

Gender and Racial Disparities in STEM Workforce: A Machine Learning Approach

Luxin Tian

June 10, 2020

Abstract

While previous studies have elaborated on the persisting gender and racial gaps in STEM workforce engagement, there lacks consensus on the classification standards of STEM occupations. In an effort to clarify the boundary of STEM fields, we construct a feature-based machine learning classifier that learns from rich features of occupations and bottom-up responses from workers that indicate their STEM engagement. The classifier achieves 84% accuracy on the test set and suggests that 7.6% of occupations are mislabelled by the official STEM classification. Using the prediction of the feature-based classifier, we examine the gender and racial gaps in STEM workforce engagement in the US in 2018. While the regression coefficients suggest female and racial minorities are significantly underrepresented in the STEM workforce, we find the estimations based on the official STEM classification are prone to overestimate the disparities. We also find the disparities are heterogeneous across states in the US.

Contents

1	Introduction	3
2	Literature Review	4
2.1	Gender and Racial Gaps in STEM Workforce	4
2.2	Classifications of STEM Occupations	6
3	Classifying STEM Occupations	8
3.1	Data	8
3.2	Classification Model	10
3.3	Classification Results	12
4	Measuring Gender and Racial Disparities	12
4.1	Data	13
4.2	Regression Model	13
4.3	Empirical Results	13
5	Conclusion	14
A	Tables and Graphs	18

1 Introduction

As Science, Technology, Engineering, and Mathematics, collectively being referred to as STEM, has been creating substantial economic impacts on human society, policymakers and employers have placed increasing emphasis on fostering the growth of STEM workforce as well as promoting equal opportunities. Although the market preference for STEM workers is generic, it has been a historical and persisting pattern that there exist significant gaps in the population distribution and the labor earnings across different groups by gender, race, and Hispanic origin in the STEM workforce.

There has been plenty of research that has elaborated on the persisting gender and racial disparities in STEM workforce engagement [Landivar \(2013\)](#); [Beede et al. \(2011\)](#); [Baird et al. \(2017\)](#); [Carnevale et al. \(2011\)](#). However, there lacks a consistent standard that classifies STEM occupations. Previous studies on the STEM labor market usually rely on two established classification standards. The official STEM occupation list is published by the Standard Occupational Classification (SOC) Policy Committee in collaboration with the Bureau of Labor Statistics. The SOC STEM occupation list provides a top-down classification that is approved by the US government, in which STEM occupations are labeled by public policy experts and researchers. However, the official classification has been criticized as it can fail to capture the potential variations in the specific STEM skills and knowledge that an occupation actually requires ([Anderson et al., 2018](#)). Another classification approach, created by the Brookings Institution, alleviates the weakness of the SOC standard by scoring each occupation in terms of its STEM intensity based on the characteristic information provided by the Occupational Information Network (O*Net) of the US Department of Labor. While this provides feature-based classifications, the scores are as well evaluated by researchers and experts and thus lack objective evidence.

This paper contributes to the existing literature by constructing a bottom-up feature-based machine learning classifier to clarify the boundary of the STEM workforce. This classifier learns from rich information about occupations and bottom-up responses from workers that indicate their STEM engagement. After performing

model selection, feature selection, and hyperparameter tuning, we adopted a Logistic classifier with cross-validation and optimize the feature matrix by Random Feature Elimination based on a Linear Support Vector Machine kernel. The classifier achieves 84% accuracy on the test set and suggests that 7.6% of occupations are mislabelled by the official STEM classification.

Using the prediction of the feature-based classifier, we re-examine the gender and racial gaps in STEM workforce engagement in the US in 2018. We find that even though the regression coefficients suggest that female and racial minorities remain significantly underrepresented in the STEM workforce, the estimations based on the official STEM classification are prone to have overestimated the disparities. The results are robust across model specifications that use either the predicted labels or probabilities of an occupation being STEM. We also examined the heterogeneity of gender disparities in STEM workforce engagement across the 50 2-states in the US.

The remaining of this paper is organized as follows. Section 2 reviews previous research. Section 3 introduces the data employed for classification and the model tuning and evaluates the results. Section 4 presents the empirical analysis into gender and racial disparities based on the classification results. Section 5 concludes.

2 Literature Review

2.1 Gender and Racial Gaps in STEM Workforce

Empirical evidence shows that, in the United States, women, African Americans, and Hispanic have been underrepresented in the STEM workforce in recent decades. As [Landivar \(2013\)](#) suggest, according to the micro-level data from the decennial US Census from 1970 to 2000 and the American Community Survey in 2011, although women had increased their presence in STEM workforce since the 1970s, they remained underrepresented in engineering and computer science occupations, which comprised over 80% of all STEM employment. According to [Beede et al. \(2011\)](#)'s calculation using the 2000 and 2009 American Community Survey (ACS) micro-level data, across

the 2000s, women’s representation in STEM jobs remains 25% lower than that of men, although they comprised nearly half of all the workers in the US, and the share of women with college-level education had increased in the workforce. US Census data also suggests that African Americans and Hispanics have been underrepresented in the STEM workforce as well(Landivar, 2013).

