

U.S. Food Environment Clustering and Analysis

STEPHANIE CUNNINGHAM, STUDENT – MASTER’S OF APPLIED DATA SCIENCE, BAY PATH UNIVERISTY

Examination of Food Environment Features and Outcomes with Machine Learning Techniques

INTRO

“Food environment” – culmination of socioeconomic and physical factors which shape the food habits of a population. A theme of food environment research is the idea of a “food desert,” generally defined as an area that is both low-income and low-access to food markets and thought to result in poor nutrition and health problems such as diabetes. Recent research calls this into question, with claims that choices drive poor health outcomes, rather than lack of nutritional access.

Goal – determine whether clustering areas by common food environment can reveal how underlying factors and health outcomes (here, diabetes rate) vary beyond food desert classification.

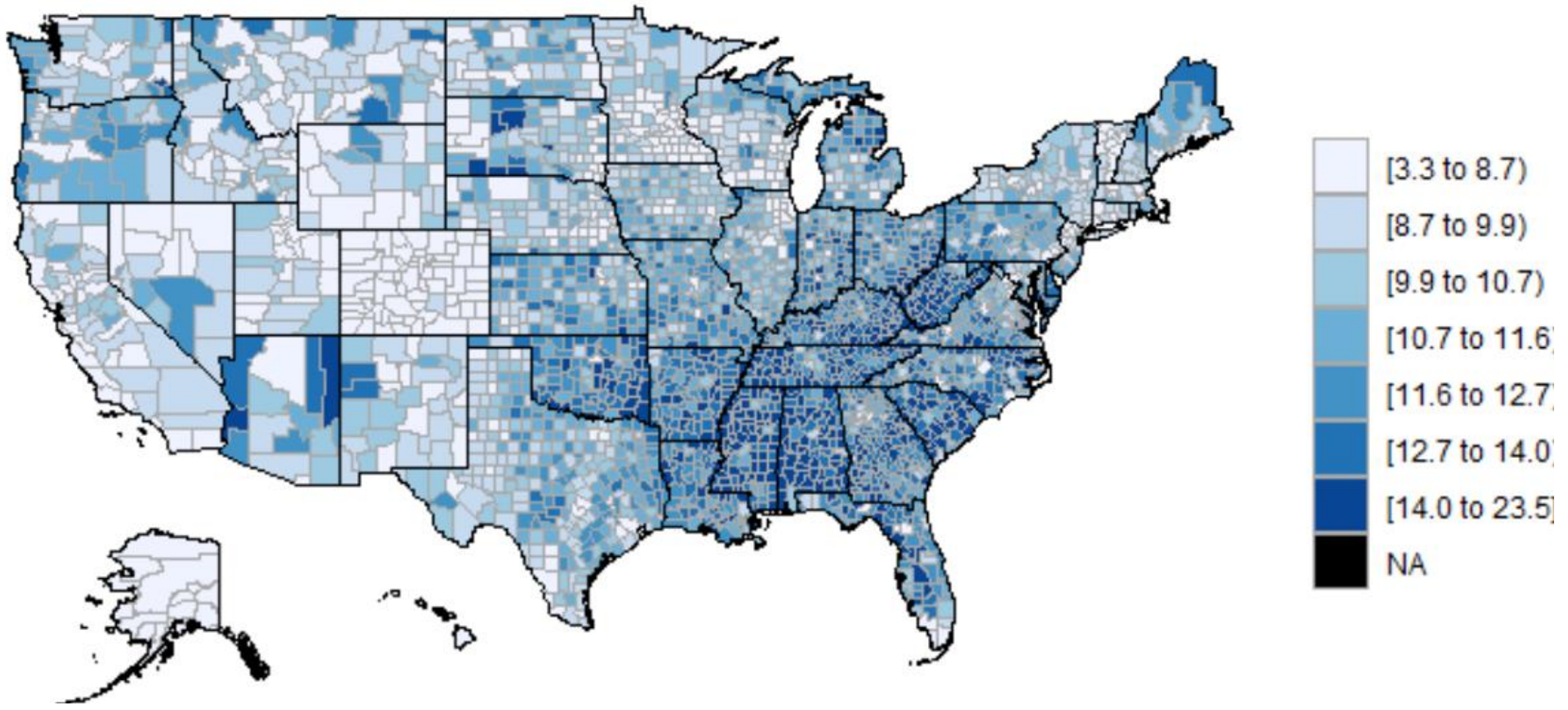
DATA

2020 Food Environmental Atlas (USDA)
2020 Food Access Research Atlas (FARA) (USDA)
2018 Labor Participation and Educational Attainment Data (US Census Bureau)

