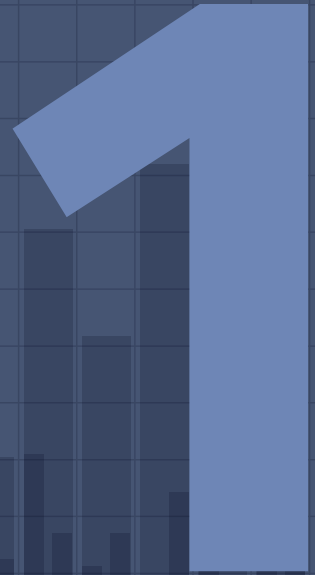


Using Feature Selection to Identify Gene Significance in Drug-Resistant Tuberculosis

Alan Qi and Powell Zhang

Overview



Overview

- Multidrug resistant tuberculosis is a form of tuberculosis that has developed resistance to *isoniazid* and *rifampin*.
- We worked with a binary dataset of resistant and susceptible tuberculosis strains.
- We utilized feature selection to determine significant genes with four models: CART, random forest, naive bayes, and genetic learning.
- We found *pncA* gene was most significant in determining drug resistance.
- These methods can be applied to other similar datasets in the future

Introduction

A large, light blue number '2' is positioned on the right side of the slide. The background is a dark blue grid with a faint, lighter blue bar chart pattern at the bottom.

Tuberculosis (TB)

- Caused by *Mycobacterium tuberculosis*
- Can be cured by administering antibiotics such as *Isoniazid, Rifampin, Pyrazinamide*, and others
- Can develop resistance to certain antibiotics, making it very difficult to cure
- Approximately 50% of patients with drug resistant tuberculosis eventually die due to the disease
- 240,000 global deaths from drug resistant tuberculosis in 2017
- 4.1% of new cases of TB and 19% of previously treated cases are drug resistant

Goals of project

- Use feature analysis to determine genes significant in determining drug resistance
- Use results to screen for drug resistance and prescribe effective antibiotics
- Expand to other diseases

Data

Biomarkers

	1a	2a	3a	4a	5a	6a	7a	8a	9a	10a	11a	12a	13a	14a	15a	16a
sample1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
sample2	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0
sample3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
sample4	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
sample5	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
sample6	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
sample7	0	0	1	0	0	0	0	0	0	0	0	0	0	0	1	0
sample8	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
sample9	0	0	1	0	0	0	0	0	0	0	0	0	0	0	1	0
sample10	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
sample11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
sample12	0	0	0	0	0	0	0	1	0	0	0	0	0	1	0	0
sample13	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
sample14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
sample15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
sample16	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0
sample17	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
sample18	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
sample19	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
sample20	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
sample21	1	0	0	0	0	0	0	0	0	0	0	0	1	0	1	0
sample22	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
sample23	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
sample24	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

480 strains of tuberculosis including drug resistant and susceptible strains
457 biomarkers showing mutated and wild type genes

Data

Each sample was labeled with an "R" or and "S" to show if it was resistant to drugs or susceptible

448a	449a	450a	451a	452a	453a	454a	455a	456a	457a	dst
0	0	0	0	0	0	0	0	1	0	S
0	0	1	0	0	0	0	0	0	1	R
0	0	0	0	1	0	0	0	0	0	R
0	0	0	0	0	0	0	0	0	0	R
0	0	0	0	0	0	0	0	0	0	S
0	0	0	0	0	0	0	0	0	0	R
0	0	0	0	1	0	0	0	0	0	R
0	0	0	0	0	0	0	0	0	0	R
0	0	0	0	1	0	0	0	0	0	R
0	0	0	0	0	0	0	0	0	0	R
0	0	0	0	0	0	0	0	0	0	R
0	0	0	1	0	0	0	0	0	0	R
0	0	0	0	0	0	0	0	0	0	R
0	0	1	0	0	0	0	1	0	0	S
0	1	0	1	1	0	0	0	0	0	R
1	0	0	0	0	0	0	0	0	0	S
0	0	0	0	0	0	0	0	0	0	S
0	0	1	0	1	0	0	0	0	0	S
0	1	0	1	1	0	0	0	0	1	S
0	0	0	0	0	0	0	1	0	0	S
0	0	0	0	0	0	0	0	0	0	R
0	0	0	0	0	0	0	0	0	0	R
0	0	0	1	0	1	0	0	0	0	S
0	0	0	0	0	0	0	0	0	0	S
0	0	0	0	0	0	0	0	0	0	S
0	0	0	0	0	0	0	0	0	0	S
0	0	0	0	0	0	0	0	0	0	S
0	0	0	1	0	0	0	0	0	0	S
0	0	0	0	1	0	0	0	1	0	S
0	0	0	0	1	1	0	0	0	0	S

