

# 95-851 Making Products Count: Data Science for Product Managers

## Homework 2: Clustering and Classification

### Fall 2020

Due 11:59 PM EST September 25, 2020

### Overview

On <https://www.countyhealthrankings.org/>, the University of Wisconsin's Population Health Institute has published rankings of all 3142 counties in the US by their health outcomes and behaviors. In class, we described some exploratory data analysis on the underlying data from 2017. In this assignment, you will apply clustering and classification techniques to the data from 2020 to determine if there are interesting groupings of counties with similar health outcomes and behaviors, and develop a predictive model to see which factors influence health outcomes. Such models can be used by individuals and governments looking to improve public health outcomes.

### Tasks

The data file for you to analyze is in the assignment in Canvas and also available at [https://www.countyhealthrankings.org/sites/default/files/media/document/analytic\\_data2020\\_0.csv](https://www.countyhealthrankings.org/sites/default/files/media/document/analytic_data2020_0.csv). Documentation of the measures in the data file is in [https://www.countyhealthrankings.org/sites/default/files/media/document/2020%20Analytic%20Documentation\\_0.pdf](https://www.countyhealthrankings.org/sites/default/files/media/document/2020%20Analytic%20Documentation_0.pdf). Use only the raw value columns (e.g. v001\_rawvalue) for your features; you don't need to analyze the numerators, denominators, confidence intervals, or various subsets of the features for this assignment. The remainder of <https://www.countyhealthrankings.org/> is worth exploring to obtain more context about how this data is collected and used to improve public health. You should turn in a Jupyter notebook with code and results including visualizations to answer the following questions (15 points total for the assignment):

- 1) What steps did you use for exploratory data analysis and subsequent data preparation on this data set? This should address summary statistics and handling of missing data and outliers (3 points)
- 2) Are there any noteworthy groupings of counties that have similar health outcomes and behaviors? You should use an unsupervised learning technique like clustering and show how you decided on the number of clusters. (4 points)
- 3) What are the five most important factors predicting premature death as shown by this data? In this data set premature death is defined as the number of years of potential life lost before age 75 per 100,000 population. Develop two models (using different

supervised learning approaches) to answer this question. Which of the two models do you believe is more accurate and how can you tell? Where would you focus public health efforts to reduce premature deaths in Allegheny County? (7 points)

- 4) Visualize the clusters you determined in #2 on a map of the US. (1 point)

## Expectations and Points Distribution:

Item	Points
Exploratory data analysis and data preparation	3
Development of clustering model	3
Determination of number of clusters	1
Visualization of clusters on map	1
Development of first supervised learning model predicting premature death	2
Development of second supervised learning model predicting premature death	2
List of 5 most important factors influencing premature death	1
Evaluating accuracy of the two supervised learning models	1
Recommendations for reducing premature death in Allegheny County	1
Total for full credit	15

## Hints/Suggestions:

1. Describe how you have chosen to handle outliers and missing values
2. Include visualizations such as histograms and boxplots where appropriate in describing the results of EDA
3. Describe how the findings from EDA influence your choice of data preparation and choice of modeling techniques
4. There are many possibilities for visualizing data on a map in Python, Geopandas <http://geopandas.org/> is one possibility Consider whether the visualization of the clusters on a map shows any regional patterns.
5. For the supervised learning models, consider how you will select the features and measure accuracy of the resulting models. How can you avoid overfitting?
6. Do not rename datafiles
7. Use relative path while importing the data

## Submission:

The submission will be a Jupyter Notebook with the name *DSPM\_HW2\_<your Andrew id>.ipynb*, and an html export of the notebook. The deadline for the homework is **September 25, 2020 11:59 PM EST.**