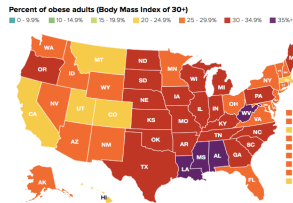


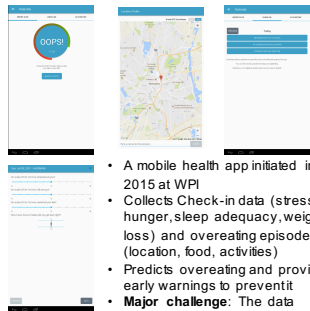
Background



Adult Obesity in the United States (2015) [1]

- 38 percent of American adults are obese
- Claimed between \$247 to \$210 billion dollars in health care spending

SlipBuddy



- A mobile health app initiated in 2015 at WPI
- Collects Check-in data (stress, hunger, sleep adequacy, weight loss) and overeating episodes (location, food, activities)
- Predicts overeating and provides early warnings to prevent it
- **Major challenge:** The data collected in the user study is limited in size

Methodology

- Generate synthetic data
- Construct predictive model with different machine learning techniques using k-fold cross validation
- Evaluate and compare the predictive performance of different machine learning models using multiple performance metrics
- Platforms: IPython Notebook and R Studio



Synthetic Data Generation

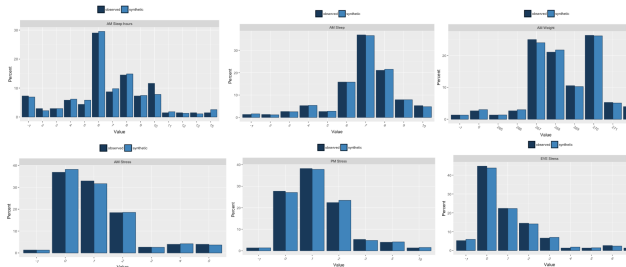
[illegible]

Using synthpop packages on R [2]:

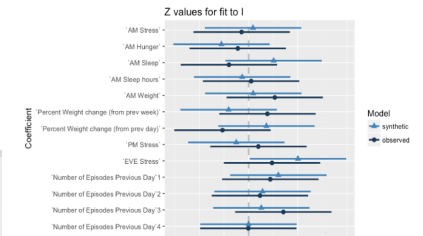
- Fit the observed data into a sequence of regressions
- First variable is randomly sampled with replacement from real data
- The values of each subsequent variable are drawn from a posterior distribution yielded by the cart regression of the previous variables

Two key properties of the original data that must be maintained:

- The probability distribution of data within each variable
- The correlation to Overeating Episode of each variables

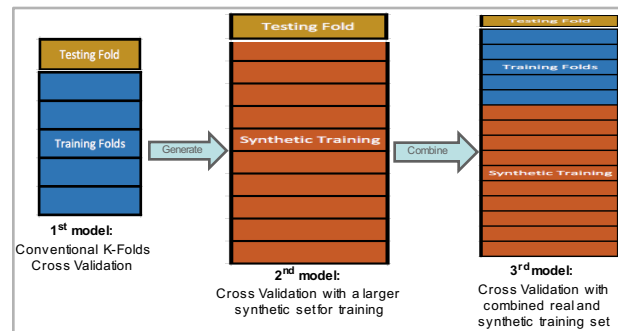


Histograms showing the similarity of relative frequency distribution between real data and synthetic data

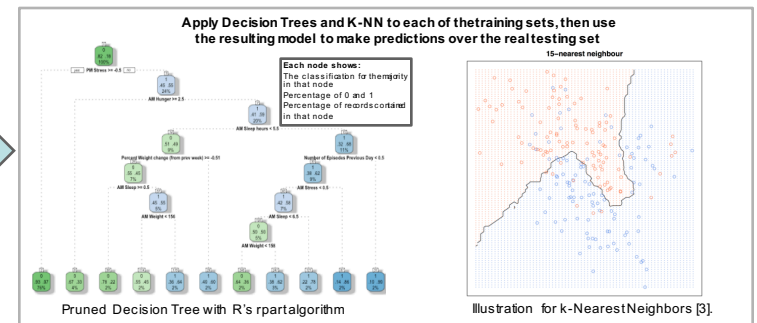


Fit a binomial generalized linear model into the real and synthetic sets to measure the consistency of correlation from each variable and Overeating Episodes

Comparing Decision Trees and K-Nearest Neighbors on Real and Synthetic Data



For each data set of CV



Performance Comparison Result and Conclusion

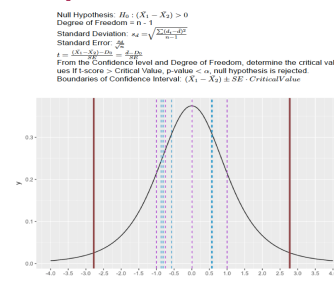
Performance metrics to evaluate and compare models:

- Accuracy
- Precision
- Recall
- Fallout
- F-Measure
- AUC

Paired t-test for the mean of each of the metrics in:

- RealTraining vs.Synthetic Training
- RealTraining vs.Real-Synth combined
- Synthetic Training vs.Real-Synth combined

- * Repeat this process for both Decision Tree and KNN results



T distribution, degree of freedom = 4, Confidence Level: 95%,
 brown line: Critical value, blue dashes: Tree measures' t score, pink dashes: KNN's Measures' t score
 As shown on the plot, none of the t statistics from either two machine learning techniques can exceed the critical value

References

- [1] Adult obesity in the United States [online]. Available: <https://stateofobesity.org/adult-obesity/>
- [2] Nowok Beata, Gillian M. Raab, and Chris Dibben. "Providing bespoke synthetic data for the UK Longitudinal Studies and other sensitive data with the synthpop package for R." *Statistical Journal of the IAOS Preprint* (2015): 1-12.
- [3] T Hastie, R Tibshirani, J Friedman. *Element of Statistical Learning*. Springer 2009.

Acknowledgement

This research was funded by the National Science Foundation (NSF) through the REU SITE: Data science research for safe, sustainable and healthy communities (CNS-1560229). This research was also conducted with the support and insightful feedback of Cole Polychronis, Brendan Foley, Quyen Hoang, Charan Sankaran, Vanshika Chowdhary, and Professor or Begisu Tulufom WPI.