

# A comparison of statistical methods for improving rare event classification in medicine

Thesis Defense

Alyssa Forber  
April 18th 2018

Colorado School of Public Health  
University of Colorado, Anschutz Medical Campus  
Department of Biostatistics & Informatics

# Outline

**Problem**

**Objectives**

**Methods**

**Case Studies**

**Simulation Study**

**Discussion**

**Conclusion**

# Imbalanced Learning Problem

- When one outcome class greatly outnumbers the other, the outcome is said to be imbalanced
- Developing a classification model for imbalanced data results in over learning the majority class and under learning the minority class, and subsequently in low sensitivity and high specificity
- Need to improve predictive performance
- Many examples in medicine: rare diseases, rare adverse reactions to medications, etc.

# Aims

- Improve predictive performance for imbalanced dataset
- Utilizing measures of sensitivity, specificity, accuracy, and ROC analysis to evaluate performance
- Compare two methods to handle imbalance
  - ▶ Informed probability cutpoints for predicted probabilities
  - ▶ Sampling techniques to balance datasets
- Evaluate these methods in two clinical data sets and in a simulation study
- Note: we are not doing cost-sensitive learning

# Methods

- Split data into training and test sets using a 2/3:1/3 temporal split
- Create balanced datasets through sampling on training set
- Fit model to sampled data
- Get predicted probabilities on the hold-out test data
- Optimize probability cutoff for outcome

## Cross-validated lasso regression - Least Absolute Shrinkage and Selection Operator

- Lasso:
  - ▶ Shrinks estimates by penalizing size of coefficients
  - ▶ Performs variable selection by shrinking some to 0

$$\hat{\beta}_{lasso} = \underset{\beta}{\operatorname{argmin}} \sum_{i=1}^N (y_i - \beta_0 - \sum_{j=1}^p x_{ij}\beta_j)^2$$

subject to  $\sum_{j=1}^p |\beta_j| \leq t$  where  $t$  is the tuning parameter.

# Cross Validation

- Find the best “tuning parameter” for determining amount of shrinkage
- Split data into  $k$  parts and then train on all except one, which is held out for testing
- Then pick the tuning parameter that minimizes error

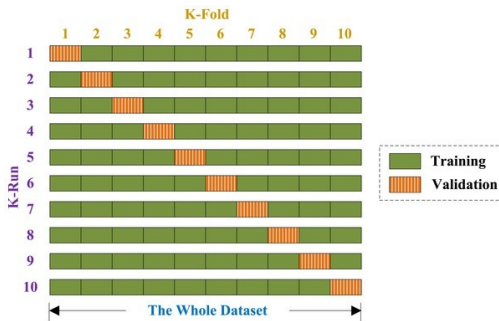


Figure 1: Source: Zhang Y, Wang S. (2015) Detection of Alzheimer's disease by displacement field and machine learning.

# Advantages and Disadvantages of Lasso

- Advantages:

- ▶ Lower variance of the predicted values
- ▶ More accurate predictions
- ▶ Reduces the number of predictors

- Disadvantages:

- ▶ Biased coefficients, inference not same as logistic regression
- ▶ No standard errors or p-values out of the model
- ▶ When two or more independent variables are highly correlated, lasso arbitrarily chooses only one to keep



# ROC & Cutoff Probabilities

ROC (with pROC package):

- Receiver Operating Characteristics
- ROC curve plots sensitivity vs 1- specificity
- Each point on curve corresponds to a decision cutoff for classification
- Youden's J statistic is used a measure to choose an optimal cutoff
  - ▶  $J = \text{sensitivity} + \text{specificity} - 1$

