

UC San Diego
SCHOOL OF MEDICINE

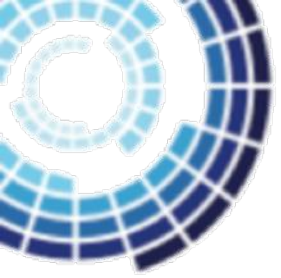
**Department of
BioMedical Informatics**

Justin Castro

Dr. Cinnamon Bloss & Dr. Jejo Koola

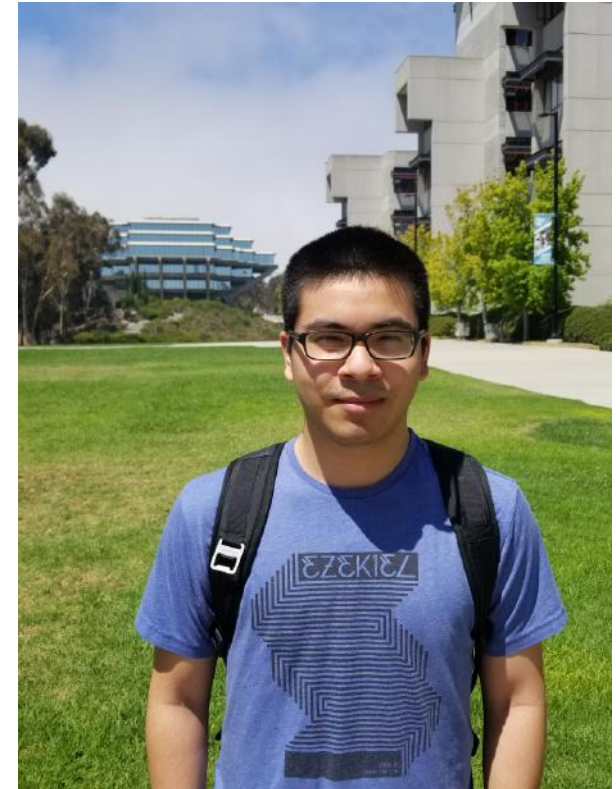


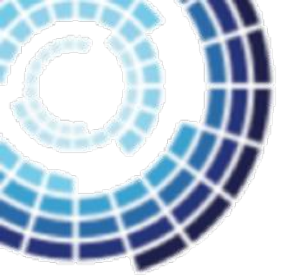
Applying Supervised Learning to Aid the Coding of Qualitative Data



About Me

- 3rd Year Data Science Major
- Worked alongside Dr. Bloss's lab.
- Applied to the DBMI program to learn more about Data Science in a clinical setting and machine learning.





Controlling Vector-Borne Disease with Genetic Engineering



<https://www.sandiegouniontribune.com/business/biotech/sd-me-darpa-ucsd-20170719-story.html>

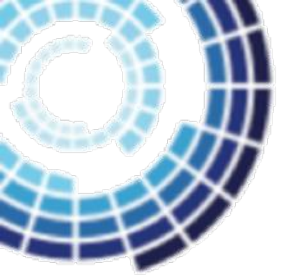
How can researchers incorporate concerns into the design of gene editing technologies?

What concerns do citizens have?

Does being more knowledgeable about GE technology change how people view it?



- Part of a larger collaboration to research solutions to vector borne disease.
 - GE Sterile Males
 - Gene Drive
- Focusing on the social and ethical perspective.



Analyzing Qualitative Data: Coding

Air Canada Review by Timmy Johnson

My spouse and I were flying to/from LHR via YUL; our outbound trip was at the end of July and home trip in early September. These seats are located immediately behind a service bulkhead and afford a lot of legroom.

The IFES and table tray are located in the armrests as is the remote control for the IFES and overhead lights etc. We found the trays easy to deploy but it took both arms and a lot of effort to get the IFES unit out of the armrest. When deployed it was difficult to view as it was somewhat offset by the armrest or the tray (if deployed).

The remote control is on wired tether which often tangled with other items and accidentally turned lights on or off or called cabin crew unnecessarily - annoying.

The seats are leather-faced and comfortable. An amenities kit is provided with pillow and blanket. The seats are close to washrooms. When cabin crew were accessing the bulkhead storage there was a tendency to hit my feet with the door.

