

Project Deliverable 4

All contractors commissioned by the state for major construction projects need to report their ethnic and gender makeup of the work forces. WGBH would like to understand the data contained in those Summary of Workforce Utilization reports. Furthermore, WGBH is interested in getting data-driven insights of the impact drawn upon specific groups of workers between 2019 to 2020. The data is given in PDF format and organized by hours spent per project per organization. Our goal is to first extract data in proper formats from the PDF files and then run some analysis.

Logistics

Weekly Meeting with the PM

- Lingyan Jiang is Thurs 11:30 AM - 1:00 PM

Weekly Meeting With WGBH

- Paul Singer, - every other Thurs 11:30 AM - 1:00 PM
- Spark Liason - Greta Bruce

Contact List

- Client Paul Singer paul_singer@wgbh.org,
- Spark Liason Greta Bruce gretab@bu.edu,
- PM Lingyan Jiang lingyanj@bu.edu,
- Students Rep Jena Jordahl jenajj@bu.edu,

Elisa Cordeiro Lopes elisacl@bu.edu, Richard Lee rlee99@bu.edu, Murtadha Bahrani murtadha@bu.edu, Carmen Sabrina Araujo sabrinaa@bu.edu.

Github accounts

elisa3lopes, rlee99, murtio, carmen-araujo, jenajjedu

This is a draft of your final report that has been reviewed by your client. It includes all visualizations, results, data, and code up to this point, along with proper documentation on how to reproduce your results, compile and use your codebase, and navigate your dataset. Your team will submit this as a PR.

Introduction

All contractors commissioned by the state for major construction projects need to report their ethnic and gender makeup of the work forces. The WGBH would like to understand the data contained in those Summary of Workforce Utilization reports. Furthermore, the WGBH is interested in getting data-driven insights of the impact drawn upon specific groups of working forces between 2019 to 2020. The data is given in PDF format and organized by hours spent per project per organization. However, these PDF files are constructed in a way where no regular PDF parser can easily extract data. They were also given in packages of hundreds of pages so that without any form of automation it would be very time consuming to be able to analyze the data.

After working through this project, our team thought of three important questions we wanted to answer. The first question was how we would extract data from our PDF files. Our PDF datasets were formatted in a way where it was extremely difficult to parse any data in a clean and precise manner. Our description of our attempts to parse this data will go into further detail below, however this portion of the project took up most of our time and a lot of time and effort was put into our PDF parser. The second question was if there was a difference between state-paid contractual hours based on color and/or sex. This was the meat of the project that we were fortunately able to get to because of the incredible work our team did to complete our parser. The third and last question we had was why all new hire hours were 0. This is highly unusual as it is quite common for workers to start a job casually whenever they see an opening.

Data

The data is collected by a Massachusetts state office, DCAMM, [Division of Capital Asset Management and Maintenance](#), and reported out to the community via an [annual report](#). WGBH requested additional documentation from the state so they could independently verify the numbers in the annual report. Through a freedom of information act request, WGBH was able to receive monthly construction workforce utilization reports. The reports are kept in PDF

format. DCAMM already provided WGBH twelve monthly reports for 2019 and in March they were to provide the data from 2020. On April 12, DCAMM released monthly reports for 2020.

Later other data may be of interest to analyze. Our team has only the construction data to analyze not the design data which is also included in the annual report. The annual report includes other data on the contractor's business location where payments were made and the location of the worksite.

Recently, WGBH let us know that we can download the database of WBE and MBE from a state site to validate the volume of contract hours reported to WBE companies and compare to the annual report. By correlating the data, we will assist WGBH in verifying the DCAMM numbers published in the yearly report, and identify new patterns.