# ROC Curve

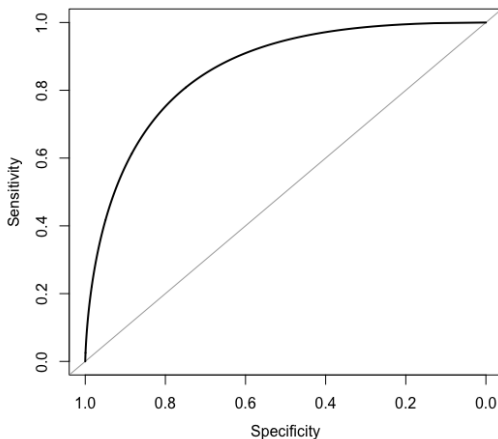


Figure 2: ROC Curve

# Confusion Matrix

**Correctly identify those w/ outcome:**

$$\text{Sensitivity} = \frac{TP}{TP + FN}$$

**Correctly identify those w/o outcome:**

$$\text{Specificity} = \frac{TN}{TN + FP}$$

**Correctly identify either group:**

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + TN + FN}$$

		Predicted class	
		P	N
Actual Class	P	True Positives (TP)	False Negatives (FN)
	N	False Positives (FP)	True Negatives (TN)

Picture Source: "Confusion Matrix." [rasbt.github.io/mlxtend/user\\_guide/evaluate/confusion\\_matrix/](https://rasbt.github.io/mlxtend/user_guide/evaluate/confusion_matrix/).

# First Approach

No Sampling, Optimize Cut-off:

- Fit original unsampled data to the lasso model
  - ▶ Predictions return probability between 0 and 1 for each observation of the test set
- Use 0.5 standard probability cutoff to compare
- Find “best” probability cutoff
  - ▶ Youden's Index

# Second Approach

## Sampling:

- Create sampled data sets that are balanced
  - ▶ Down sample
  - ▶ Up sample
  - ▶ SMOTE
- Predict and use Youden's Index as cutoff

# Down Sampling

- Under-sample majority to equal minority

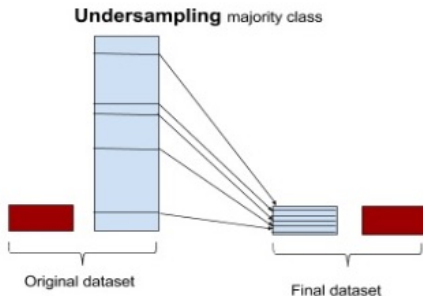


Figure 3: Down sample

Picture Source: "Learning from Imbalanced Classes." Silicon Valley Data Science, 25 Sept. 2017, [svds.com/learning-imbalanced-classes/](https://svds.com/learning-imbalanced-classes/).

# Up Sampling

- Over-sample minority to equal majority

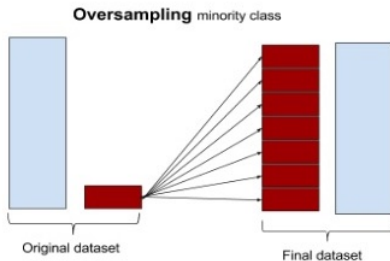


Figure 4: Up sample

Picture Source: "Learning from Imbalanced Classes." Silicon Valley Data Science, 25 Sept. 2017, [svds.com/learning-imbalanced-classes/](https://svds.com/learning-imbalanced-classes/).

# SMOTE

- Synthetic Minority Over-sampling Technique

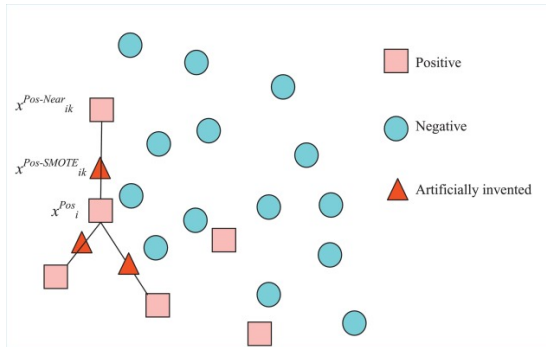


Figure 5: SMOTE

Picture Source: Sun, Jie, et al. "Imbalanced Enterprise Credit Evaluation with DTE-SBD: Decision Tree Ensemble Based on SMOTE and Bagging with Differentiated Sampling Rates."



# Models

- Model 1: No Sampling
  - ▶ 0.5 cutoff
  - ▶ Youden's Index
- Model 2: Down Sampling
  - ▶ Youden's Index
- Model 3: Up Sampling
  - ▶ Youden's Index
- Model 4: SMOTE
  - ▶ Youden's Index

# Case Studies

We used these methods on two case studies to see how results may differ between two real world examples with rare outcomes.