Also this is a service area for crew so it can be noisy as can be the sounds of passengers walking to and from the washrooms. Overhead storage is plentiful. Meals were excellent as was the crew service. For overseas flights it is worthwhile paying extra for the comfort.

If we had to do this again we would probably choose seats farther back in this 24-seat cabin to avoid the bulkhead closet crew activity and washroom sounds.

Seating 22

Entertainment system 14

Entertainment system 14

Seating 22

Seating: Amenities kit 1

Seating 22

Seating 22

Crew 3

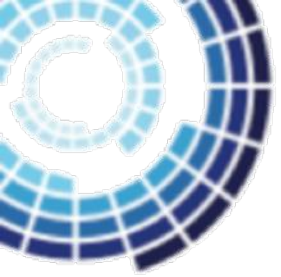
Crew 3

Meal service 6

Seating 22

Seating 22

- Coding facilitates more extensive analysis.
 - Indexing/Mapping/Tagging
 - Labeling your text allows you to revisit it in the analysis process
 - Makes the actual analysis much easier
- 30+ codes and 2300+ chat messages over 13 FG's.
- **Can we use ML to aid this tedious process?**



Preprocessing

- Qualitative data is complicated and messy

130: people on **nex** door **ocmplaining**
aobut being bit during day - i think so **cal**
peopel just **arent'** used to mosquitoes

711: i am 711 am i in
714: 711 has me craving a slurpee. ;)
711: thats funny

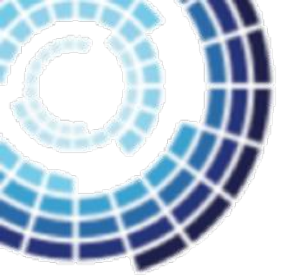
- Are all words equally important?

“724: I am not sure of the risks associated with **Gene Driving**. It might be a solution or it **might also create another problem.**”

“702: I'm for the most **cost effective and like that they are** not locally confined so I guess Sustained **Gene Drive**”

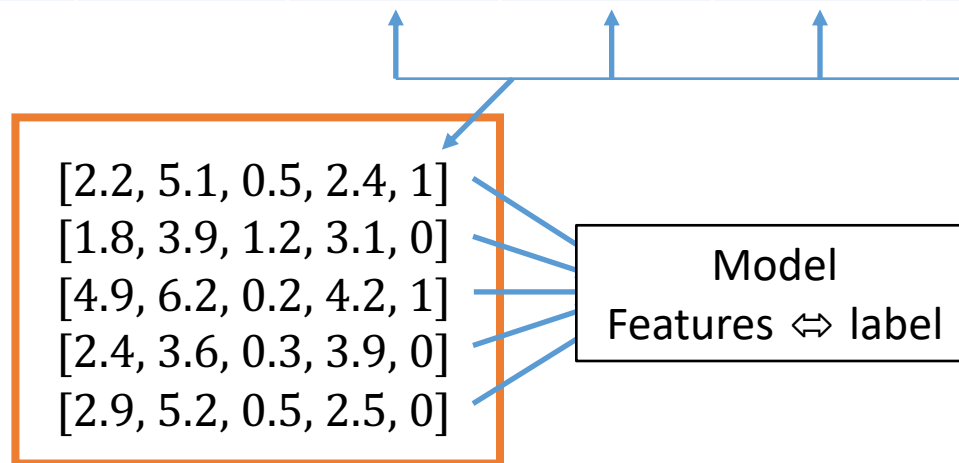
“711: would the **gene drive still** let them be a **part of the chain**”

