

PREDICTION OF SUCCESSFUL DISCHARGE

CHILDREN AND YOUTH IN BEHAVIORAL HEALTH TREATMENT: APPLYING
MACHINE LEARNING TO TREATMENT OUTCOMES

JUDITH CALVO, PhD

JUNE 25, 2020



SUCCESSFUL DISCHARGE

- HYPOTHESIZED PREDICTORS: AGE, DIAGNOSIS AT ADMISSION, LENGTH OF STAY, INITIAL SCORES ON STANDARD INSTRUMENT, ETHNICITY, LANGUAGE AND LEVEL OF CARE
- OUTCOME OF INTEREST: CHILD SUCCESSFULLY COMPLETES TREATMENT PROGRAM



WHY PROJECT WAS CHOSEN

- EFFECTIVE TREATMENT DEPENDS ON ABILITY TO CREATE INTERVENTIONS THAT ARE LIKELY TO LEAD TO SUCCESSFUL DISCHARGE
- IMPORTANCE OF ASCERTAINING WHAT LEADS TO SUCCESS IN NATURALISTIC SETTINGS WHERE RANDOMIZED CLINICAL TRIALS ARE NOT FEASIBLE. MACHINE LEARNING ALLOWS FOR USE OF DATA COLLECTED IN SUCH NATURALISTIC SETTINGS.
- MACHINE LEARNING MODELS ARE VERSATILE AND HIGHLY RECOMMENDED AT MANY LEVELS OF HEALTH CARE DELIVERY: DIAGNOSIS, PREDICTION AND TREATMENT PLANNING
- IMPORTANT TO IDENTIFY HIGH RISK CLIENTS AT ONSET AND INTERVENE EARLY AND AT THE APPROPRIATE LEVEL OF CARE TO AFFECT OUTCOMES



DATA SOURCE

- AGENCY EMR
 - LEGACY SYSTEM WITH LIMITS ON HOW DATA CAN BE EXTRACTED
 - USE OF CRYSTAL REPORTS AND SQL TO EXTRACT DATA
 - CHILDREN AND YOUTH BETWEEN 6-25 YEARS OLD
 - EMR CONTAINS DE-IDENTIFIED CLIENT DATA:
 - DEMOGRAPHICS
 - EPISODIC DATA (ADMISSIONS AND DISCHARGES)
 - DIAGNOSIS AT ADMISSION (PRIMARY DIAGNOSIS ONLY)
 - PROGRAM AS A PROXY FOR LEVEL OF CARE
 - DISCHARGE REASON-OUTCOME VARIABLE DICHOTOMIZED TO REFLECT SUCCESSFUL VS, NOT SUCCESSFUL REASON
 - SCORES ON A MEASURE OF NEEDS AND STRENGTHS AT INTAKE (CHILD AND ADOLESCENT NEEDS AND STRENGTHS--CANS)



QUESTIONS WE HOPE TO ANSWER

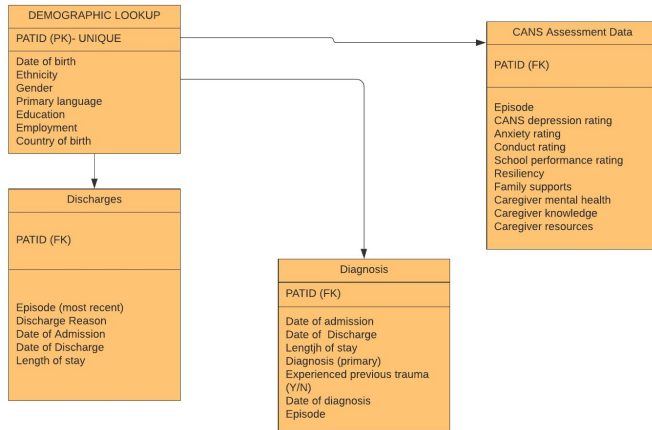
- WHAT FACTORS LEAD TO SUCCESSFUL COMPLETION OF TREATMENT?
- WHAT CLIENT CHARACTERISTICS ARE INDICATIVE OF HIGH RISK AND FAILURE TO COMPLETE TREATMENT?
- CAN WE CREATE A STREAMLINED, PARSIMONIOUS MODEL THAT CAN BE REPLICATED TO PREDICT TREATMENT OUTCOMES ACROSS OTHER DATASETS, TIMES AND GROUPS OF CLIENTS?