While women and racial minorities are underrepresented in STEM fields, their engagement in STEM education and workforce turns out to be an important force driving the alleviation of earning gaps across genders and races. Beede et al. (2011) and Baird et al. (2017) find that education backgrounds in STEM fields and working in a STEM occupation can render a large increase in monetary earnings, and women’s wage premium is significantly larger than that of men on average for working in a STEM occupation. Similar mechanisms are also found between racial groups. Carnevale et al. (2011) finds that, between 2005 and 2009, earning gaps across racial minority groups are much narrower in the STEM workforce than that in the non-STEM workforce, after controlling for the age variable.

Existing research has elaborated on the direct evidence from demographical data of gender and racial gaps that persisted in the US labor market. Given potentially complicated self-selection processes underlying people’s decisionmaking on education and employment, there is very little empirical research making efforts on identifying the causal relationship between gender and racial identity and engagement in STEM occupations and employment outcomes. On the contrary, many researchers tried to understand the gender and gaps related to the STEM workforce qualitatively, from the perspectives of economics, psychology, sociology, as well as cognitive science, and empirical evidence has been presented to verify the hypothesis. The gender and racial gaps can form from an early stage, even from as early as high school and college when students choose their subjects and field of studies, and the selection can compound and reinforce over years in their career path (Fryer Jr and Levitt, 2010; Hill et al., 2010; Jaeger et al., 2017). Furthermore, Kanny et al. (2014) review the explanations of gender gaps in STEM fields in the past forty years and identify five dominant narratives that emerged in previous research: “Based on a systematic review of 324 full

texts spanning the past 40 years of scholarly literature, five dominant meta-narrative explanations emerged: individual background characteristics; structural barriers in K-12 education; psychological factors, values, and preferences; family influences and expectations; and perceptions of STEM fields”. Wang and Degol (2017) further summarize six explanations for the minor presence of women in STEM fields from biological and sociocultural perspectives: “(a) cognitive ability, (b) relative cognitive strengths, (c) occupational interests or preferences, (d) lifestyle values or work-family balance preferences, (e) field-specific ability beliefs, and (f) gender-related stereotypes and biases”. These explanations can provide evidence on model specifications for research that examine the role of gender and racial identity in employment status and outcomes.

2.2 Classifications of STEM Occupations

While plenty of research has tried to estimate the inequalities related to the STEM workforce, there has been a lack of consensus on the classification standards of STEM occupations (Landivar, 2013; Anderson et al., 2018). Many efforts have been made to clarify the boundary of STEM fields by both the government and researchers.

As is summarized by Anderson et al. (2018), there are two dominant existing methods for defining STEM occupations at different levels. Most of the previous academic and policy research employs the official standards published by the Standard Occupational Classification Policy Committee (SOCPC), which collaborated with nine federal agencies for the 2010 revision of the Standard Occupational Classification manual by the Bureau of Labor Statistics. As is introduced by Landivar (2013), the SOC manual, which provides 539 specific occupational categories associated with 23 occupational groups, classifies occupations based on the type of work performed, and the SOCPC creates the classifications of STEM occupations based on the SOC codes. This method is characterized by its top-down approach of identifying work features, as the classification standard is suggested by experts and researchers from general to discrete occupation categories. However, there have been criticisms of the SOCPC approach considering that it does not reflect on the potential variations in

the extent to which an occupation actually requires STEM skills and knowledge but rather hinges on top-down subjective perceptions by the public sector representatives and experts(Anderson et al., 2018).

Another classification approach that alleviates the concerns about the SOCP approach is created by the Brookings Institution (Rothwell, 2013). This approach uses occupation characteristics data provided by the Occupational Information Network (O*Net) of the US Department of Labor. Each occupation in the SOC occupation codes is evaluated and scored by its level of knowledge required from one of the four STEM fields, which gives each occupation four scores from 0 to 7 corresponding to the extent in which this occupation requires knowledge in each STEM field. Then, occupations with a score at least 1.5 standard deviations above the mean in at least one STEM field would be classified as a high STEM occupation, and, those whose average score of the four scores for the four fields is at least 1.5 standard deviations above the mean would be classified as super STEM occupations(Anderson et al., 2018; Rothwell, 2013). This approach is based on the characteristics of occupations and thus is expected to be more accurate to reflect on the extent to which occupation is STEM-like. Nevertheless, the evaluation process is still somewhat subjective and only based on a limited number of characteristics, and the threshold, 1.5 standard deviations, is also somewhat an arbitrary cut-off.