- Case Study 1: Predicting chronic opioid therapy in hospitalized patients (5%)
- Case Study 2: Identifying surgical site infections in hospitalized patients (3.4%)

# Case Study 1

- Data from Denver Health electronic health records (EHR) 2008 to 2014
- Definition of Chronic Opioid Therapy (COT) one year following the index hospital discharge:

*Receipt of  $\geq 90$ -day supply of opioids with  $< 30$ -day gap in supply over a 180-day period or receipt of  $\geq 10$  opioid prescriptions over one year.*

- 27,705 patients where 5% progressed to COT within a year
- 35 explanatory variables
  - ▶ Ex: age, race, history of chronic pain, discharge diagnosis

# Case Study 1

Variable	Yes COT 1,457 (5%)	No COT 26,248 (95%)	p-value
Past Year Benzodiazepine Receipt	16%	5%	<.0001
Past Year Receipt of Non-opioid analgesics	43%	16%	<.0001
Number of opioid prescriptions:			
0	38%	80%	
1	17%	11%	
2	14%	4%	
3	9%	2%	
4-9	23%	3%	<.001
Receipt of opioid at discharge	56%	28%	<.001
> 10 MME per hospital day	80%	52%	<.001

Some of these predictors are modifiable factors by clinicians, like number of prescriptions, or milligrams of morphine equivalents

# Case Study 1: Results

Table 1: Results for Chronic Opioid Therapy

Model	Threshold	Sensitivity	Specificity	NPV	PPV	Accuracy	AUC	Coefficients
Unsampled 0.5	0.5	8	99	96	35	96	86	31
Unsampled	0.043	85	73	99	12	73	86	31
Down sampled	0.401	85	73	99	12	74	86	34
Up sampled	0.399	85	74	99	12	74	87	34
SMOTE	0.472	74	84	99	17	84	86	33

# ROC Plot

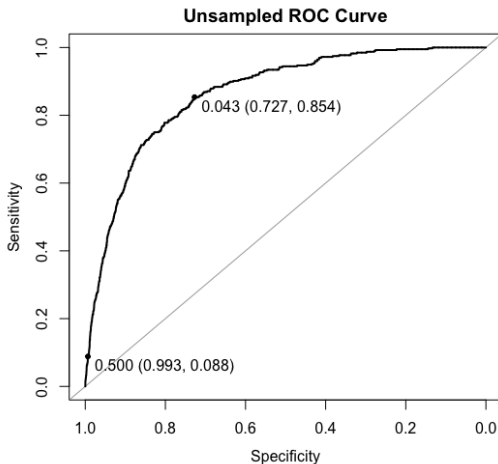


Figure 6: ROC for Original Data: Youden's Index and 0.5 cutoffs

## Case Study 2

- Need to identify post-operative complications to supplement manual chart reviews by nurses
- Surgical site infections (SSI) are the most common postoperative complication
- >50% of SSIs occur after patient discharge
- Data from 6,840 patients at the University of Colorado Hospital from 2013-2016 where 3.4% develop SSI
- 136 independent variables, mostly binary indicators
  - ▶ Ex: Antibiotic prescriptions, procedure codes, ICD-9 codes, risk of operation

## Case Study 2: Results

Table 2: Results for Surgical Site Infections

Model	Threshold	Sensitivity	Specificity	NPV	PPV	Accuracy	AUC	Coefficients
Unsampled 0.5	0.50	23	100	97	88	97	89	35
Unsampled	0.04	80	90	99	24	90	89	35
Down sampled	0.48	82	87	99	20	87	89	20
Up sampled	0.45	79	91	99	24	90	89	123
SMOTE	0.15	89	79	99	14	80	88	88



# Simulation Study

- We conducted a simulation study to look at a greater range of prevalences
- We chose 3%, 5%, 10%, 20%, 40%, and 50% prevalences
- Goals:
  - ▶ Evaluate performance of methods described previously
  - ▶ Determine at what point it's no longer necessary to worry about imbalance