FIRST STEP: ERD DATA RELATIONSHIPS IN POSTGRES

ERD of My Database: Machine

judith.calvo | Learning Tables



- Final data did not include Country of Birth or Employment as these did not vary much. Most youth were born in the U.S. and were not employed at the time of admission
- Caregiver scores on the CANS were also removed in the final model as they did not show significant relationship to the outcome. They did not predict nor play a role in classification based on testing several models



Tools Used

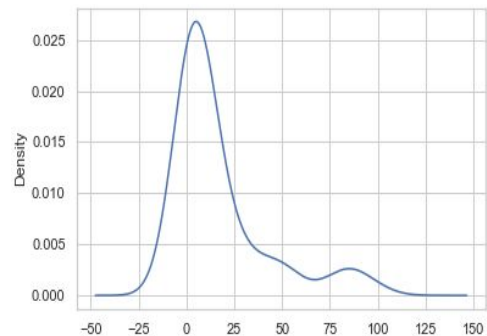
- DATA PREPROCESSING IN PANDAS
 - ENCODING
 - DUMMY VARIABLE CREATION
 - MATPLOTLIB AND SEABORN FOR VISUALIZATIONS
 - SMOTE TO CORRECT IMBALANCE AS NEEDED
- STATSMODEL TO CREATE FULL LOGISTIC MODEL WITH BETAS AND STATISTICAL SIGNIFICANCE
- CALCULATION OF ODDS AND PROBABILITY FOR KEY FEATURES USING NUMPY EXPONENTIATION FORMULAS
- LOGISTIC MODEL FITTING USING SCIKIT LOGISTIC REGRESSION LIBRARY
- ROC CURVE USING SCIKIT AUC LIBRARY
- RANDOM FOREST TO COMPARE RESULTS
- PCA AND HVPLLOT FOR INTERACTIVE GRAPHING
- POWER BI FOR INTERACTIVE DASHBOARD USING PYTHON SCRIPTING



DATA EXPLORATION

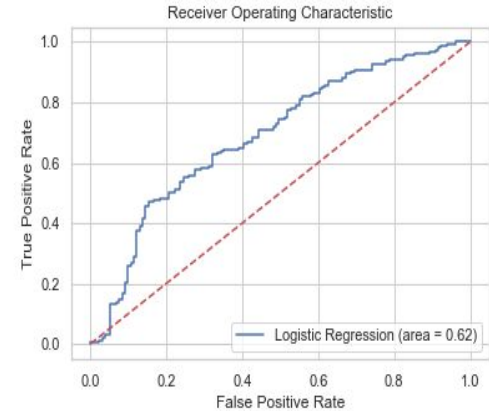
- DESCRIBE WITH SUMMARY STATISTICS
- LIST CATEGORICAL VARIABLES TO NARROW DOWN
- ENCODING OF CATEGORICAL VARIABLES
 - SCIKIT ENCODER
 - DENSITY PLOTTING TO ASCERTAIN BINNING CATEGORIES
 - CUSTOM ENCODING OF LEVEL OF CARE BASED ON YOUTH'S PROGRAM
 - CREATED DUMMY VARIABLES USING PANDAS GET DUMMIES FUNCTION
- GRAPHS
 - VISUALS OF DATA CHARACTERISTICS
- SMOTE TO CORRECT IMBALANCE