In light of such caveats of using these two classification standards, Anderson et al. (2018) suggests a new bottom-up approach based on a micro-level survey that interviewed workers in the real workforce whether their job requires STEM knowledge and skills. This survey also includes the categorical label of interviewee's occupation in terms of the SOC standard codes, and thus this self-reported classification can be compared with the existing "occupation-based classifications". Anderson et al. (2018) identified that 10% to 15% occupations that actually require STEM knowledge and skills are not counted by the traditional STEM classifications. With classifications more sensitive to practical variations across workers, they find that women are more likely to engage in jobs that are not consistently classified as STEM with their self-reported outcome. Furthermore, they find that there is no significant gender gap in

STEM workforce engagement, suggesting that the results of previous research had overestimated the inequality. They also verified that the wage premium for women working in STEM fields that are consistently classified by different methods is larger than other women working in periphery STEM fields that are inconsistently classified. This paper contributes to the literature, although being rich already, of the STEM workforce and provides strong evidence that a more accurate classification method that takes into accounts the variations across specific jobs.

Therefore, by integrating multiple data sources and employing a machine learning approach, this paper is expected to construct a more accurate classification of STEM occupations, which, although remains at the occupation level, can quantitatively reflect on the probability that one occupation should be counted as STEM workforce. As is suggested by [Anderson et al. \(2018\)](#), such enhancement in the reliability of classification can potentially challenge the existing literature by providing a more unbiased estimation of gender and racial inequalities in the STEM workforce.

3 Classifying STEM Occupations

3.1 Data

To address the limitation of previous classification standards, we take a feature-based approach that integrates multiple data sources to construct an eclectic classification system for predicting STEM occupations. Concretely, we merge the Occupational Information Network (O*NET) database with the American Life Panel (ALP) Survey data. The merged data matches the bottom-up responses from workers who indicate whether the specific job they are doing should be considered as STEM with a comprehensive set of feature variables that characterize all the occupations coded in SOC categories.

The ALP surveys, conducted by RAND, interview a panel of more than 6,000 nation-wide representatives through the internet each year who are sampled based on probabilistically models ([Pollard and Baird, 2017](#)). We mainly select the responses to

two surveys, MS480 and MS436, which are available on RNAD’s website. In MS480, interviewees were directly asked to report whether they are working in a STEM job with the definition of a STEM job described in the question. In MS480, interviewees were also asked to indicate their occupations coded in the SOC categories. Since the panel representatives continuously participate in the surveys, the records on the individual level can be matched across different surveys in the ALP dataset. The two questions in MS480 and MS436 are described in the following way (Pollard and Baird, 2017):

MS436 Q3 occupation: In which category is your current occupation?

MS480 C2 Whether or not you have a STEM degree, are you currently working in a STEM job? Note that this relates to the tasks you do, and not the industry you work in. For example, an engineer for a bioengineering research firm would be in a STEM job, but an administrative assistant at the same bioengineering research firm would NOT be in a STEM job. (1 Yes, I am currently working in a STEM job; 2 No, I am not currently working in a STEM job.)

After merging the responses of these two questions and dropping any records with missing data, we obtained a set of labeled data that consists of 1424 records.

To extend the bottom-up STEM labels to out-of-sample occupations, we match the ALP survey data with the O*NET database. The O*NET contains a comprehensive set of descriptive variables that characterizes all aspects of features of the nearly 1,000 occupations in the US economy. The variables are organized into the O*NET Content Model (Peterson et al., 1999) that constitutes 42 files belonging to six main categories of features in regards to occupations and workers. We selected three files that are relevant to identifying STEM occupations from the database, Knowledge, Skills, and Work Activities, and joined the feature variables on the SOC codes of each occupation. For each feature variable, the original data file stores the value of the feature and a series of descriptive statistics of the samples that are used to quantify the feature. We reshaped the data to a wide panel and group any multiple measures

of each feature of each occupation by taking the mean of the values in each group and obtained a matrix of 774 occupations with 112 feature variables in total.

Finally, we match the ALP data with the O*NET to construct a labeled dataset for training a machine learning model. After dropping any rows with missing data, the final dataset contains 1279 records with 109 feature variables, and 24.7% of the responses are labeled as STEM.

3.2 Classification Model

We build a pipeline for training a machine learning classifier on the merged labeled data and predicting the STEM labels for each occupation in the O*NET dataset. The pipeline contains three steps: standardization, feature selection, and classification.

First, we standardize the feature variables by removing the mean value of each variable and then dividing the difference by the standard deviation of the variable. This step ensures that all the feature variables are centered around zero and maintain variance in similar orders. Otherwise, some variables in higher orders of magnitude can dominate the cost function in the supervised learning algorithm and prevent the model from learning other features.

Second, we perform feature selection with Recursive Feature Elimination (RFE) to prevent overfitting. RFE selects feature variables to be included in the feature matrix for further classifications by recursively eliminate feature variables and train an estimator on each of the selected subsets of features.