# Simulation Methods

- Using the data from case study 1, we selected 10 of the strongest predictors to fit a logistic regression model
- Generate simulated COT outcome using coefficients from that logistic regression
- The new outcome was simulated with a logistic distribution

$$F(x) = \frac{e^x}{1 + e^x}$$

- Controlled prevalence by adjusting the intercept
- Included an additional 20 predictors (30 total) to implement sampling and lasso regression
- Average results for methods run 1000 times

# Simulation Results Part 1

Table 3: Results for 3, 5 and 10%

	Threshold	Sensitivity	Specificity	Accuracy	AUC	Coefficients
<b>3%</b>						
Unsampled 0.5	0.50	2	100	97	79	6
Unsampled	0.03	70	75	74	79	6
Down Sampled	0.49	70	74	74	78	9
Up Sampled	0.48	70	74	74	78	27
SMOTE	0.41	70	74	74	78	15
<b>5%</b>						
Unsampled 0.5	0.50	4	100	95	78	7
Unsampled	0.05	69	74	74	78	7
Down Sampled	0.49	69	74	74	78	9
Up Sampled	0.49	69	74	74	78	23
SMOTE	0.41	69	74	74	78	16
<b>10%</b>						
Unsampled 0.5	0.50	8	100	91	77	8
Unsampled	0.10	68	74	73	77	8
Down Sampled	0.49	68	74	73	77	9
Up Sampled	0.49	68	74	73	77	18
SMOTE	0.42	68	73	73	77	19

# Simulations Results Part 2

Table 4: Results for 20, 40, and 50%

	Threshold	Sensitivity	Specificity	Accuracy	AUC	Coefficients
<b>20%</b>						
Unsampled 0.5	0.50	21	97	82	76	8
Unsampled	0.20	66	73	72	76	8
Down Sampled	0.49	66	73	72	76	9
Up Sampled	0.49	66	73	72	76	13
SMOTE	0.42	66	72	71	75	23
<b>40%</b>						
Unsampled 0.5	0.50	49	84	70	74	9
Unsampled	0.39	66	71	69	74	9
Down Sampled	0.49	66	71	69	74	9
Up Sampled	0.49	65	71	69	74	10
SMOTE	0.41	64	71	68	73	28
<b>50%</b>						
Unsampled 0.5	0.50	63	72	68	73	9
Unsampled	0.49	63	72	68	73	9

# Sensitivity and Specificity by Prevalence

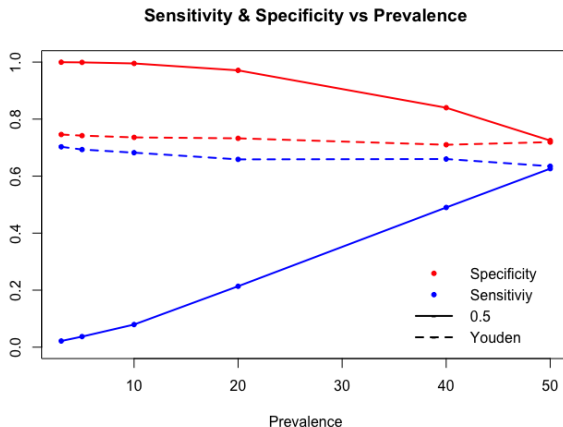


Figure 7: Comparing Youden's Index vs 0.5 cutoff results across prevalence

# Youden's Index by Prevalence

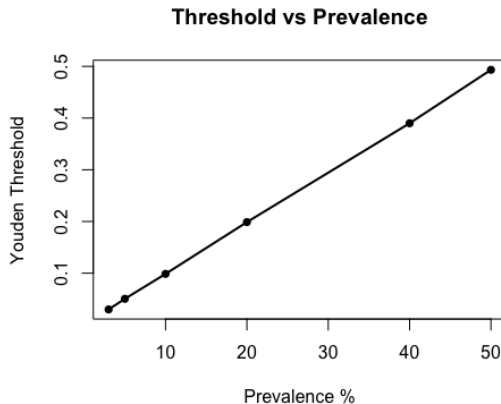


Figure 8: Youden's Index matches prevalence

# Discussion

- Similar performance between Youden's + sampling vs. Youden's alone
- One handles issue before and one after model fitting
- Over sampling and SMOTE had highest number of coefficients
- Threshold equals prevalence
- Sensitivity low even in less extreme imbalanced data for 0.5 cutoff