The 2019 and 2020 dataset of construction data is organized as tables of projects summaries by month per contractor, trade and level of experience, such as bridges, buildings, etc. The important statistics per company includes, their types of workers listing the number of hours worked by race, sex, and ethnicity. For this project, no additional datasets are required to be extracted, but our team is open to get any other information as it seems relevant to analyze. An example of a file is April 2019:

<https://drive.google.com/file/d/1brxGTjfkhwKRXPAbzDwHI4bP6J08Xwtz/view?usp=sharing>

Creating the Parser

Phase 1: We used Python libraries PyPDF2 and Tabula to scrape the data from the PDF files and then used acrobat to save the files in CSV format. Each method produced the same misalignments between the hourly data rows and the company/trades header data. This issue stemmed from the PDF merging cells to pretty print the data for human readability. Initially we had no idea what was necessary so we attempted to set up Grobid on SCC to parse the PDFs into XML files. After our first 2 weeks we had tried 7 different methods (real time, tabula software, tabula library, Grobid, PyPDF2, managing data after tabula, and transforming data back into PDFs)

Phase 2: We were soon able to build a parser that extracted individual project names, project codes, and the contractors/companies involved in each. However, this parser still missed a few contractors/companies in each listing. We had many issues using commas as a splitter as it divided the numbers containing commas into two.

Phase 3: Successfully created a parser that could extract and create a pandas dataframe with columns: project id & name, contractor name, construction trade name, and detail lines per trade level name. Our next issue was extracting numbers with ". ". When separating numbers into distinct entries, we were getting arrays like: ['0.000.000.000.00', '2.800.00'] , instead of: [0.0, 0.0, 0.0, 0.0, 0.0, , 2.8, 0.0].

Phase 4: After our first parser, the team spent a lot of time trying to find alternate solutions to our problems because nothing was working out. Fortunately, we were finally able to create a program that produces clean dataframes. This was done by reading the data directly from the PDF into pandas dataframes and using a Json format to make a custom shape for our PDFs to be extracted from. However, we still ran into many issues such as: inconsistencies in the data reading empty numbers, misaligned rows and columns, two columns for one number, rows being chopped off and placed in a different row, and white space characters embedded in column heading fields.

Phase 5: All bugs listed in Phase 4 were fixed and our parser successfully worked on different datasets.

In Depth Description of our Working Parser:

We used PyPDF2 to extract the number of pages, the interactive tabula tool to create bounding box x,y coordinates, and tabula python library which utilizes pandas directly when it chunks out the file page by page. Next, we created a python script to read each pdf file in the input directory and produce a CSV file into a second directory. The file contains a denormalized model of the monthly project and contractor workforce hours performed per ethnicity and gender. One output CSV file is created per input PDF file. To keep the data appropriately marked, we added month and year data to the dataframe from the filename being parsed.

By generating the monthly columnar CSV files, we can build a mini-data mart for querying four different hierarchical trees. One tree for Project/Contractor/Trade/Experience Level, one for time series(month and year) and two others for Ethnicity and Gender. All arms of the tree tie count hours worked per contract per month. The structure of organizing the data is commonly called a data cube and the schema strucalled a star schema. Finally, to do the analysis work, we read our file-based data cube into a single pandas DataFrame again using a script. From the combined dataset, we could easily execute the group-by statements to compute percentages of the money received per ethnicity and per gender.

Analysis

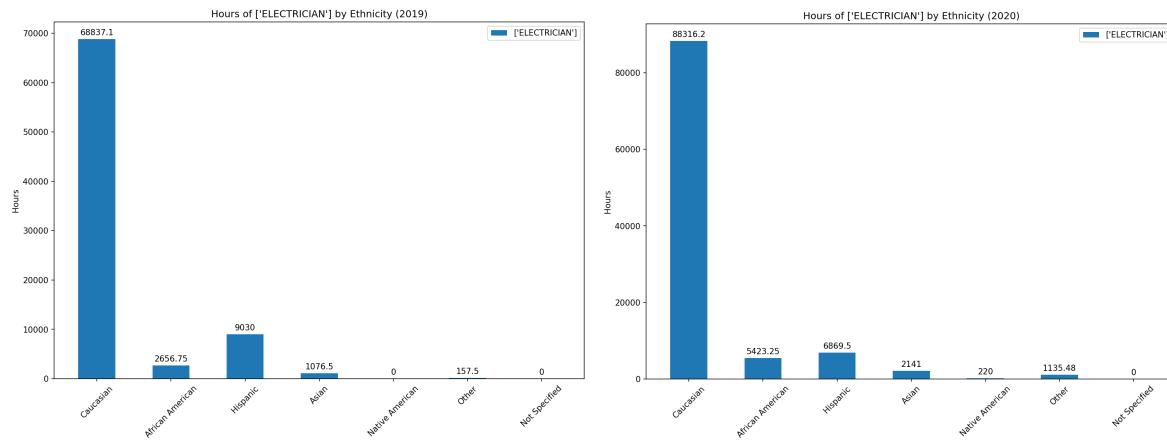
- The color of money and the trade of Electricians.