We use a Support Vector Classifier (SVC) with a linear kernel as the estimator for the RFE algorithm. An SVC is a Support Vector Machine classifier that constructs a maximal margin hyperplane to separate the training observations of different classes such that the perpendicular distances from the observations to the separating hyperplane are maximized. Concretely, an SVC estimator is the solution to find the maximal margin M and a hyperplane defined by $\mathbf{x}\beta^T$, where \mathbf{x} is the feature matrix and β is the parameter vector $(\beta_0, \beta_1, \dots, \beta_p)$ in (1).

$$\begin{aligned}
& \max_{\beta_0, \beta_1, \dots, \beta_p, M} M \\
& \text{s.t.} \quad \sum_{j=1}^p \beta_j^2 = 1 \\
& y_i (\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) \geq M \quad \forall i = 1, \dots, n
\end{aligned} \tag{1}$$

For each subset of features, the SVC gives a vector β , which is orthogonal to the separating hyperplane. Each element in the vector β specifies the importance of the corresponding feature, based on which the features with the least importance are eliminated from the current subset of features. After recursively repeating the elimination until the prespecified number of features are left in the matrix. We further use 5-fold cross-validation to select the optimal number of features to be included for further classification. The cross-validation result of feature selection is shown in Figure 1. We also tried other feature selection methods and selected RFE with Linear SVC based on its joint performance with the Logistic Regression Classifier over other pipeline designs. The model selection scores are shown in Table 1.

Third, we train a classification model on the merged data with bottom-up responses, which is then used for predicting the STEM labels of all the unlabelled occupations in the O*NET dataset. We experimented with several classifiers, the Logistic Regression Classifier, the Support Vector Machine Classifier, and tree-based classifiers such as the Decision Tree Classifier and the Random Forest Classifier. We randomly split the merged data into a 0.8 training set and a 0.2 test set. For each experimental classifier, we performed a grid search with 5-fold cross-validation to tune the hyperparameters. The cross-validation scores for each pipeline are shown in Table 1.

Finally, we selected the Logistic Regression Classifier, which achieves mean accuracy scores of 0.8407 and 0.8438 on the training and test sets, respectively. Concretely, a Logistic Regression fits the following model 2: $\forall i \in 1, 2, \dots, n$:

$$\begin{aligned}
Z_i &= \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i \\
\Pr(\text{STEM}_i = 1) &= \text{sigmoid}(Z_i)
\end{aligned} \tag{2}$$

We train the Logistic Classifier with the optimal parameter from the grid search and predict the STEM labels for all the occupations in the O*NET dataset. With the default 0.5 activation threshold of the sigmoid function, the classifier does not perform well, as it labels far fewer STEM occupations than the SOC official classification. We experimented with different probability thresholds to label the STEM occupations and plot the number of STEM labels against the thresholds value in Figure 2. The slope of the curve indicates the sensitivity of the selection of the threshold value. We select 0.43 as the threshold value where the sensitivity shows a sharp swift, indicating there is can exist a sharp gap between STEM and non-STEM occupations. We obtained a set of 774 occupations associated with predicted probability and labels.

3.3 Classification Results

Using the classifier we construct, we can predict the score of each occupation that indicates the extent to which it is STEM-like. As is shown in Figure 3, our feature-based classifier gives a similar result to the SOC official classification. However, 5.5% occupations that are identified as STEM by SOC are not labeled as STEM by our classifier, and 2.1% occupations that are not approved as STEM are actually showing STEM features instead. These occupations are referred to as periphery STEM occupations by Anderson et al. (2018), which reflects the limitations of the traditional classification standard in identifying STEM occupations that are emerging and fading out.

4 Measuring Gender and Racial Disparities

Our empirical analysis of the gender and racial disparities in the STEM workforce replicates several previous studies, which adopt either the SOC official classification or the Brookings scoring approach. We demonstrate the potential of our feature-based classification of STEM occupations by evaluating the difference in results given by different classification approaches. Concretely, we fit a linear Logistic Regression model with the American Community Survey (ACS) data by the US Census Bureau.

4.1 Data

We use the 2018 one-year data of the ACS, which contains 3,214,539 individuals across the United States, with sample weights specified. We select sex, racial codes, and SOC code variables to match our classification result on the SOC code. After dropping any records with missing data, we obtained a cross-section that consists of 1,247,722 samples. The descriptive statistics is shown in Table 2.

4.2 Regression Model

As a demonstration of using our classification results, we fit reduced-form simple linear models to identify the gender and racial disparities in the STEM workforce, without examining the causal effects and mechanisms. Concretely, we specify a Linear Regression model and a Logistic Regression model in (3) and (4).

$$Z = \alpha + \beta_1 \text{Gender} + \beta_2 \text{Race} + \beta_3 \text{Gender} \times \text{Race} + \gamma X + \epsilon \quad (3)$$

$$Y = \text{sigmoid}(Z) \quad (4)$$

Gender is a dummy variable which equals 1 when the individual is male and 0 if female. Race is a matrix of dummy variables indicating the racial groups. X is a matrix of control variables, including age, a binomial term of age, and school attainments. Y is the dependent variable. When Y is a discrete binary variable, such as the predicted STEM label or the official STEM label, we will fit the Logistic Regression that combines (3) and (4), and when Y is the natural log of the predicted probability of being STEM, we fit a weighted OLS regression model (3). We assign the samples with the weights in the ACS data.