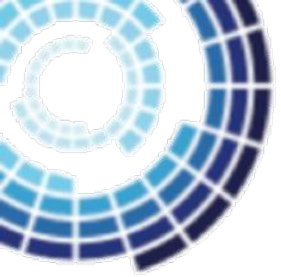
- Computers do not understand language



Feature Selection and Vectorization: Quantitative Data

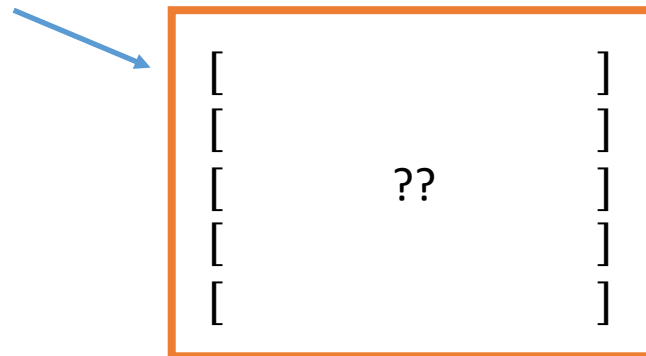
Flower type	Collector ID	Petal Width	Petal Length	Sepal Width	Sepal Length
sunflower	2	2.2	5.1	0.5	2.4
dandelion	2	1.8	3.9	1.2	3.1
sunflower	4	4.9	6.2	0.2	4.2
dandelion	4	2.4	3.6	0.3	3.9
dandelion	4	2.9	5.2	0.5	2.5

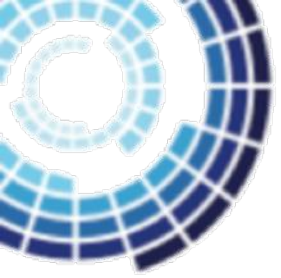




Feature Selection and Vectorization: Qualitative Data

Response	Cost	Diseases	...	Ecosystem
"724: I am not sure of the risks associated with Gene Driving. It might be a solution or it might also create another problem."	0	0	...	0
"702: I'm for the most cost effective and like that they are not locally confined so I guess Sustained Gene Drive"	1	0	...	0
...
"711: would the gene drive still let them be a part of the chain"	0	0	...	1





Feature Selection and Vectorization: Qualitative Data

a_1	a_2	a_3	...	a_n
1	4	1	...	2
0	3	1	...	1
0	2	0	...	0
1	2	0	...	1
1	2	0	...	0
0	3	0	...	1
0	2	0	...	0
...
0	1	0	...	0

- Bag of Words representation, Term Document Matrix, Document Vector Matrix.
- One way of representing documents mathematically.
 - Word Embeddings (Kevin)

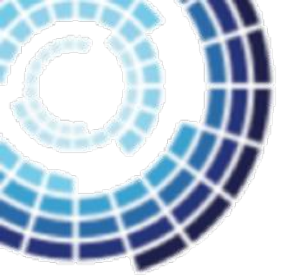
Columns = All unique words across ALL responses (the corpus)

Rows = Individual responses

Values = Times that word appears in the document

Responses

Unique Words



Feature Selection and Vectorization: Qualitative Data

a_1	a_2	a_3	...	a_n
.4251	.4947	.56924928
0	.4256	.56924297
0	.3985	0	...	0
.4251	.3985	04297
.4251	.3985	0	...	0
0	.4256	04297
0	.3985	0	...	0
...
0	.2718	0	...	0

- Term Frequency – Inverse Document Frequency (TF-IDF) Vectorization
- Gives more weight to words that are unique to that document.

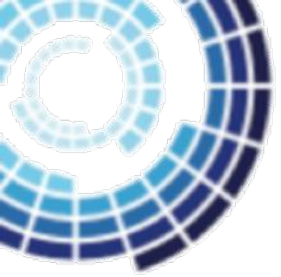
$$tf_{i,d} \times \log \frac{\text{Total \# of Documents}}{df_i}$$

tf = term frequency

df = document frequency

Responses

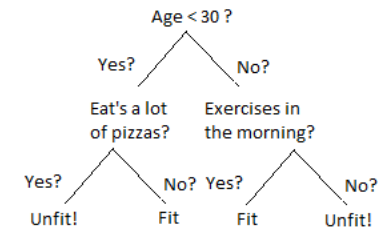
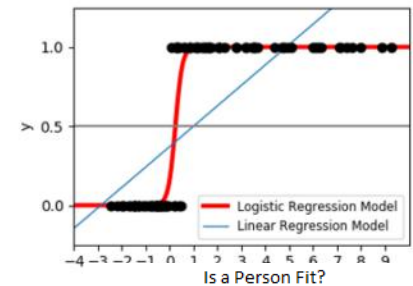
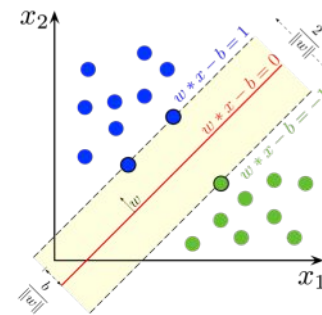
Unique Words