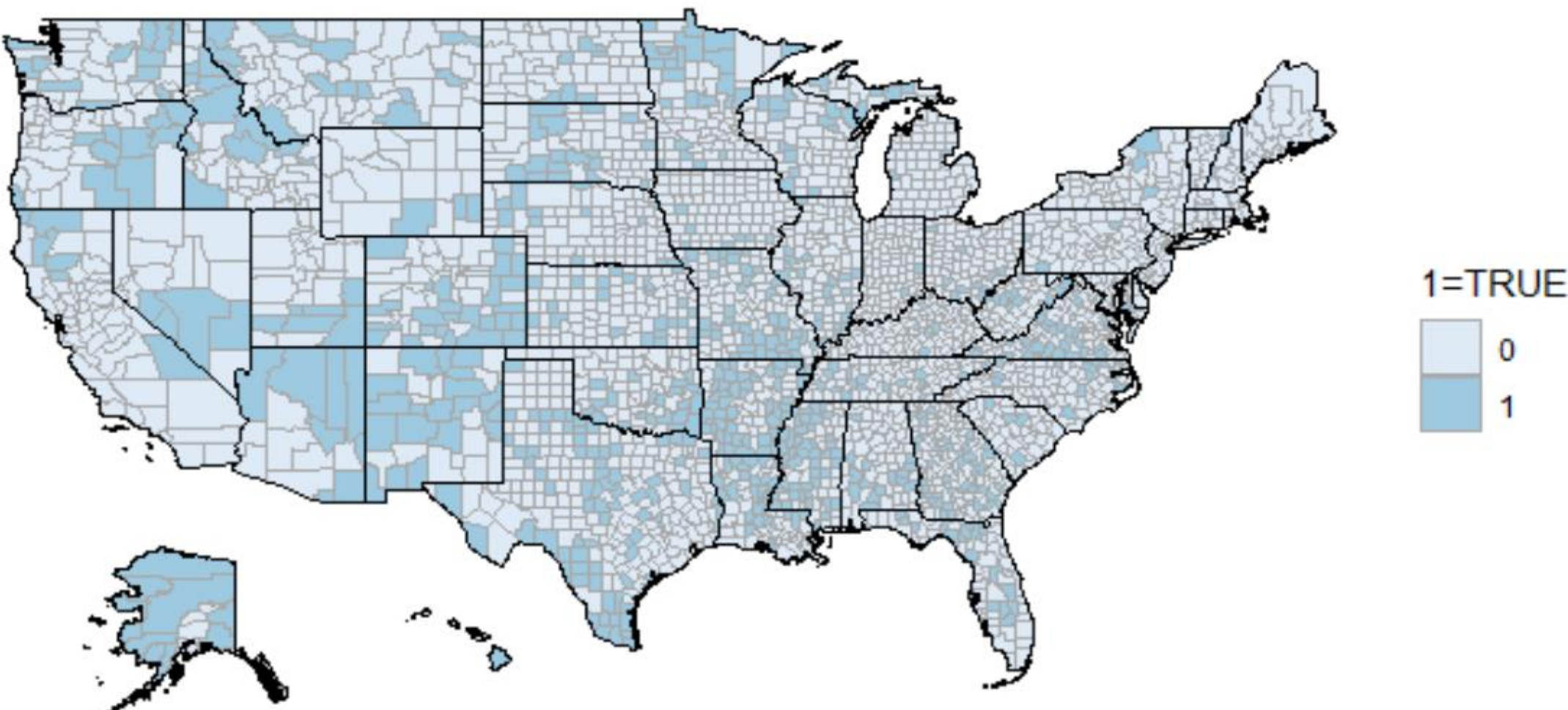
290 total features on all 3143 counties

Outcome “food desert status” labeled for counties with >1/3 tracts defined as food deserts based on FARA.

