

# Statistical Learning for Public Policy I

## Problem Set 2

Oct 12, 2022

This lab uses data from Fatehkia et al. (2019). Their paper investigates the predictability of crime rates across urban neighborhoods with the help of interests recorded on Facebook. They use the Facebook Advertising API to collect data on different interests of Facebook users of a particular ZIP code area. They combine the collected user data with demographic information. Then, they fit various statistical models to predict assault, burglary, and robbery rates as a function of the Facebook and census demographic variables.

A subset of the authors' data are stored in `fbk.csv`. The table below provides a description for each variable in the dataset. All variables with `acs_` are from the American Community Census (2015) and all variables with `fb_` are from Facebook.

Name	Description
<code>zip</code>	ZIP code
<code>state</code>	State
<code>city</code>	City name
<code>assaults17</code>	Total reported assaults in 2017 per 100,000 capita
<code>acs_pop15_19</code>	% of population aged 15-19
<code>acs_pop18_24</code>	% of population aged 18-24
<code>acs_med_age</code>	Median age
<code>acs_white</code>	% of population one race White
<code>acs_black</code>	% of population one race Black or African-American
<code>acs_foodstamp</code>	% of households on food stamp benefits
<code>acs_med_inc</code>	Median family income
<code>acs_inc150k</code>	% households with income > 150K
<code>acs_inc25k</code>	% households with income <= 25K
<code>acs_pop18_24_nohs</code>	% of population 18-24 with less than high school degree
<code>acs_pop18_25_college</code>	% of population 18-24 with bachelors or higher degree
<code>acs_pop25_nohs</code>	% of population 25+ with less than high school degree
<code>acs_pop25_college</code>	% of population 25+ with bachelors or higher degree
<code>fb_hip_hop_music</code>	% of Facebook users (older than 18) interested in Hip Hop
<code>fb_...</code>	% of Facebook users (older than 18) interested in ...

Before working on this problem set, you might consider working through the Linear Regression Lab in ISL (p. 109-114).

### Question Set 1

- How many observations are in the data? How many variables are there?
- What share of the total number of observations is within the state of California?
- In which zip codes are the most (the least) assaults per capita?

## Question Set 2

Fit two bivariate linear regressions: In model A we wish to predict assaults using information on the share of Facebook users interested in Hip Hop music. In model B we wish to predict assaults using information on the share of Facebook users interested in first-person shooter games.

- a) Which of the two predictors accounts for more variation in the outcome variable?
- b) Predict the assault rate for a zip code in which the share of Facebook users interested in in first-person shooter games or hip hop music is equal to the average in the data. Construct a confidence interval for your predicted value.
- c) Construct a model C that leverages both the share of Facebook users interested in Hip Hop music and users interested in first-person shooter games to predict assaults. What is the increase in  $R^2$  relative to the best-performing bivariate linear model?
- d) Suppose that the share of Facebook users interested in Hip Hop music increases from 20 percent to 30 percent, i.e., a 10 percentage point increase. How many more assaults should we expected according to model C? What about an increase from 50 to 60 percent?

## Question Set 3

Construct a linear regression model that uses all continuous variables in the dataset.

- a) What is the increase in  $R^2$  relative to the best-performing bivariate linear model?
- b) Predict the assault rate for a typical zip code (i.e., a zip code for which all variables are centered on their mean) and compute the standard error. What percent of the uncertainty for the predicted value is due to estimation uncertainty?