: Best Kaggle ROC Score: 0.76067
Place: 279/581
Submissions: 20
- winner 0.82812
48th percentile
51/top10, 161 #1

- Improvement in FP rate: 1.8%/15% → 12% improvement

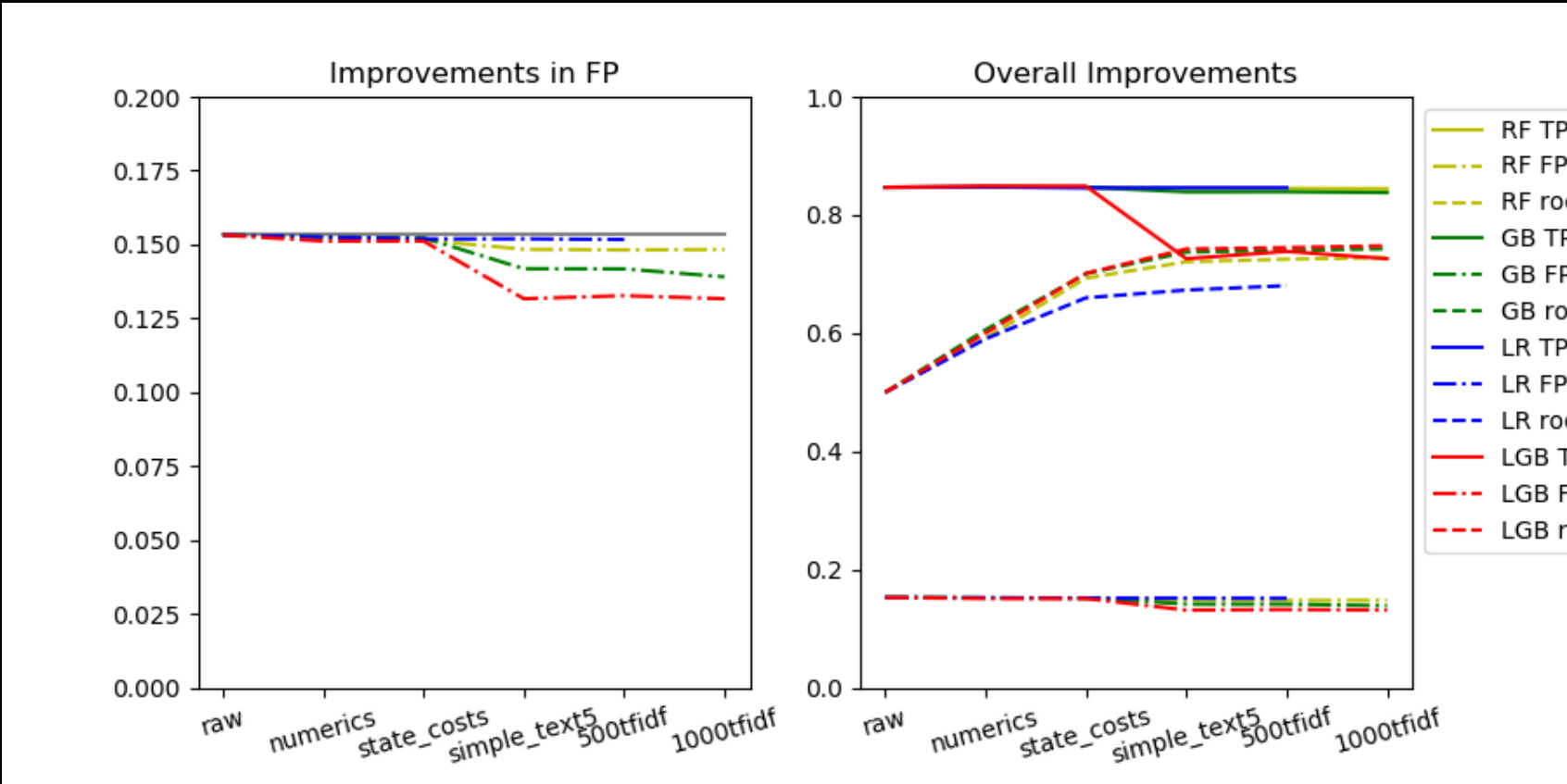
- Confusion Matrix: Before and After

		Actual	
		0	1
Prediction	0	0	0
	1	15%	85%

		Actual	
		0	1
Prediction	0	2.0%	12.3%
	1	13.2%	72.5%



A PERSPECTIVE ON RESULTS





CLIENT RECOMMENDATIONS

- Before using the models, change your training programs, with an increased focus on Math/Science and Special Needs. High leverage in bringing those rejection ratios to the mean.
- Develop guidelines around the simple text screens. That was the single biggest grouping for improvement in the FP rate.
- Make more geographic segmentation available. There is almost certainly district to district differences that outweigh the state to state variations.



IF I HAD IT TO DO OVER:

- Different decisions on how to structure text features. Specifically, I would try separating the resources and essays.
- I would use larger AWS instances (or spark.) Some of the classifiers did not complete with more than 1.5K features.
- I would focus on computationally (and memory) efficient classifiers.

SUMMARY:

- PROBLEM: Binary classification, 180K training datapoints, total features (post tfidf) were >5000. Required critical decisions at each step.
- Achieved a 12% improvement in the target parameter. Decent from a data science perspective, but pragmatically, this would not benefit the client proportionally. More work is needed.
- Light Gradient Boost was by far the most effective classifier that I tried due mainly to its efficiency.