4.3 Empirical Results

We fit the regression models using both the nation-level data and state-level data. Through the nation-level regressions, we identify the gender and racial disparities in

STEM workforce across the United States in 2018. The result indicates that females remain significantly underrepresented in the STEM workforce, using as the dependent variable either the official STEM labels, our feature-based labels, or the natural log probability of being engaged in STEM occupations given by our classifier. The result is robust across model specifications with different sets of covariates. However, the coefficient of the gender variable in the with our feature-based STEM labels is lower than that of the model with the official labels, indicating the traditional classification approach may have resulted in researchers overestimating the gender disparities in STEM workforce engagement.

We also identified that racial minority groups such as African American, American Indian, Alaska Native, American Indian, and Hawaiian or Pacific Islanders are significantly underrepresented in STEM workforce compared to White Americans, and females in these minority groups are further underrepresented in comparison to male. However, Asian people show a significantly higher probability of engaging in STEM workforce than White Americans, but female Asians are still less likely to engage in STEM occupations than male Asians. Again, we find that the results may have been overestimated by researchers using the traditional classification approach. Full nation-level regression results are shown in Table 3, 4, and 5.

Using the state-level data, we only fit the WLS regressions with the predicted probability as a STEM score of each occupation, with all control variables included. The result shows that gender disparities are heterogeneous across the states in the US. We plot the regression coefficients in Figure 4. States with the most significant gender disparities are South Dakota, Vermont, Iowa, and North Dakota.

5 Conclusion

Based on the survey data as STEM labels and the occupation feature data, we constructed a bottom-up feature-based classifier for predicting STEM occupations. The classifier achieves 84% accuracy on the test set and suggests that 7.6% of occupations are mislabelled by the official STEM classification. Using the prediction of the

feature-based classifier, we examine the gender and racial gaps in STEM workforce engagement in the US in 2018. While the regression coefficients suggest female and racial minorities are significantly underrepresented in the STEM workforce, we find the estimations based on the official STEM classification are prone to overestimate the disparities. The results are robust across model specifications that use either the predicted labels or probabilities of an occupation being STEM. We also find the disparities are heterogeneous across states in the US. States with the most significant gender disparities are South Dakota, Vermont, Iowa, and North Dakota. The bottom-up feature-based classifier is expected to provide more accurate STEM labels for occupations that better captures the variations in the extent to which STEM tasks are involved in a specific job. Therefore, our empirical results should be more unbiased in comparison to previous empirical findings as there should be less measurement errors in the delineation of STEM occupations.

References

- Anderson, Drew, Matthew Baird, and Robert Bozick**, “Who Gets Counted as Part of America’s STEM Workforce? The Implications of Different Classification Approaches for Understanding the Gender Gap in STEM,” 2018.
- Baird, Matthew D, Robert Bozick, and Mark Harris**, “Postsecondary Education and STEM Employment in the United States: An Analysis of National Trends with a Focus on the Natural Gas and Oil Industry,” 2017.
- Beede, David N, Tiffany A Julian, David Langdon, George McKittrick, Beethika Khan, and Mark E Doms**, “Women in STEM: A gender gap to innovation,” *Economics and Statistics Administration Issue Brief*, 2011, (04-11).
- Carnevale, Anthony P, Nicole Smith, and Michelle Melton**, “STEM: Science Technology Engineering Mathematics.,” *Georgetown University Center on Education and the Workforce*, 2011.
- Hill, Catherine, Christianne Corbett, and Andresse St Rose**, *Why so few? Women in science, technology, engineering, and mathematics.*, ERIC, 2010.
- Jaeger, Audrey J, Tara D Hudson, Penny A Pasque, and Frim D Ampaw**, “Understanding how lifelong learning shapes the career trajectories of women with STEM doctorates: The life experiences and role negotiations (LEARN) model,” *The Review of Higher Education*, 2017, 40 (4), 477–507.
- Jr, Roland G Fryer and Steven D Levitt**, “An empirical analysis of the gender gap in mathematics,” *American Economic Journal: Applied Economics*, 2010, 2 (2), 210–40.
- Kanny, Mary Allison, Linda J Sax, and Tiffani A Riggers-Piehl**, “Investigating forty years of STEM research: How explanations for the gender gap have evolved over time,” *Journal of Women and Minorities in Science and Engineering*, 2014, 20 (2).
- Landivar, Liana Christin**, “Disparities in STEM employment by sex, race, and Hispanic origin,” *Education Review*, 2013, 29 (6), 911–922.
- Peterson, Norman G, Michael D Mumford, Walter C Borman, P Jeanerret, and Edwin A Fleishman**, *An occupational information system for the 21st century: The development of O* NET.*, American Psychological Association, 1999.
- Pollard, Michael and Matthew D Baird**, “The RAND American Life Panel: Technical Description,” 2017.
- Rothwell, Jonathan**, *The hidden STEM economy*, Metropolitan Policy Program at Brookings, 2013.