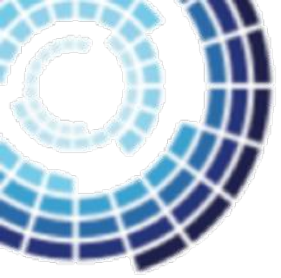


Machine Learning Models

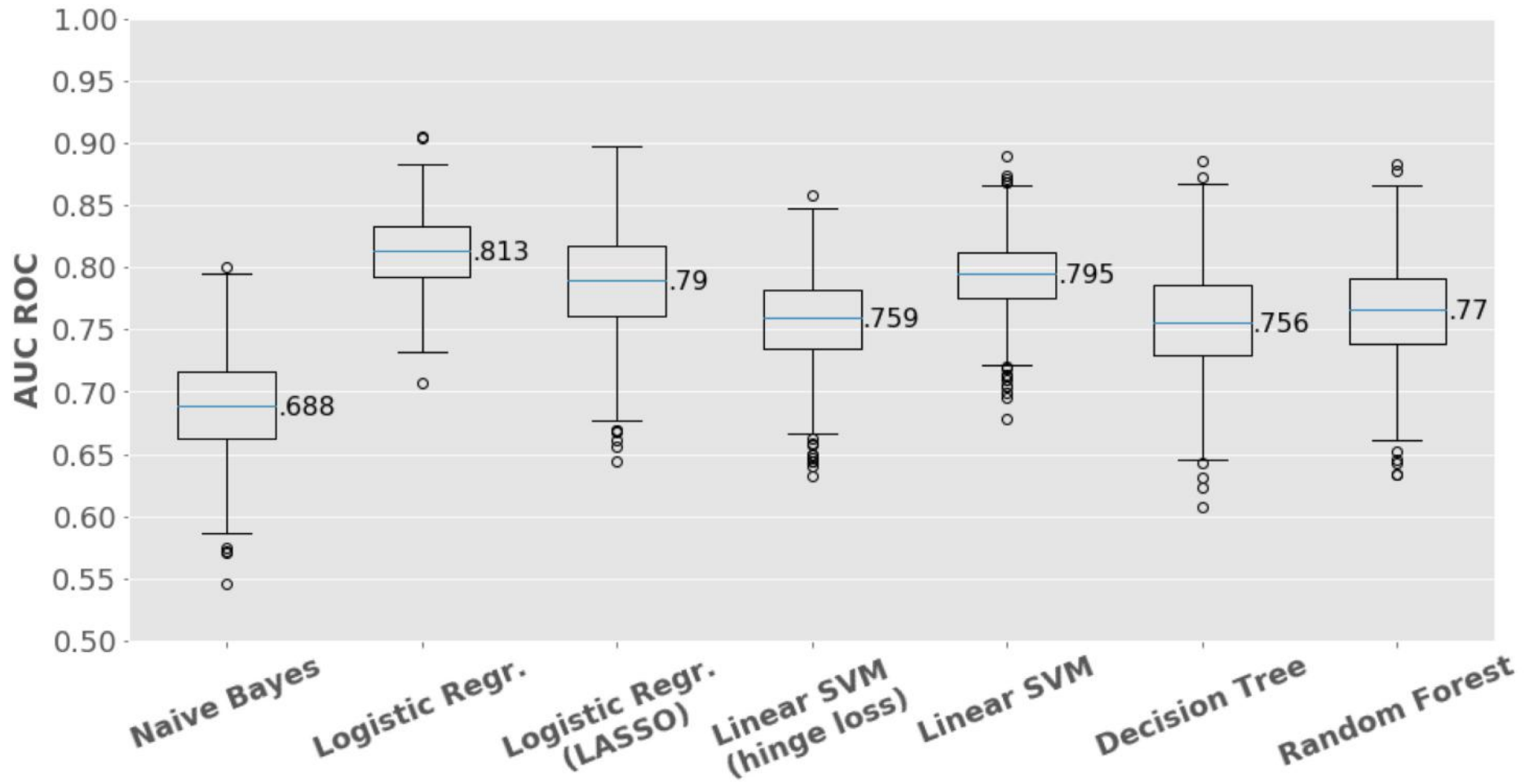
- Disease vs. Not Disease
 - Mentions of diseases or negative health consequences.
- Models
 - Naïve Bayes
 - Logistic Regression
 - Support Vector Machines
 - Decision Trees
- Internally Validated with bootstrapping.
- Externally Validated with FDA comment dataset.
- Area Under ROC Curve (AUC ROC) primary metric of model discrimination.

