## Thesis Proposal

### Predictive Modeling with Imbalanced Data

Alyssa Forber

University of Colorado, Anschutz Medical Campus

December 2017

# Outline

Problem

Objectives

Methods

Case Studies

Results

Simulation

Conclusion

# Imbalanced Learning Problem

- Presents a problem of imbalanced data
- Poor sensitivity with rare outcomes
- Need to improve predictive performance

# Aims

- Accurate predicting $\rightarrow$ improving sensitivity and specificity for imbalanced outcome
- Using and comparing methods of probability cut-points and sampling

# Methods

- Create sampled datasets
- Run model on sampled data
- Get predicted probabilities on the test data
- Optimize probability cutoff for outcome

# Model

- Roughly 2/3 temporal split of data to get train and test set
- Cross validated lasso regression

# Lasso

- Lasso:
  - Shrinks estimates
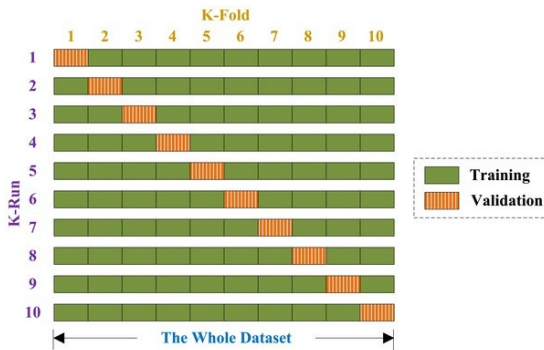  - Performs variable selection when shrunk to 0

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

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

# Cross Validation

- Cross validation:
  - Find the best "tuning measure" for model selection which determines amount of shrinkage of estimates
  - Split data into k parts and then train on each of those except one you validate against
  - Then pick the tuning measure that minimizes error

# Advantages and Disadvantages

- 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

# ROC & Cutoff Probabilities

ROC (with pROC package):

- ROC curve plots sensitivity vs specificity for each cut-off
- Top left corner is ideal
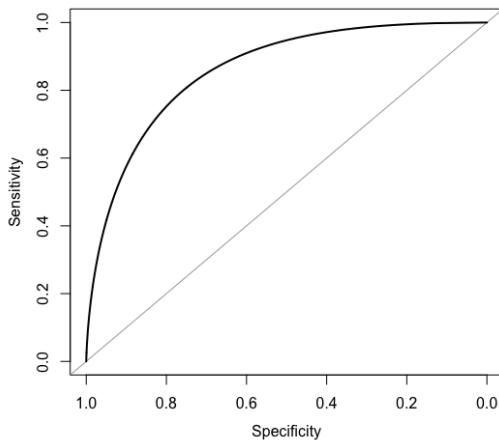- Youden Index is the furthest upper left corner or "max"

# ROC Curve



Figure 2: ROC Curve

## Confusion Matrix

**Correctly identify those w/ outcome:**

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

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

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

**Correctly identify either group:**

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

|  | Predicted class | |
| --- | --- | --- |
|  | $P$ | $N$ |
| $P$ | True Positives (TP) | False Negatives (FN) |
| $N$ | False Positives (FP) | True Negatives (TN) |

Actual Class

# First approach

No Sampling, Optimize Cut-off:

- Use original unsampled data and get predictions from the lasso model
  - Predictions return probability between 0 and 1 for each observation
- Use 0.5 standard probability cutoff to compare
- Find "best" probability cutoff
  - Youden Index

# Second Approach

Sampling:

- Create sampled data sets that are balanced
    - Down sample
    - Up sample
    - SMOTE

- Predict and use both standard 0.5 and Youden Index as cutoff

# Down Sampling

- Under-sample majority to equal minority
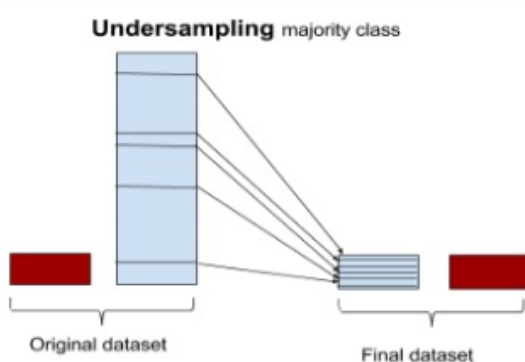