Adult Diabetes Rate 2013



Counties with >1/3 Tracts as Food Deserts

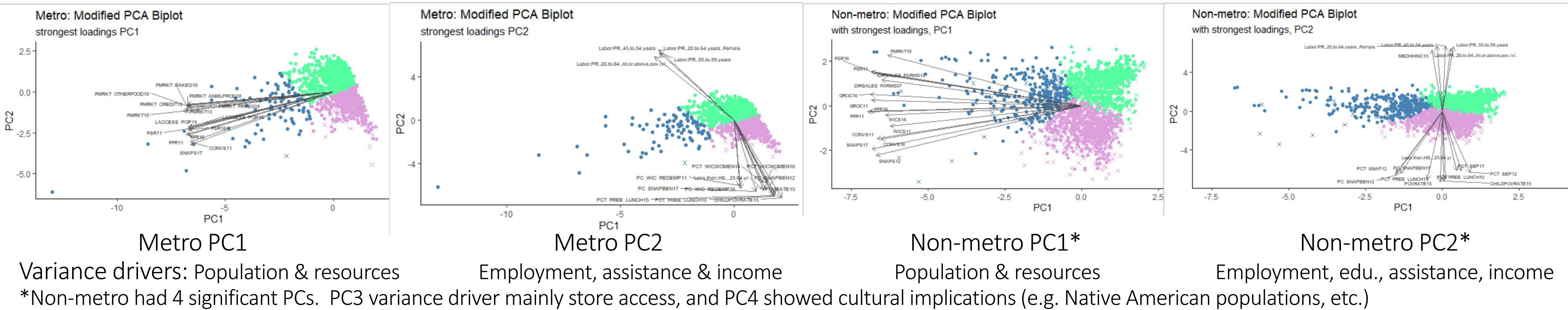


CONCLUSIONS

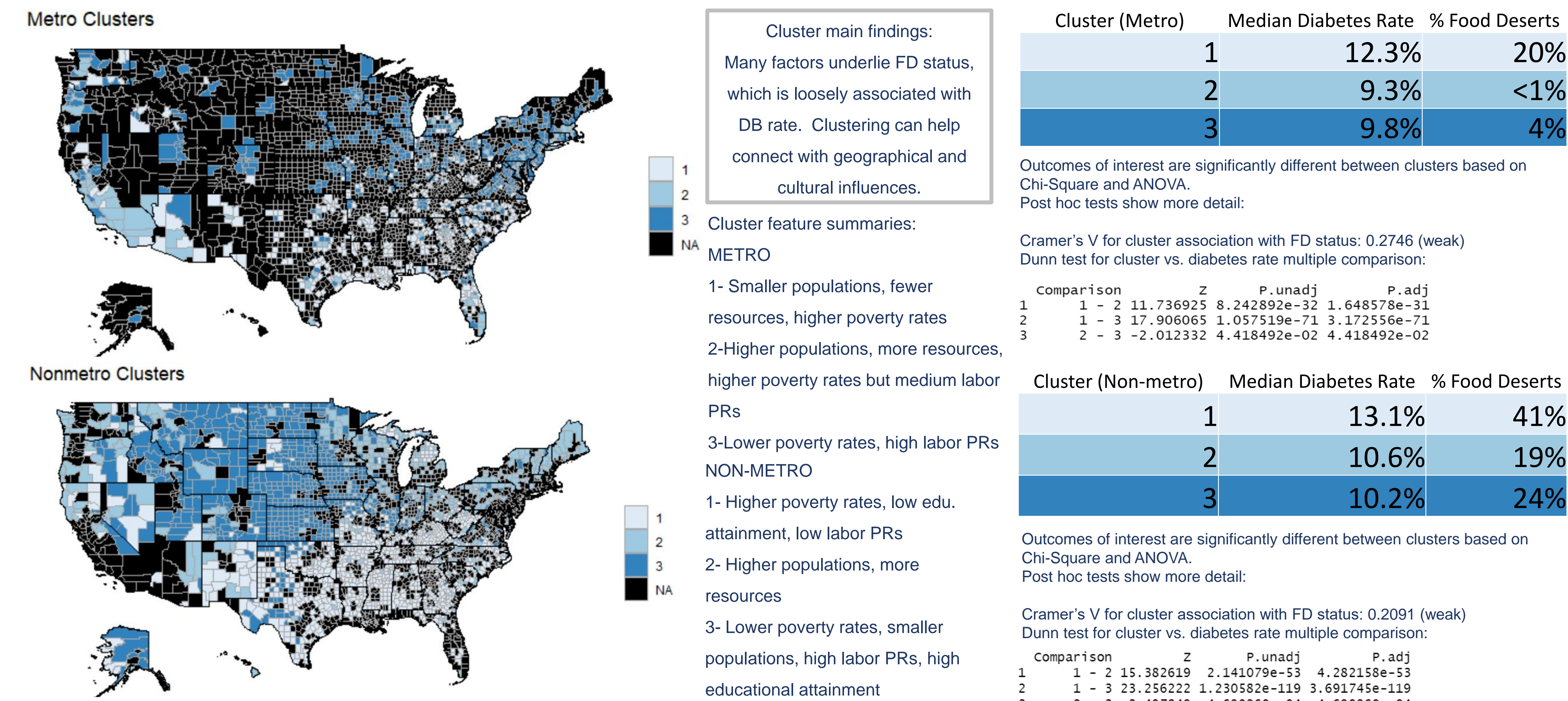
While food environment, food desert status, and diabetes rate are inter-related, certain factors associated with food desert status are more correlated with poor health outcomes than others. While poverty and access are important to food environment, this is not a simple causal relationship and further clustering work could reveal more personalized solutions for populations in need.

