

# Machine Learning - Mini-Project 4

PPHA 30546 - Professor Clapp  
Winter 2024

This assignment must be handed in via Gradescope on Canvas by **11:45pm Central Time on Wednesday, March 6th**. You are welcome (and encouraged!) to form study groups (of no more than 2 students) to work on the problem sets and mini-projects together. But you must write your own code and your own solutions. Please be sure to include the names of those in your group on your submission.

You should submit your answers in one of two ways:

1. As a single PDF containing BOTH a write-up of your solutions that directly integrates any relevant supporting output from your code (e.g., estimates, tables, figures) AND your code appended to the end of your write up. You may type your answers or write them out by hand and scan them (as long as they are legible). Your original code may be a Python (\*.py) or Jupyter Notebook (\*.ipynb) file converted to PDF format. OR
2. As a single PDF of a Jupyter Notebook (\*.ipynb) file with your your solutions and explanations written in Markdown.<sup>1</sup>

Regardless of how you submit your answers, be sure to make it clear what question you are answering by labeling the sections of your write up well and assigning your answers to the appropriate question in Gradescope. Also, be sure that it is immediately obvious what supporting output from your code (e.g., estimates, tables, figures) you are referring to in your answers. In addition, your answers should be direct and concise. Points will be taken off for including extraneous information, even if part of your answer is correct. You may use bullet points if they are beneficial. Finally, for your code, please also be sure to practice the good coding practices you learned in PPHA 30537/8 and comment your code, cite any sources you consult, etc.

You are allowed to consult the textbook authors' websites, Python documentation, and websites like StackOverflow for general coding questions. You are not allowed to consult material from other classes (e.g., old problem sets, exams, answer keys) or websites that post solutions to the textbook questions.

## 1 Overview

Ongoing concerns about election fraud have led to a wave of new policies designed to ensure the integrity of the vote. Unfortunately, there are concerns that these policies may also cause voter suppression. This has led to calls to make Election Day a national holiday as a partial way to

---

<sup>1</sup>Converting a Jupyter Notebook to PDF is not always straightforward (e.g., some methods don't wrap text properly). Please ensure that your PDF is legible! We will deduct points if we cannot read your PDF (even if you have the correct answers in your Notebook).

address those concerns.<sup>2</sup> The intuition behind this proposed reform is simple: making Election Day a holiday will improve voter turnout by giving people the time to vote without needing to take time off from work. But is there empirical evidence that suggests work constraints prevent people from voting?

In order to inform potential policy, you are tasked with answering the following related question: is having flexible work hours associated with being more likely to vote? At your disposal are two datasets, "vote.csv" and "work.csv," that we will call `vote_df` and `work_df`.

Both come from the U.S. Census Bureau and U.S. Bureau of Labor Statistics' **Current Population Survey (CPS)**.<sup>3</sup> In addition to being used to calculate monthly labor force statistics, the CPS provides information about a variety of economic and social well-being topics via supplemental questions that are asked to a rotating subset of the CPS sample. Because of the rotating nature of these supplements, individuals who are asked questions for one supplement are usually not asked asked the questions in other supplements.

Table 1 contains the names of the variables available in each dataset and their definitions. The `vote_df` dataset has a binary variable, which is appropriately called `vote`, that indicates the voting status of each individual in the data. Meanwhile, the `work_df` dataset has its own binary variable called `work` that indicates whether individuals have flexible work schedules. The two datasets share a set of core variables.

Table 1: Variable Names and Definitions

Variable	Definition	Dataset	
		Vote	Work
<code>vote</code>	Person voted in the last election	✓	
<code>work</code>	Person has a flexible work schedule		✓
<code>prtage</code>	Age	✓	✓
<code>pesex</code>	Sex	✓	✓
<code>ptdtrace</code>	Race	✓	✓
<code>pehspnon</code>	Hispanic origin	✓	✓
<code>prcitshp</code>	U.S. citizenship status	✓	✓
<code>peeduca</code>	Highest level of schooling	✓	✓

Individuals in one dataset are almost assuredly different than the individuals in the other dataset. As a result, for any individual in our data, we will either know their voting status or their work schedule, but we cannot know both simultaneously. Hence, we have a missing data problem. The plan to overcome this challenge is as follows:

- First, explore and clean the data.
- Train a Support Vector Machine (SVM) classifier on `work_df` that uses the demographic variables to forecast whether someone has flexible work hours.

<sup>2</sup>See, for instance, this Brookings Institution [blog post](#).