**Undersampling** majority class

Original dataset

Final dataset

Figure 3: Down sample

# Up Sampling

- Over-sample minority to equal majority
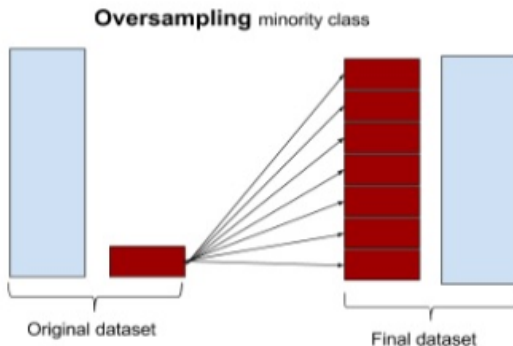


**Oversampling** minority class

Figure 4: Up sample

# SMOTE
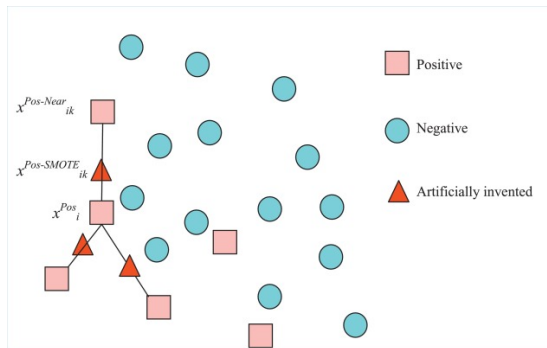
- Synthetic Minority Over-sampling Technique



Figure 5: SMOTE

## 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 1: Results

Table 1: Results

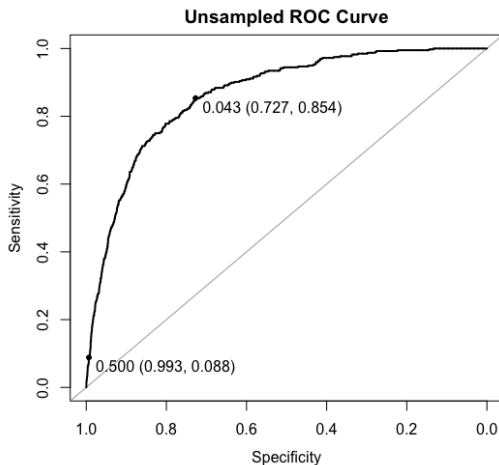| Data | Threshold | Specificity | Sensitivity | NPV | PPV | Accuracy | AUC | Covariates |
|------|-----------|-------------|-------------|-----|-----|----------|-----|------------|
| **Unsampled 0.5** | 0.5 | 99 | 8 | 96 | 35 | 96 | 86 | 31 |
| **Unsampled** | 0.043 | 73 | 85 | 99 | 12 | 73 | 86 | 31 |
| **Down sampled 0.5** | 0.5 | 81 | 75 | 99 | 15 | 81 | 86 | 34 |
| **Down sampled** | 0.401 | 73 | 85 | 99 | 12 | 74 | 86 | 34 |
| **Up sampled 0.5** | 0.5 | 82 | 75 | 99 | 15 | 82 | 87 | 34 |
| **Up sampled** | 0.399 | 74 | 85 | 99 | 12 | 74 | 87 | 34 |
| **SMOTE 0.5** | 0.5 | 86 | 71 | 99 | 17 | 85 | 86 | 33 |
| **SMOTE** | 0.472 | 84 | 74 | 99 | 17 | 84 | 86 | 33 |

# ROC Plot



Figure 6: ROC for Original Data: Younden and 0.5 cutoffs

# Case Study 2

# Case Study 2: Population?

# Case STudy 2: Results

# Simulation

# Simulation Results

# Conclusions

# Questions?

Questions