Is there any discrimination? To answer this question we first need to find out the mean work hours of each group across trades.

	Caucasian	African American	Hispanic	Asian	Native American	Other	Not Specified
mean	6461	221	824	81	3	108	2
ELECTRICIAN	68837	2656.75	9030	1076.5	0	157.5	0
HVAC (ELECTRICAL CONTROLS)	617.25	0	0	0	0	0	0

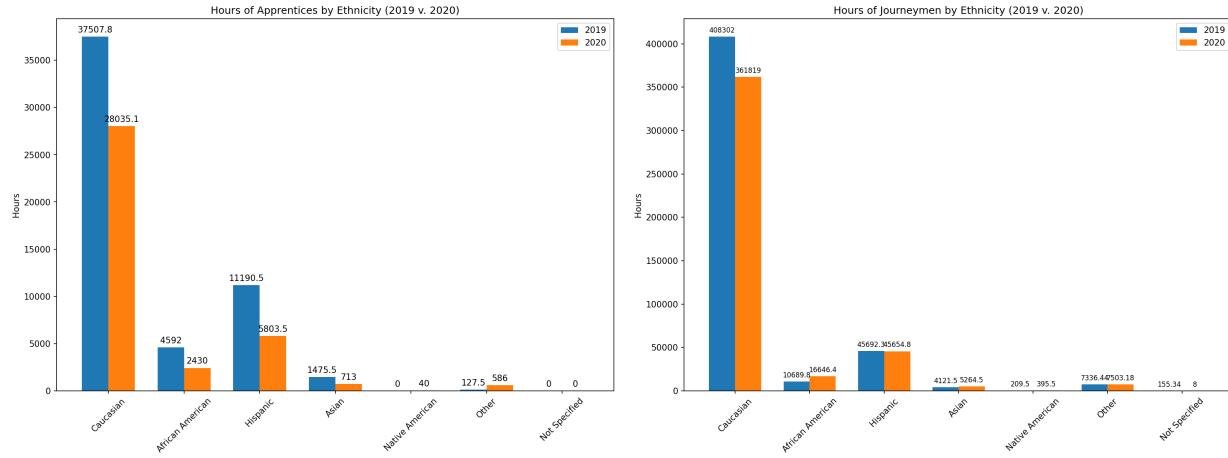
HVAC (TESTING AND BALANCING - AIR)	366	0	0	0	0	0	0
HVAC (DUCTWORK)	226	0	8	0	0	0	0

The next step is to show how the number of hours of some groups, for a given trade, is above or below its mean hours. We consider the 'ELECTRICIAN' trade in the 2019 reports for comparison. The number of work hours assigned to each work group is above the group's mean hours. However, we have reservations about how the name of the trade is playing a role here. For example, when investigating the number of hours for other related trades, such as 'HVAC (ELECTRICAL CONTROLS)', 'HVAC (TESTING AND BALANCING - AIR)', and 'HVAC (DUCTWORK)', we can see that they are dominated by Caucasian work groups. When adding these trades under 'ELECTRICIAN' trades, discriminatory patterns might be observed. We hypothesize that categorizing work hours under too many overlapping trades is an obstacle in investigating the color of money. For more charts related to trades see supplementary materials.