Wang, Ming-Te and Jessica L Degol, “Gender gap in science, technology, engineering, and mathematics (STEM): Current knowledge, implications for practice, policy, and future directions,” *Educational psychology review*, 2017, *29* (1), 119–140.

A Tables and Graphs

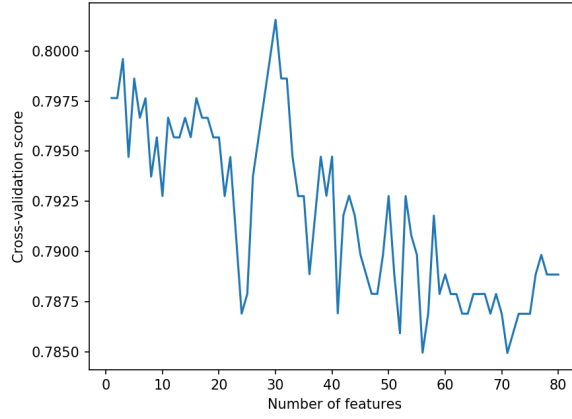


Figure 1: Cross-validation scores of RFECV with Linear SVC

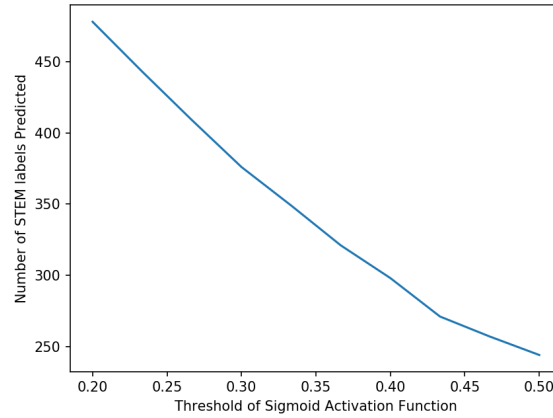


Figure 2: Sensitivity of Number of STEM Labels Predicted on Activation Threshold Value

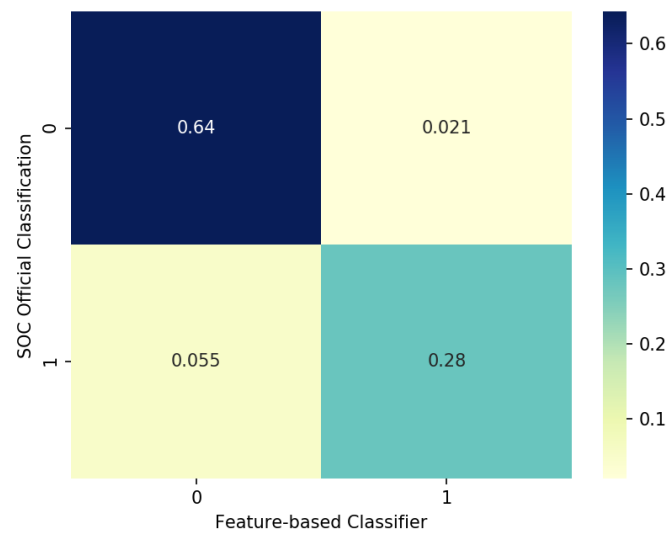


Figure 3: Confusion Matrix of Feature-based Classification and SOC Official Classification

Table 1: Performance of Experimental Feature Selection Methods and Classifiers

Pipeline	Feature Selection	Classifier	Best CV Score on Training/Test Data	Precision by Labels	Recall by Labels	F1 by Labels
1	None	Random Forest	0.86706, 0.80859	0.88, 0.58	0.86, 0.63	0.87, 0.61
2	LASSO with CV	Random Forest	0.86706, 0.79685	0.87, 0.57	0.87, 0.57	0.87, 0.57
3	Ridge with CV	Random Forest	0.86706, 0.79297	0.87, 0.56	0.86, 0.57	0.86, 0.56
4	Logistic with CV	Random Forest	0.86706, 0.79688	0.87, 0.57	0.87, 0.57	0.87, 0.57
5	RFECV with Linear SVC	Random Forest	0.86706, 0.8125	0.88, 0.63	0.87, 0.57	0.87, 0.60
6	None	Logistic with CV	0.82111, 0.84375	0.86, 0.75	0.95, 0.50	0.90, 0.60
7	RFECV with Linear SVC	Logistic with CV	0.84066, 0.84376	0.87, 0.72	0.93, 0.55	0.90, 0.62