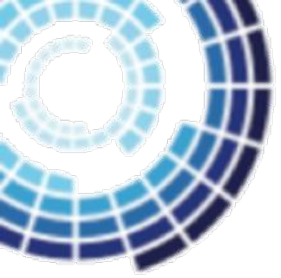
$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$





Comparing Model Performance





Logistic Regr. and Naïve Bayes: a closer look

Logistic Regr.

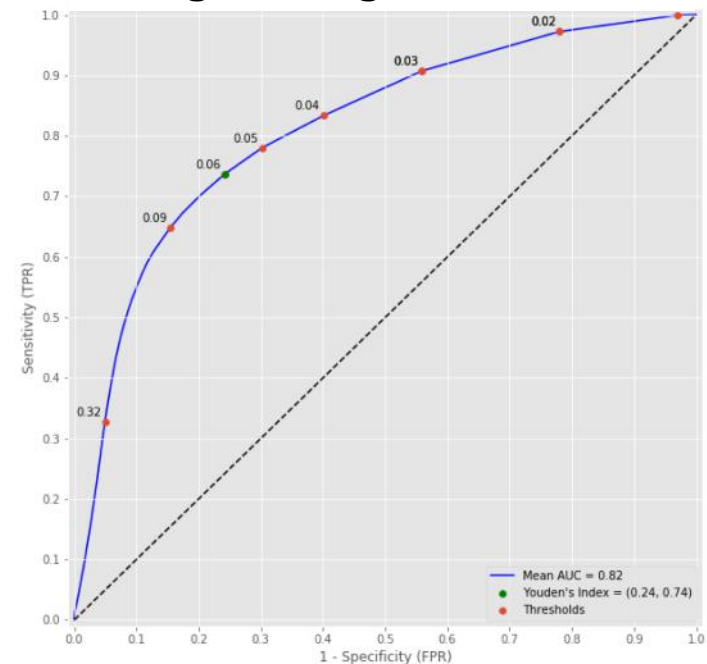
word	β
disease	4.69
health	2.18
carry	1.81
life	1.01
eliminate	1.01
zika	1.01
malaria	1.01
risk	.992
spread	.955

Naïve Bayes

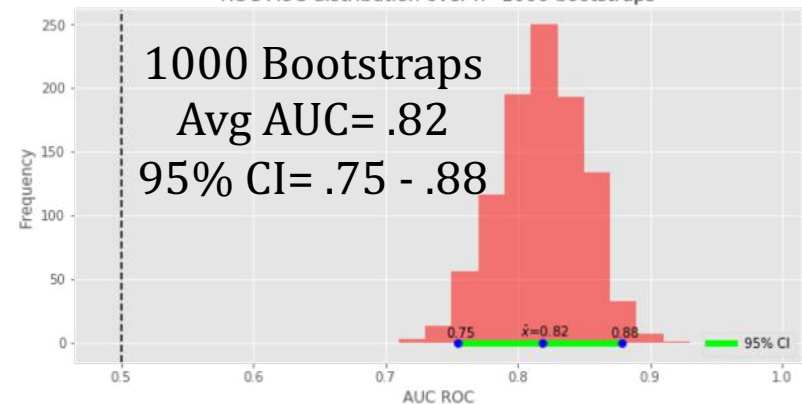
word	weight
disease	-5.07
health	-5.93
carry	-5.98
mosquito	-6.19
get	-6.24
control	-6.26
spread	-6.30
eliminate	-6.31
people	-6.33

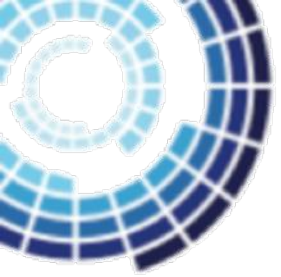
Higher = more indicative of "Disease"