# Limitations

- Within methods:
  - ▶ Not cost sensitive learning
  - ▶ Lasso dropping highly correlated variable
  - ▶ Temporal split is considered external, but still could be applied to outside dataset
  - ▶ Choosing cutoff in test set may be biased
- Within dataset:
  - ▶ Patients acquiring prescriptions outside of Denver Health uncaptured
  - ▶ Nurses incorrectly reviewing patients



# Acknowledgments

Thesis Adviser Katie Colborn, thank you for all your time and mentoring  
Committee members Elizabeth Juarez-Colunga and Susan Calcaterra, thank  
you for your time and expertise  
Thank you for attending.

# Key References

- ① Chawla NV, Bowyer KW, Hall LO, et al. SMOTE: Synthetic Minority Over-sampling Technique. *Journal of Artificial Intelligence Research* 2002; 16: 321-357.
- ② Japkowicz NS, Shaju. The class imbalance problem: A systematic study. *Intelligent Data Analysis* 2002; 6: 429-449.
- ③ Kuhn, Max, and Kjell Johnson. *Applied Predictive Modeling*. Springer, 2016.
- ④ Tibshirani, Robert. "Regression Shrinkage and Selection via the Lasso." *Journal of the Royal Statistical Society. Series B (Methodological)*, vol. 58, no. 1, 1996, pp. 267-288. JSTOR, JSTOR, [www.jstor.org/stable/2346178](http://www.jstor.org/stable/2346178).
- ⑤ Youden WJ. Index for rating diagnostic tests. *Cancer* 1950; 3: 32-35.

# Supplemental Slides

# SMOTE

- 1 For each minority observation, calculate the  $k$  nearest neighbors (default 5)
- 2 Depending on amount of over sampling, pick nearest neighbor (for 200%, pick 2 of 5 at random)
- 3 For each attribute, calculate the difference between the original sample and it's selected neighbor
- 4 Compute a random number between 0 and 1 (the gap)
- 5 Create synthetic sample by adding to the original sample the difference times the gap

## Example:

- Age of sample is 40 and age of neighbor is 60
- Difference =  $60 - 40 = 20$
- Gap chosen at random to be 0.8
- Synthetic age is  $40 + 20 * 0.8 = 56$

# Case Study 1

- Design: Denver Health retrospective analysis electronic health record (EHR) data from 2008 to 2014.
- Patients: Hospitalized patients at an urban, safety-net hospital.
- Definition of Chronic Opioid Therapy (COT) one year following the index hospital discharge:

*Receipt of  $\geq 90$ -day supply of opioids with  $< 30$ -day gap in supply over a 180-day period or receipt of  $\geq 10$  opioid prescriptions over one year.*

# Case Study 1: Patient Population

- 27,705 patients
- Majority had incomes  $<185\%$  of the Federal Poverty Level
- 70% were ethnic minorities
- 5% with COT
- Excluded Patients:
  - ▶  $<15$  or  $>85$  years old
  - ▶ Those in prison, jail, or police custody
  - ▶ Those who died within one year following their index hospitalization
  - ▶ Patients with  $<2$  healthcare visits to Denver Health three years preceding their index hospitalization
  - ▶ Undocumented persons receiving emergent hemodialysis
  - ▶ Obstetric patients

## Case Study 2

One continuous measure- CPT-specific SSI event rate for initial operation (a measure indicating the risk of the operation, given past experiences of patients from the entire NSQIP dataset of more than 5.4 million patients)

CPT stands for common procedural terminology

ICD-9 stands for international classification of diseases version 9

# Logistic Distribution

Cumulative Distribution Function:

$$F(x) = \frac{e^x}{1 + e^x}$$

Probability Density Function:

$$f(x|\mu, \sigma) = \frac{e^{\frac{-(x-\mu)}{\sigma}}}{\sigma(1 + e^{\frac{-(x-\mu)}{\sigma}})^2}$$