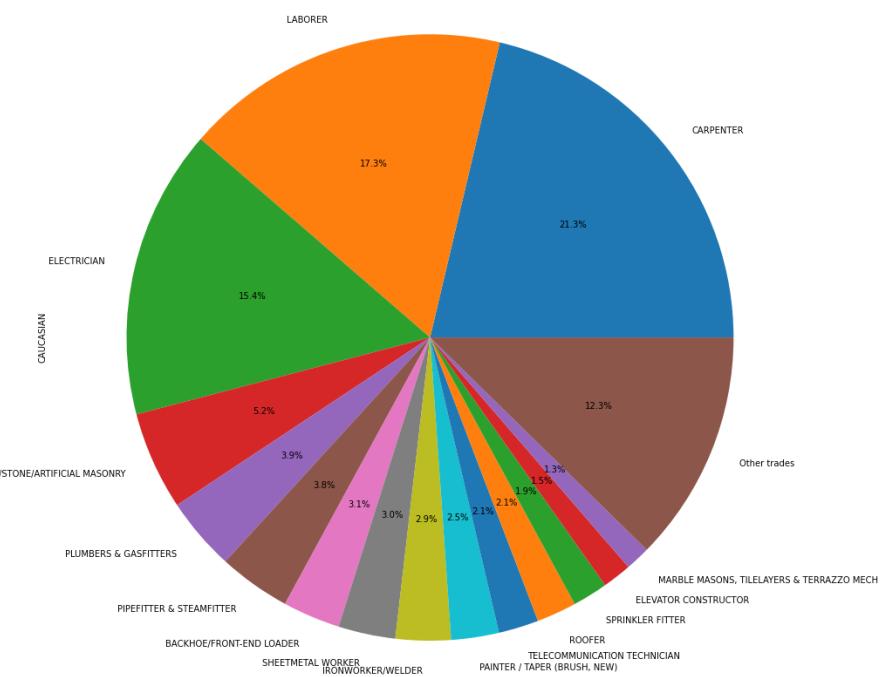
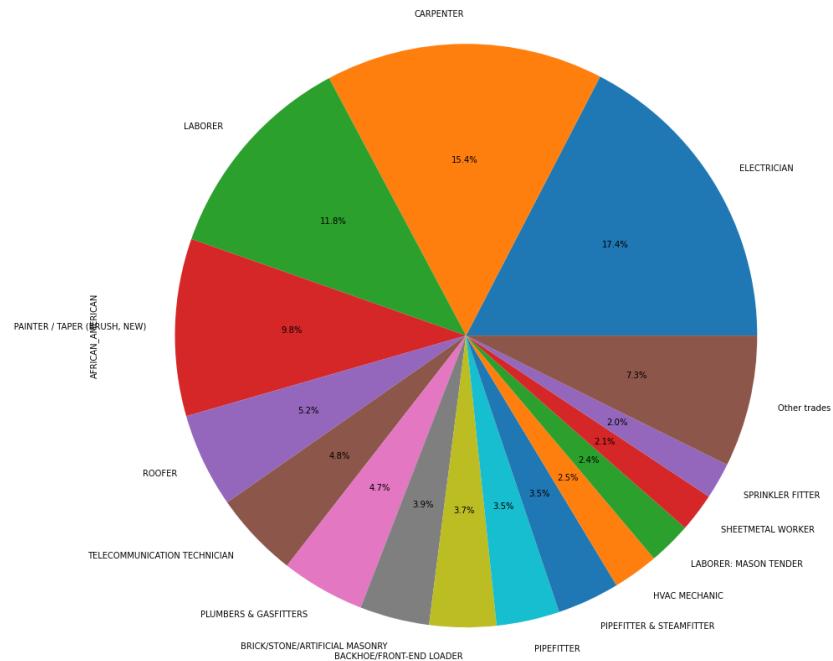
- Trends of hiring across years.

The question we want to investigate is how the number of workers differ among groups across years. The number of apprentices between 2019 and 2020 has witnessed the same patterns across groups. All groups, except for 'Other', have less apprentices in 2020 than 2019. For journeymen, it is slightly different. The 'Caucasian' work group witnessed a drop between 2019 and 2020, while 'Hispanic', 'Asian' and 'Other' numbers were almost the same between the two years. The only noticeable increase can be seen in the 'African American' work groups. For more charts of the number of workers see supplementary materials.