Sample Density Plot



DATA ANALYSIS

- LOGISTIC REGRESSION WITH FEATURE OUTPUT OF BETA WEIGHTS
- CALCULATION OF ODDS AND PROBABILITY FOR KEY FEATURES
- LOGISTIC MODEL FITTING
- METRICS TO EVALUATE MODEL
- ROC CURVE AFTER MODEL FIT (EXAMPLE)
- RANDOM FOREST TO COMPARE RESULTS
- PCA IN UNSUPERVISED LEARNING TO CREATE YOUTH GROUPINGS/PROFILES BASED ON ENTERING CHARACTERISTICS
- DROP AND ADD NEW FEATURES WITH NEXT RUN TO IMPROVE MODEL FIT AND PREDICTION, REPEAT ANALYSIS



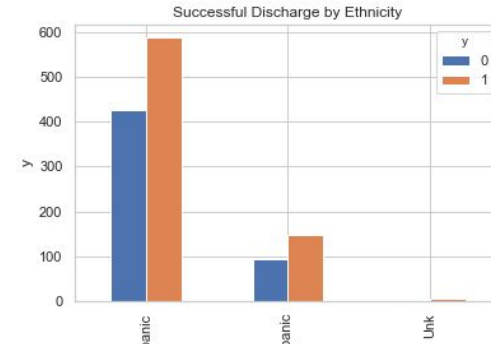
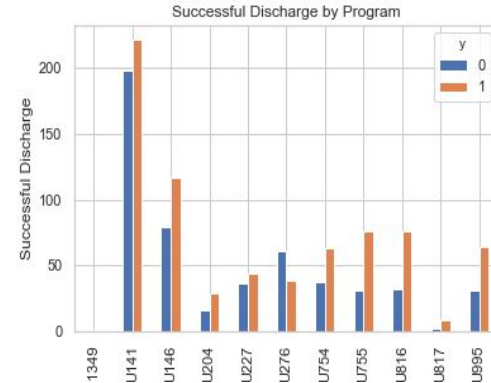
DASHBOARD TOOLS AND INTERACTIVE ELEMENTS

- POWER BI FOR INTERACTIVE DASHBOARD USING EMBEDDED PYTHON SCRIPTING
- INTERACTIVE ELEMENTS
 - SLICERS FOR EXPLORATORY DATA BY DISCHARGE STATUS, LEVEL OF CARE AND LEVELS ON INITIAL CANS SCORES FOR KEY FEATURES (I.E., LEVELS OF OPPOSITIONAL BEHAVIOR AT INTAKE, CHILD RESILIENCE, FAMILY STRENGTH) THAT WERE POSITED TO BE PREDICTIVE OF OUTCOME
 - PYTHON SCRIPT EMBEDDED IN POWER BI USING MATPLOTLIB FOR GRAPHING
 - IMPORTING OF DATASETS USING A PYTHON SCRIPT WITH PANDAS TO ALLOW FOR INTERACTIVITY OF VISUALIZATIONS.



STORYBOARD--DATA EXPLORATION PHASE

- DATA STORED IN POSTGRES AND .CSV
- CLEANED AND ENCODED
- SUCCESSFUL DISCHARGE BY ORIGINAL PROGRAM (1= SUCCESSFUL) PRIOR TO CUSTOM ENCODING TO CREATE LEVELS OF CARE
- SUCCESSFUL DISCHARGE BY ETHNICITY (NO STRONG DIFFERENCE IN PROPORTIONS BASED ON ETHNICITY--STILL HIGHER PROPORTIONS OF SUCCESS)



STORYBOARD--ANALYSIS

- LOGISTIC REGRESSION RESULTS
- LOGISTIC MODEL FIT
 - CLASSIFICATION REPORTING
- RANDOM FOREST
 - CLASSIFICATION AND FEATURE IMPORTANCE
- PCA CLASSIFICATION (POSSIBLY EXCLUDE)



PROPOSED DASHBOARD LAYOUT

Dashboard/Story Board Layout

Judith Calvo | June 12, 2020

Logistic Regression (second tab will be Random Forest)