Table 2: Descriptive Statistics of Variables

[illegible]

Table 3: Weighted Logistic 1: Gender and Racial Disparities in STEM Workforce Engagement

VARIABLES	Logistic: Z = SOC Official STEM Labels				
	I	II	III	IV	V
Male	1.0634*** (0.0046)	1.0647*** (0.0047)	1.0623*** (0.0047)	0.9741*** (0.0051)	1.0245*** (0.0057)
Age		0.1046*** (0.0009)	0.1036*** (0.0009)	0.0529*** (0.0010)	0.0529*** (0.0010)
Age-squared		-0.0010*** (0.0000)	-0.0010*** (0.0000)	-0.0005*** (0.0000)	-0.0005*** (0.0000)
African American			-0.1776*** (0.0081)	0.0050 (0.0090)	-0.0902*** (0.0181)
American Indian			-0.6622*** (0.0322)	-0.2725*** (0.0358)	-0.0960 (0.0635)
Alaska Native			-1.1766*** (0.1617)	-0.5372*** (0.1783)	-0.5636 (0.3763)
American Indian or Alaska Native			-0.7308*** (0.0763)	-0.2433*** (0.0850)	0.1750 (0.1349)
Asian			0.5535*** (0.0083)	0.2870*** (0.0092)	0.6425*** (0.0148)
Hawaiian/Pacific			-0.7973*** (0.0699)	-0.4503*** (0.0766)	-0.1780 (0.1292)
Other Races			-1.1044*** (0.0173)	-0.4970*** (0.0191)	-0.3140*** (0.0332)
Two or More Races			-0.2013*** (0.0152)	-0.1577*** (0.0167)	0.0665** (0.0286)
Female African American					0.1233*** (0.0209)
Female American Indian					-0.2542*** (0.0766)
Female Alaska Native					0.0285 (0.4276)
Female American Indian					-0.6428*** (0.1722)
Female Asian					-0.5575*** (0.0186)
Female Hawaiian/Pacific					-0.4020** (0.1596)
Female of Other Races					-0.2669*** (0.0403)
Female of Two or More Races					-0.3287*** (0.0350)
const	-1.7948*** (0.0039)	-4.1593*** (0.0192)	-4.1092*** (0.0193)	-5.0293*** (0.0564)	-5.0611*** (0.0564)
Controls	N	N	N	Y	Y
Num of observations	1,247,722	1,247,722	1,247,722	1,247,722	1,247,722

Standard errors in parentheses. *p < 0.1; **p < 0.05; ***p < 0.01.

Table 4: Weighted Logistic 2: Gender and Racial Disparities in STEM Workforce Engagement

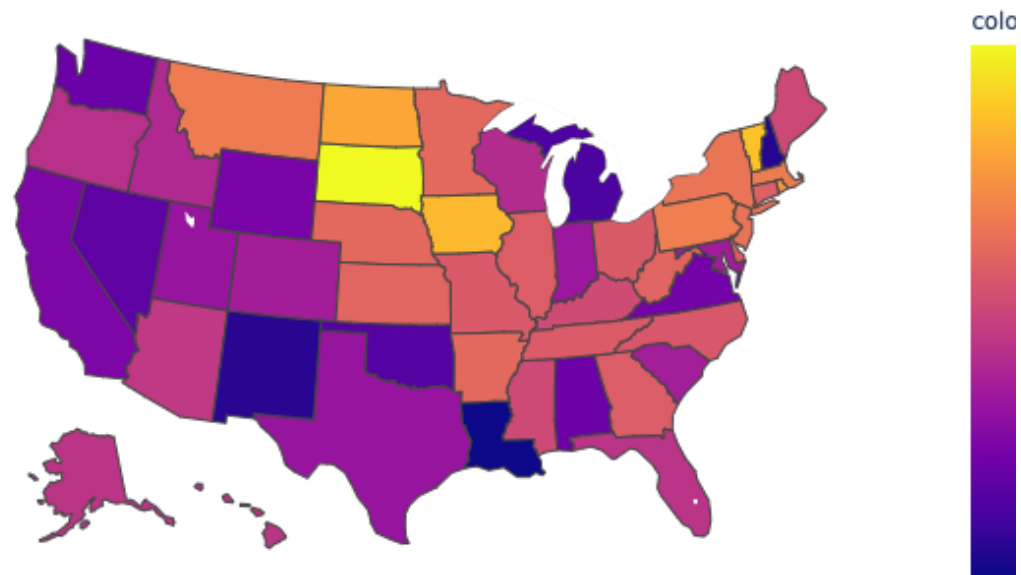
VARIABLES	Logistic: Z = Featured-based Predicted Label				
	I	II	III	IV	V
Male	0.9616*** (0.0048)	0.9604*** (0.0048)	0.9578*** (0.0048)	0.8433*** (0.0053)	0.8896*** (0.0059)
Age		0.1028*** (0.0009)	0.1018*** (0.0009)	0.0474*** (0.0011)	0.0474*** (0.0011)
Age-squared		-0.0010*** (0.0000)	-0.0010*** (0.0000)	-0.0004*** (0.0000)	-0.0004*** (0.0000)
African American			-0.2552*** (0.0085)	-0.0759*** (0.0095)	-0.1722*** (0.0192)
American Indian			-0.7168*** (0.0343)	-0.3044*** (0.0382)	-0.0560 (0.0659)
Alaska Native			-1.0376*** (0.1615)	-0.3001* (0.1815)	0.0914 (0.3216)
American Indian or Alaska Native			-0.6945*** (0.0789)	-0.1662* (0.0887)	0.2721** (0.1378)
Asian			0.5710*** (0.0084)	0.2894*** (0.0093)	0.5885*** (0.0151)
Hawaiian/Pacific			-0.7825*** (0.0730)	-0.4007*** (0.0803)	-0.0575 (0.1301)
Other Races			-1.1592*** (0.0187)	-0.5305*** (0.0206)	-0.3212*** (0.0349)
Two or More Races			-0.2006*** (0.0157)	-0.1502*** (0.0173)	0.0881*** (0.0294)
Female African American					0.1240*** (0.0222)
Female American Indian					-0.3620*** (0.0806)
Female Alaska Native					-0.5457 (0.3871)
Female American Indian					-0.6896*** (0.1784)
Female Asian					-0.4674*** (0.0189)
Female Hawaiian/Pacific					-0.5218*** (0.1644)
Female of Other Races					-0.3093*** (0.0429)
Female of Two or More Races					-0.3529*** (0.0362)
const	-1.8634*** (0.0040)	-4.2024*** (0.0199)	-4.1476*** (0.0200)	-5.0016*** (0.0605)	-5.0329*** (0.0605)
Controls	N	N	N	Y	Y
Num of observations	1,247,722	1,247,722	1,247,722	1,247,722	1,247,722