- Distribution of groups over trades

The distribution of workers within a group over the trades is different from one group to another. For example, The dominant trades for ‘African American’ work groups are ‘ELECTRICIAN’, ‘CARPENTER’, ‘LABORER’, and ‘PAINTER’, respectively. On the other hand, the dominant trades for ‘Caucasian’ work groups are ‘CARPENTER’, ‘LABORER’, ‘ELECTRICIAN’, ‘BRICK/STONE/ARTIFICIAL MASONRY’, respectively. While we cannot draw firm conclusions about the color of money, due to the same reasons stated in observation 1, we can notice that there are trades more popular within certain groups. For example, the ‘Asian’ work group has nearly 30% of its workers in ‘TELECOMMUNICATION TECHNICIAN’. See supplementary materials for more information.



- As a final note, in our initial analysis of the data we found interesting patterns when viewing African Americans and Women and so here we did a detailed analysis of trades and companies to try to understand what type of jobs were attracting African Americans and/or Women and what companies were employing these two groups. All the graphs will be under Supplementary Materials. Here, we will only highlight a few key findings.

- First, I went ahead and ranked **trades** based on two criteria: total number of hours that went to African Americans or Women and percentage of total hours for that trade that went to African Americans or Women. I will begin with the African American grouping.
 - When I ranked trades based on the **total number of Apprentice hours** worked by **African Americans** in 2019, I found that these trades had the following percentage of African American hours out of all the ethnicities:

For ELECTRICIAN (1880.0 hours) : 8.5874%

For CARPENTER (1640.5 hours) : 19.1961%

For PIPEFITTER (536.0 hours) : 12.9861%

For ROOFER (337.5 hours) : 16.7245%

For LABORER (118.0 hours) : 3.2855%

- Conversely, when I ranked trades based on the **percentage of hours of Apprentice** work that belong to African Americans for 2019, I found that these trades had the following total hours worked by African Americans:

For CARPENTER (19.196%) : 1640.5 hours

For ROOFER (16.724%) : 337.5 hours

For BRICK/PLASTER/CEMENT MASON (16.667%) : 8.0 hours

For PIPEFITTER (12.986%) : 536.0 hours

For IRONWORKER (11.88%) : 32.0 hours

- When I ranked trades based on the **total number of Journeymen hours** worked by African Americans in 2019, I found that these trades had the following percentage of African American hours out of all the ethnicities:

For LABORER (1689.75 hours) : 1.85809%

For PAINTER / TAPER (BRUSH, NEW) (1502.0 hours) : 9.009%

For ELECTRICIAN (776.75 hours) : 1.297%

For TELECOMMUNICATION TECHNICIAN (728.5 hours) : 6.4606%

For CARPENTER (715.0 hours) : 0.73%

- Conversely, when I ranked trades based on the **percentage of hours of Journeymen** work that belong to African Americans for 2019, I found that these trades had the following total hours worked by African Americans:

For ASPHALT RAKER (42.857%) : 9.0 hours

For LABORER: HAZARDOUS WASTE/ASBESTOS REMOVER (37.342657%) :
267.0 hours

For FENCE & GUARD RAIL ERECTOR (28.082%) : 20.5 hours

For HVAC MECHANIC (25.5806%) : 385.5 hours

For TELEDATA WIREMAN/INSTALLER/TECHNICIAN (9.1575%) : 25.0 hours

- As you can see, both the percentages and totals are higher for Apprentices than for Journeymen. Of course, most are not the same trades, but even when comparing the same trades (the ones highlighted), you can see a higher hours and percentages (with the exception of LABORER) for Apprentices compared to Journeymen. Interestingly, the trades that have a large share of African Americans are ASPHALT RAKER and HAZARDOUS WASTE/ASBESTOS REMOVER, both jobs with occupational exposure that could affect the workers' health in the long run.