Logistic Regr. ROC curve



ROC AUC distribution over n=1000 bootstraps





But is the Model Actually Generalizable? External Validation

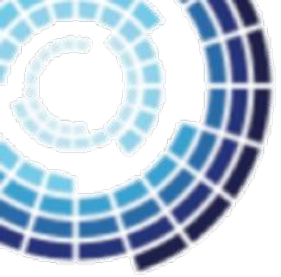
- FDA Comments Dataset
 - 2600 individual comments sent to the FDA in 2014.
 - Potential impact of Oxitec OX513A Mosquito deployment in Monroe County, Florida.
 - Would release modified male mosquitoes whose offspring would carry a gene that produced a lethal protein.

“Use the GMO mosquitos NOW, before there's thousands of cases of microcephaly here in Florida. Waiting is not an option.”

“PLEASE don't release these GMO mosquitoes. Please STOP this madness. Nature knows best not science.”

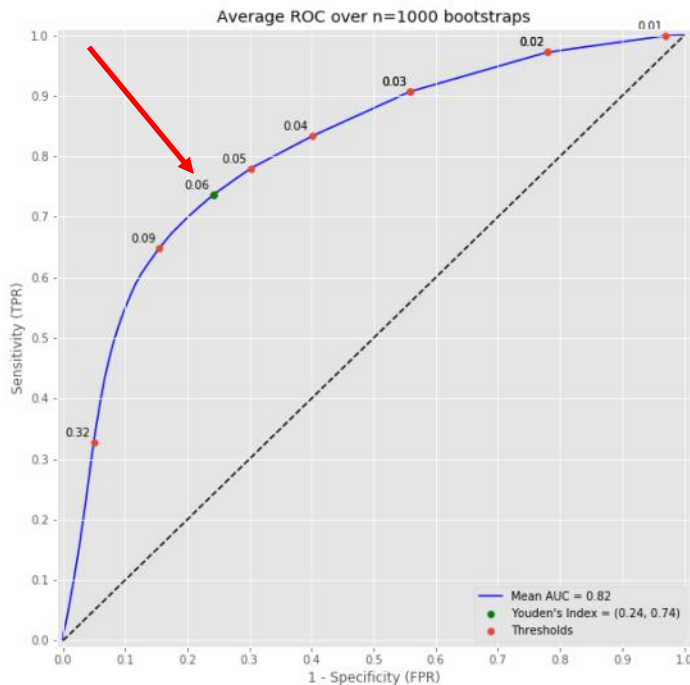


OXITEC



Generalizing a Model: FDA comments

TP = True Positive FP = False Positive
FN = False Negative TN = True Negative

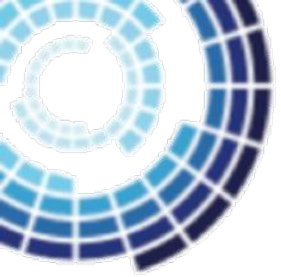


Expected:
FPR = .24
TPR = .74

Ground Truth

		Ground Truth	
		Disease	Not Disease
Predicted	D	62 TP	38 FP
	ND	27 FN	73 TN

67.5% Accuracy
F1 Score = .66
FPR = .34
TPR = .697



Takeaways and the Future

- Garbage in, Garbage out.
 - Preprocessing and Feature Selection is key.
- Is machine learning a reliable coding method?
 - Depends on aims of qualitative analysis.
 - Depends on how distinct codes are, distinguishing features are key.
- Limitations of BOW representations.
- New territory for ML aided coding.



UC San Diego

SCHOOL OF MEDICINE

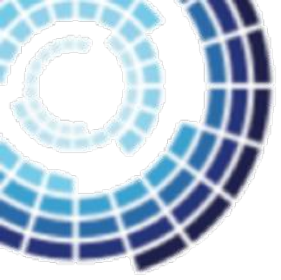
Department of BioMedical Informatics

Acknowledgements

- My Mentors:
 - Dr. Cinnamon Bloss
 - Dr. Jejo Koola
- Kevin Ngo
- Cynthia Schairer Ph.D
- The Bloss Lab



Questions?



Logistic Regr. Calibration with/without Isotonic Regression