<sup>3</sup>Specifically, the `vote_df` dataset is based on the **CPS Voting and Registration Supplement**. The `work_df` dataset is based on the **CPS Work Schedules Supplement**. Note that these semi-synthetic datasets were created for pedagogical reasons, so results should be viewed accordingly.

- Apply the SVM classifier from the previous step to `vote_df` to predict whether the people in that dataset have flexible work hours.
- Regress voting status on the predictions obtained in the previous step.
- Adjust our regression estimate to account for measurement error in the imputed flexible work hours measure.

## 2 Data Analysis

1. Report the data type of each variable in each dataset.
2. Notice that all the variables in our dataset except for `prtag` are categorical, meaning that they take discrete values as opposed continuous values. Scikit-learn's `SVC()` function only accepts numerical values for training.
  - (a) First convert the target variables (`vote` and `work`) to binary forms by mapping them to 0s and 1s.
  - (b) For each of the remaining categorical variables compare the categories in the version of the variable in the `vote_df` and `work_df` datasets. Are there any discrepancies?
  - (c) Convert the categorical variables using a technique called “one-hot encoding” by creating multiple binary variables corresponding to each category using pandas' `get_dummies` function.
  - (d) We want each core variable between our two datasets to have the same structure for our prediction exercise, so whenever there's a discrepancy between the categories reported for a given variable, adjust the data to ensure that the core variables have the same structure between the `vote_df` and `work_df` datasets.

For instance, `work_df['prcitship']` has five categories while `vote_df['prcitshp']` has four. So add a column that's all zeros for the “FOREIGN BORN, NOT A CITIZEN OF” citizenship category after applying `get_dummies` to `vote_df['prcitshp']`.

3. Now that the datasets are set up, train a SVM classifier on the `work` data that fits the flexible work variable as a function of the core variables (and a constant). Be sure to standardize the predictors before fitting the model.<sup>4</sup> There are several choices to make when fitting a SVM, mainly the regularization or *cost hyperparameter* ( $C$ ) that penalizes observations that violate the margin/hyperplane and the *kernel* that introduces non-linearity to the SVM. Consider four values of hyperparameter  $C$ : 0.1, 1, 5 and 10 and three kernels: “linear”, “poly” and

---

<sup>4</sup>Since SVMs calculate the perpendicular distance from data points to the separating hyperplane, they require a constraint ( $\sum_{j=1}^p \beta_j^2 = 1$ ) on the parameters that define the hyperplane to ensure that the distances calculated are actually perpendicular. Since the values of these parameters will depend on how the data is scaled, it's important to standardize the predictors before estimating the SVM.

“sigmoid.” Use 5-fold cross-validation (5FCV) to determine which cost and kernel to use.<sup>5</sup> Report the cross-validation error rates of all 12 SVM models.<sup>6</sup>

4. Pick and report the value of  $C$  and kernel that minimize the 5FCV error rate. Use this model for the rest of the exercises.
5. What is the accuracy score of the model that you decided on when fitting to the `work_df` data in the previous question?
6. With the SVM model that you fit on `work_df`, impute the work schedules using the core variables from `vote_df`. The result is the imputed work flexibility measure needed for the main analysis. Compute and report descriptive (summary) statistics for the imputed measure.
7. Regress voting status on the imputed work schedule. Use age, age squared, and sex as predictors in addition to the imputed work schedule. Report, briefly interpret, and discuss the results.
8. Since we imputed the work schedules, there are likely to be some forecasts that are incorrect. To account for the bias this measurement error causes, we will need to divide the estimate of the parameter of interest by a scaling function

$$M(a, b) = \frac{1}{1 - 2b} \left( 1 - \frac{(1 - b)b}{a} - \frac{(1 - b)b}{1 - a} \right),$$

where  $a = \Pr(\widehat{work} = \text{"flexible"})$ ,  $b = \Pr(\widehat{work} \text{ is incorrectly labeled})$ , and  $\widehat{work}$  is the imputed value of the work variable from the SVM.<sup>7</sup> Write a simple function to compute  $M(a, b)$ . For the value of  $a$ , find the proportion of imputed work schedules that are flexible. For the value of  $b$ , use the cross-validation error rate from Question 3. Report  $a$ ,  $b$ , and  $M(a, b)$ .

9. Correct for the attenuation bias in your results from Question 7. Is the bias corrected version larger or smaller? Does the bias-correction change your previous result? Briefly explain.

---

<sup>5</sup>Shuffle the data randomly for splitting, and set `random_state=26`.

<sup>6</sup>Note that it may take some time to run this code.

<sup>7</sup>Note that this scaling function  $M(\cdot)$  is different from the margin defined in the SVM slides.