- When I ranked trades based on the **total number of Apprentice hours** worked by **Women** in 2019, I found that these trades had the following percentage of Women hours out of both genders:

For LABORER (3231.0) : 89.96%
 For ELECTRICIAN (1695.0) : 7.74%
 For CARPENTER (730.5) : 8.548%
 For PIPEFITTER (256.0) : 6.2023%
 For IRONWORKER (104.25) : 38.72%

- Conversely, when I ranked trades based on the **percentage of hours of Apprentice work** that belong to **Women** for 2019, I found that these trades had the following total hours worked by Women:

For LABORER (89.96%) : 3231.0 hours
 For IRONWORKER (38.71866%) : 104.25 hours
 For BRICK/PLASTER/CEMENT MASON (33.33%) : 16.0 hours
 For CEMENT MASONRY/PLASTERING (23.8095%) : 10.0 hours
 For MARBLE & TILE FINISHERS (15.7068%) : 15.0 hours

- When I ranked trades based on the **total number of Journeymen hours** worked by African Americans in 2019, I found that these trades had the following percentage of African American hours out of all the ethnicities:

For LABORER (3351.44 hours) : 3.685%
 For IRONWORKER/WELDER (765.25 hours) : 5.249%
 For MARBLE MASONS, TILELAYERS & TERRAZZO MECH (488.0 hours) : 4.85%
 For ELEVATOR CONSTRUCTOR HELPER (474.5 hours) : 8.86%
 For LABORER: MULTI-TRADE TENDER (464.0) : 47.6876%

- Conversely, when I ranked trades based on the **percentage of hours of Apprentices work** that belong to **Women** for 2019, I found that these trades had the following total hours worked by Women:

For ROLLER/SPREADER/MULCHING MACHINE (100.0%) : 64.0 hours
 For LABORER: MULTI-TRADE TENDER (47.68756%) : 464.0 hours
 For EQUIPMENT OPERATOR (Class B CDL) (41.376%) : 445.0 hours

For HVAC MECHANIC (**16.1579%**) : **243.5 hours**

For ASBESTOS REMOVER - PIPE / MECH. EQUIPT (**9.8734 %**) : **39.0 hours**

- Something of note is that the percentages and even some of the total hours are higher for Women than for African Americans. For example, if we look at one trade that they have in common, LABORER, we see that not only do Women Apprentices and Women Journeyman have more hours, they also account for a larger share of the workforce than African Americans. In fact, Women make up almost the entire labor force for LABORER. Additionally, we don't see any trade being dominated by African Americans, or with such a large percentage of African Americans, as we see with Women.

Conclusions

Three important findings stand out showing the differences in the amount of journeymen work, work year-over-year given the pandemic is one of the years, and the kind of work with the highest percentages given to ethnicities.

There appears to be a bias towards hiring caucasian men for journeymen work on MA construction contracts. The evidence shows that ethnicities and women get lower paying apprentice work but not the better and more long-term journeymen positions.