Methods



3

CART (Classification and regression trees)

- Commonly known as decision trees
- Samples are divided based on certain characteristics using a binary tree
- Important genes are determined by looking at important leaf nodes

Random Forest

- Continuation of CART
- Creates many uncorrelated decision trees. Each tree decides on what the outcome should be, and the most commonly chosen outcome is returned



Naive Bayes

- Algorithm that estimates the probability of a strain being resistant/susceptible based on prior probabilities of predictors

The diagram shows the Naive Bayes formula with arrows pointing from labels to the corresponding parts of the equation:

- Likelihood** points to $P(x|c)$
- Class Prior Probability** points to $P(c)$
- Posterior Probability** points to $P(c|x)$
- Predictor Prior Probability** points to $P(x)$

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$

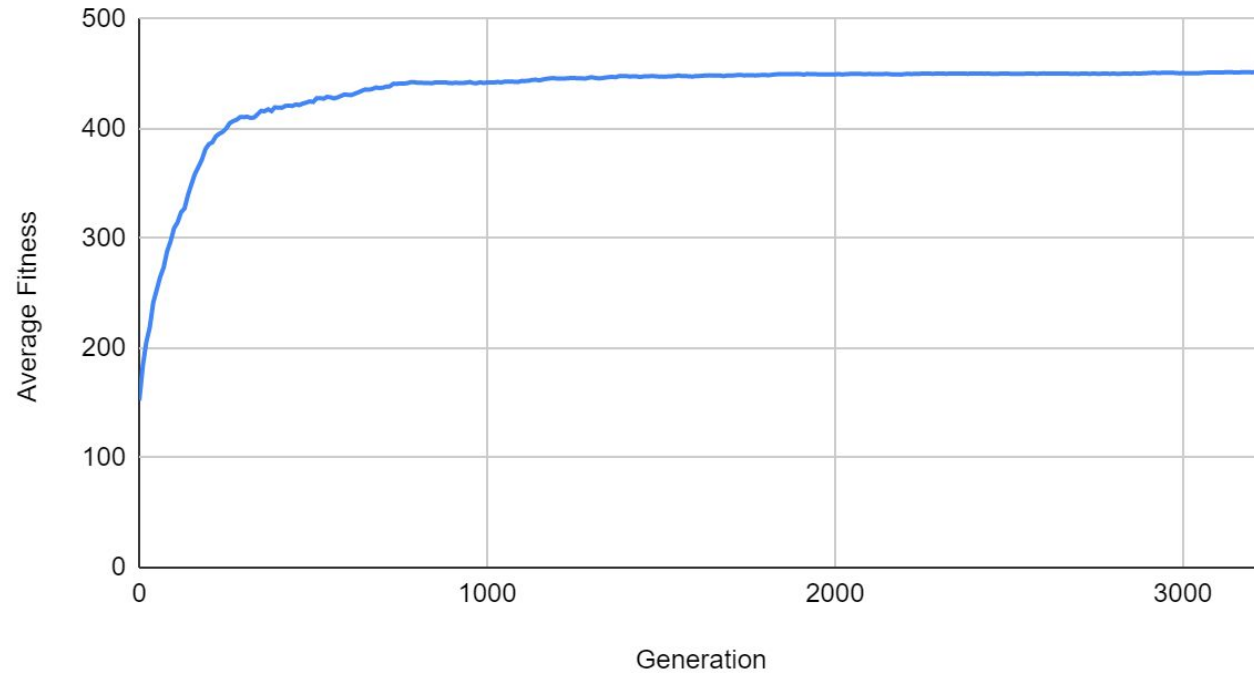
Below the diagram, the joint probability formula is given:

$$P(c|X) = P(x_1|c) \times P(x_2|c) \times \dots \times P(x_n|c) \times P(c)$$

Genetic Learning

- Randomly generate 100 lists of weights. Each weight corresponds to a gene.
- For each strain, weights of mutated genes are added up. If the total exceeds 10, the strain is categorized as resistant, otherwise it is categorized as susceptible.
- The lists are rated on how many strains they categorize correctly.
- The best lists are used to produce new lists with a combination of their weights.
- After many generations, the lists become better at categorizing strains. We choose the first list that can correctly categorize 95% of the strains and analyze its weights.
- Higher weights represent more important genes.

Average Fitness vs. Generation

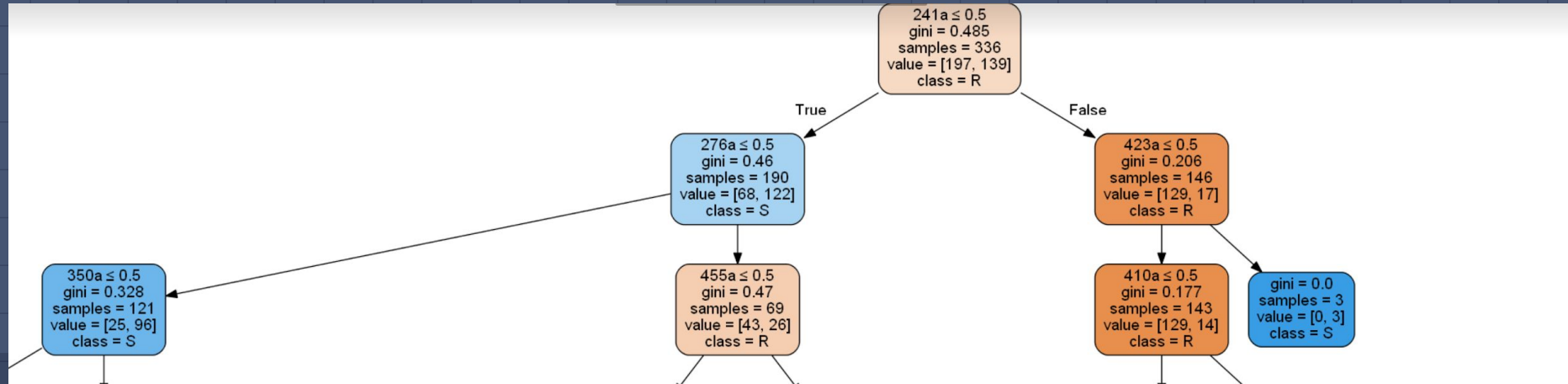


Results

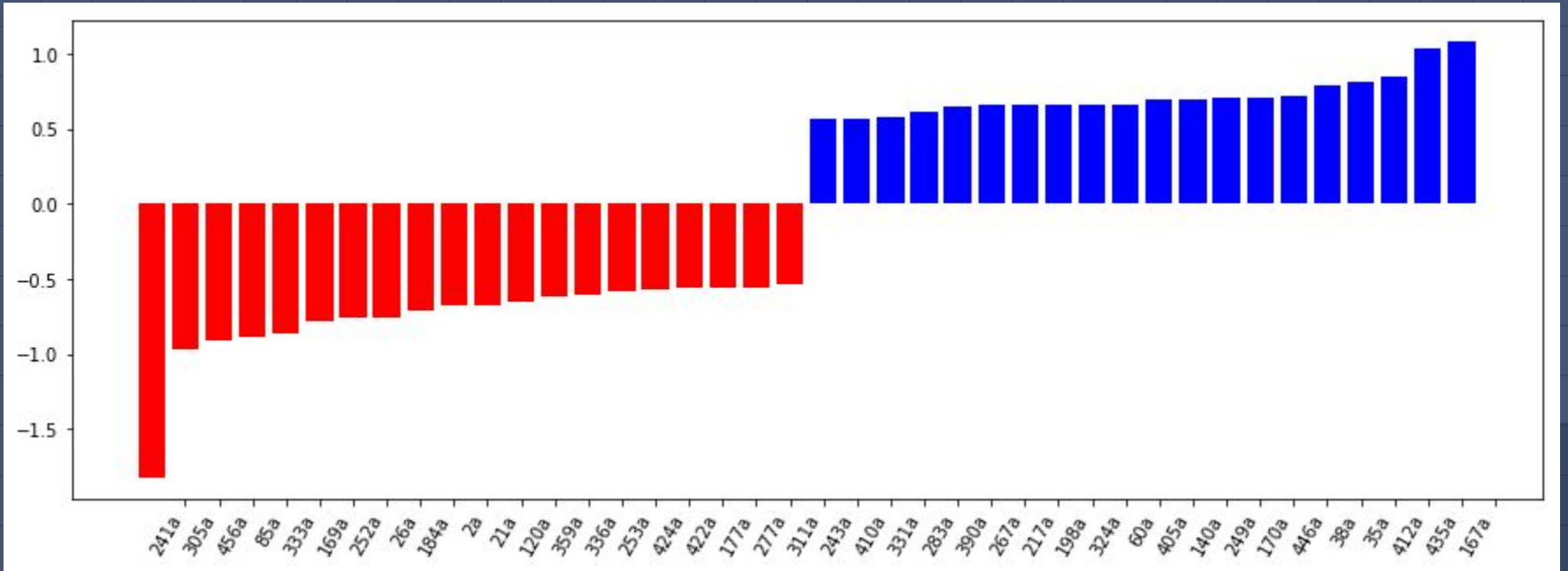


4+

Zoomed in on Decision Tree

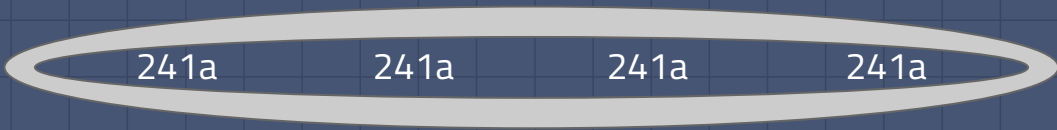


Naive Bayes Plot of Significant Features



Comparing Model Results side by side

Random Forest Naive Bayes Genetic Learning CART



241a	241a	241a	241a
276a	305a	120a	423a
455a	456a	208a	410a
409a	85a	396a	334a
3a	333a	289a	159a
277a	169a	2a	202a
346a	252a	455a	299a
272a	26a	184a	399a
219a	184a	218a	53a
275a	2a	305a	112a

Most Significant



Gene 241a is pncA

Functional pncA gene converts *Pyrazinamide*, a common tuberculosis drug, into its active form, pyrazinoic acid, which accumulates inside tuberculosis cells kills them.

Mutated pncA gene is strongly correlated with Pyrazinamide resistant *Mycobacterium tuberculosis*, indicated by previous studies.

Thus, feature selection is successful for finding gene of interest.

Conclusions and Future Study

Feature Selection is successful for this dataset

Useful in combating emerging crisis of antibiotic resistance

Applicable to other datasets of similar nature with further research

Acknowledgments

- Dr. Alterovitz, Insung Na
- Ling Teng
- MIT PRIMES
- Our Parents