Standard errors in parentheses. *p <0.1; **p <0.05; ***p <0.01.

Table 5: WLS: Gender and Racial Disparities in STEM Workforce Engagement

VARIABLES	log(Pr{STEM=1})				
	I	II	III	IV	V
Male	0.2399*** (0.0016)	0.2323*** (0.0016)	0.2263*** (0.0016)	0.1151*** (0.0014)	0.1075*** (0.0016)
Age		0.0499*** (0.0003)	0.0497*** (0.0003)	0.0258*** (0.0003)	0.0257*** (0.0003)
Age-squared		-0.0005*** (0.0000)	-0.0005*** (0.0000)	-0.0002*** (0.0000)	-0.0002*** (0.0000)
African American			-0.1299*** (0.0025)	-0.0636*** (0.0022)	-0.1522*** (0.0034)
American Indian			-0.2356*** (0.0105)	-0.0811*** (0.0092)	-0.0326** (0.0133)
Alaska Native			-0.3062*** (0.0492)	-0.0608 (0.0433)	-0.0153 (0.0603)
American Indian or Alaska Native			-0.3250*** (0.0215)	-0.1392*** (0.0190)	-0.0929*** (0.0270)
Asian			0.1850*** (0.0033)	0.0234*** (0.0029)	0.0396*** (0.0045)
Hawaiian/Pacific			-0.2702*** (0.0190)	-0.1256*** (0.0167)	-0.1001*** (0.0247)
Other Races			-0.4396*** (0.0038)	-0.1604*** (0.0035)	-0.0964*** (0.0047)
Two or More Races			-0.0775*** (0.0049)	-0.0556*** (0.0044)	-0.0470*** (0.0065)
Female African American					0.1513*** (0.0045)
Female American Indian					-0.0934*** (0.0184)
Female Alaska Native					-0.0961 (0.0866)
Female American Indian					-0.0917** (0.0379)
Female Asian					-0.0265*** (0.0059)
Female Hawaiian/Pacific					-0.0474 (0.0336)
Female of Other Races					-0.1372*** (0.0068)
Female of Two or More Races					-0.0158* (0.0087)
const	-1.9956*** (0.0012)	-3.1434*** (0.0061)	-3.0968*** (0.0061)	-3.1684*** (0.0096)	-3.1644*** (0.0096)
Controls	N	N	N	Y	Y
R-squared	0.0175	0.0501	0.0653	0.2757	0.2767
Num of observations	1,247,722	1,247,722	1,247,722	1,247,722	1,247,722

Standard errors in parentheses. *p < 0.1; **p < 0.05; ***p < 0.01.



The values are the coefficients of the state-level WLS regression with the natural log of the predicted probability of being engaged in STEM as the dependent variable. States with the most significant gender disparities are South Dakota, Vermont, Iowa, and North Dakota.

Figure 4: State-level Gender Disparities in STEM Workforce Engagement in the US (2018)