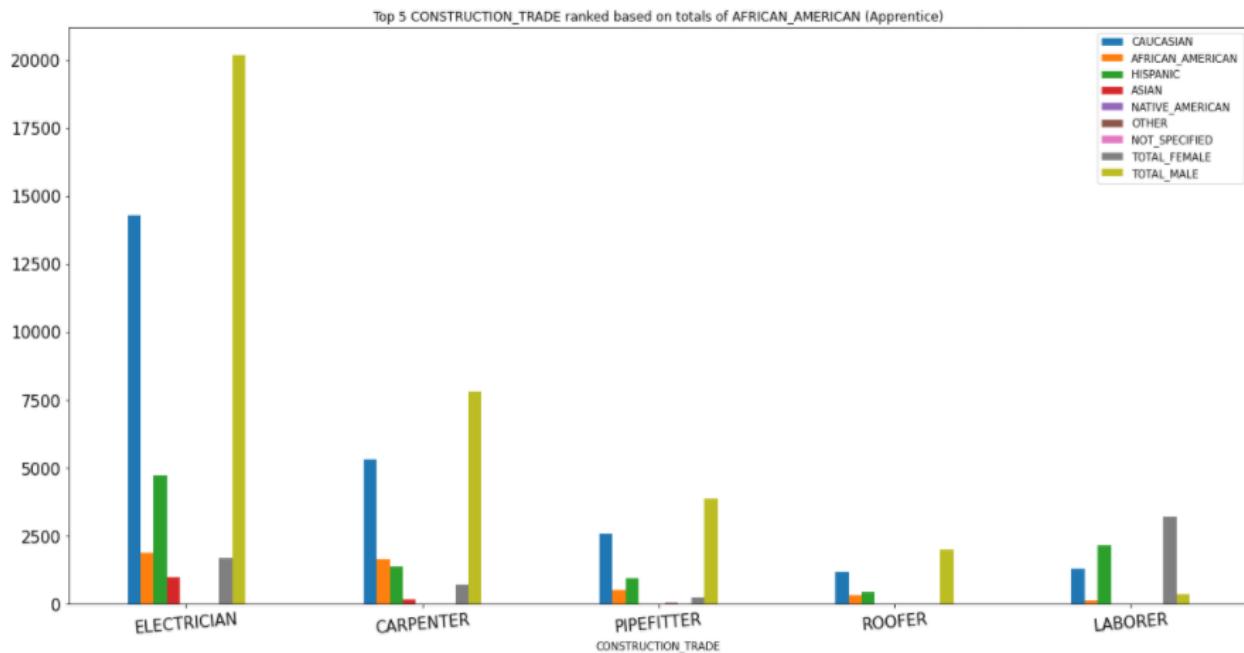
With respect to the pandemic's impact on the workforce, there is a drop in utilization hours in 2020 than 2019, and not as low as one might expect. Surprisingly, the ethnic makeup of the workforce is stable between the years. The consistency in very different economic times shows a tenacity for favoring caucasian males in construction work in the state.

More concerning is the nasty kind of work that pops to the top of the graph when the population is not split out by apprentice and journeymen: **Top 5 trades based on percentage of African Americans, shown in screenshots.** Asphalt Raker and Laborer Hazardous Waste are the top job percentages for African-Americans. For both African-Americans and Hispanics, Laborer Hazardous Waste work is more often tasked to them than caucasions.