RESULTS

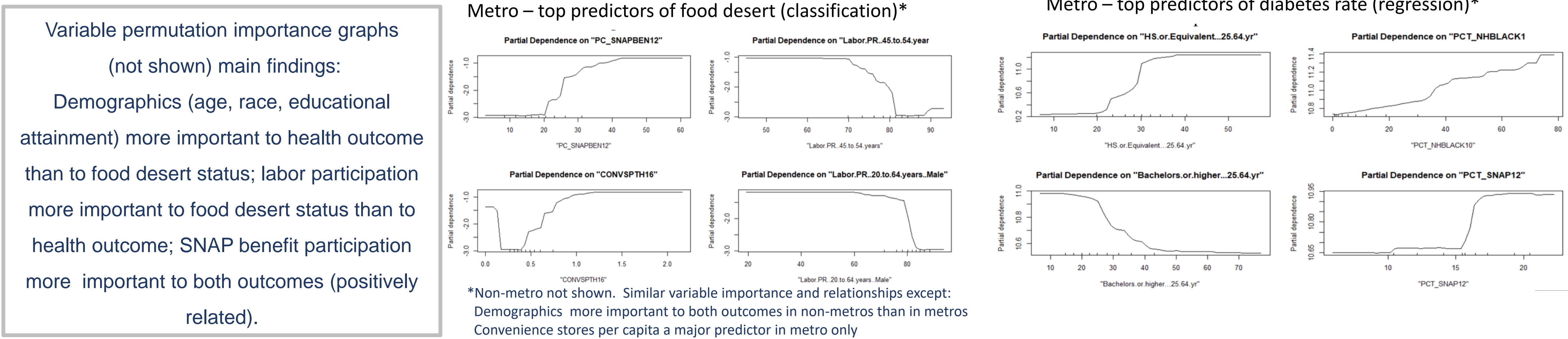
PCA & Clustering (with strongest loadings by PC)



Cluster Mapping and Outcome Comparison



Random Forest – Variable Importance & Partial Dependence



METHODS

Bayesian PCA

Selected for dimensionality reduction due to its proficiency at handling missing value imputation. Counties were separated by metro status to reduce influence of population over PCA results.

K-Means Clustering

3 food environment clusters were identified for both Metro and Non-metro counties. Strongest loadings were plotted on each principal component to impart meaning to clusters.

Statistical Analysis

Significance (Chi-Sq, ANOVA) and post hoc tests (Cramer's V and Dunn's, respectively) – (how) do clusters associate with food desert status and with health outcomes (diabetes rate)?

Random Forest

This supervised learning portion does not involve the clusters, only attributes and outcomes. Classification of food desert status and regression on diabetes rates. Used to compare predictive importance of variables (excepting those which directly define food desert e.g. poverty rate)

FURTHER RESEARCH

Additional clustering/stratifying work to identify more specific needs – variance in this research still driven largely by population size and then by income levels.

2017 policy change for SNAP-authorized stores to encourage healthier offerings – data not yet reflective. In both 2012 and 2017 data, partial dependence showing authorized SNAP stores per capita is a strong positive predictor of diabetes rate, and patterns follow closely with number of convenience stores. Research on policy effects is needed.