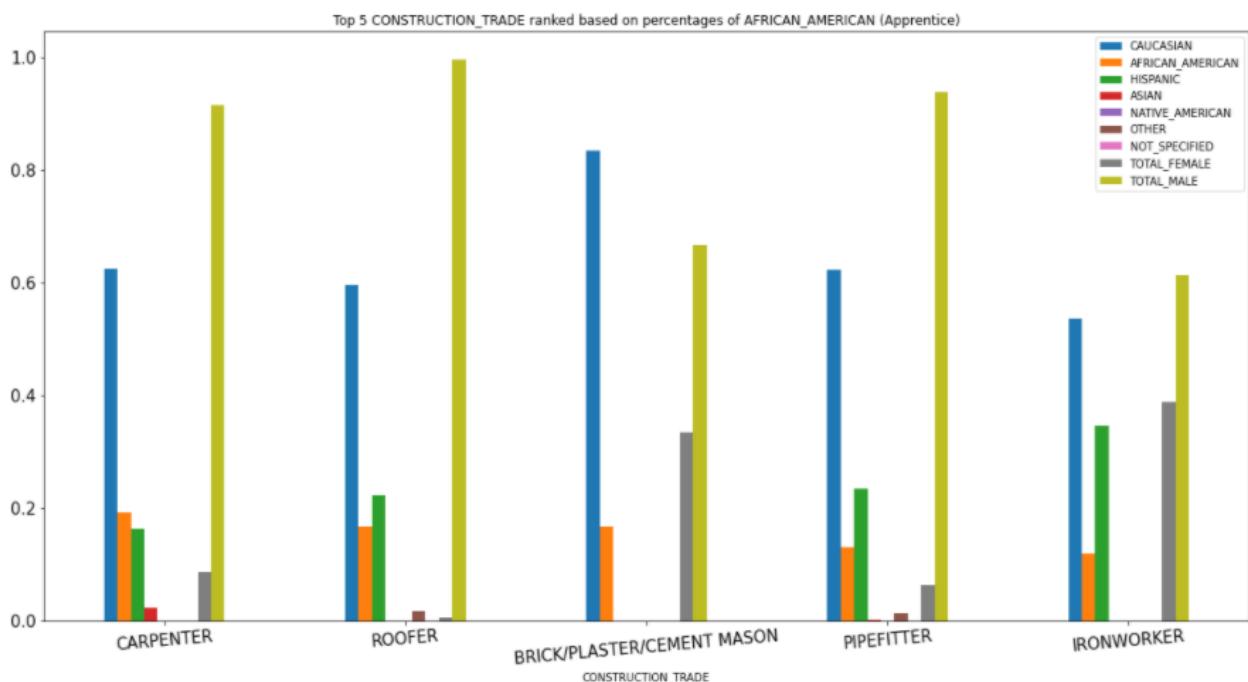
Limitations and Risks

Due to the past two years being very unusual, there is a risk in making assertions about the behavior of the state related to ethnicities and women. With more than 5 years of data, the information will have greater import. Still the data is hard evidence of journeymen work being greatly reduced for the ethnicities and women compared to caucasian men. The kinds of work appear to be more hazardous. The risk involved in making claims of racial or gender bias in the current social climate is that someone is bound to get upset with talking about the subject even if the numbers are accurate. Therefore for the final, we are collecting data published in other places on the internet about the distribution of people of various ethnicities that live in the state (from the census) and that perform construction work (from the department of labor statistics). These outside sources provide a backdrop from which to compare and contrast why workers who are present in the state's workforce are not working on state contracts.

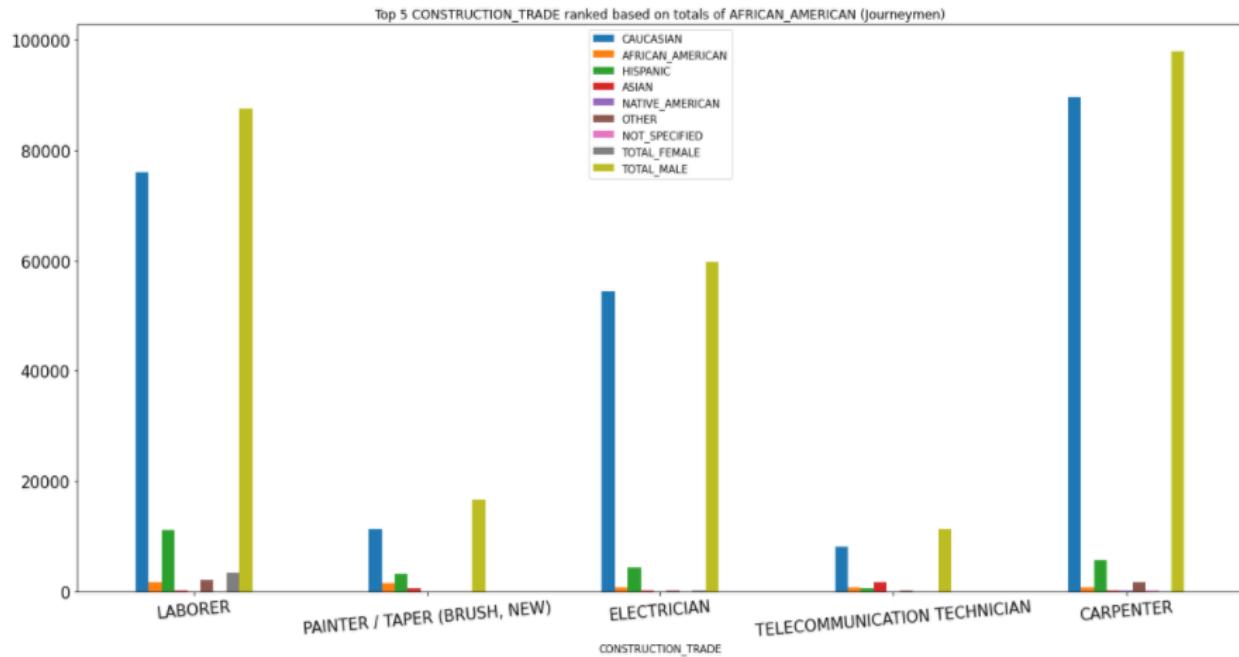
SCREENSHOTS



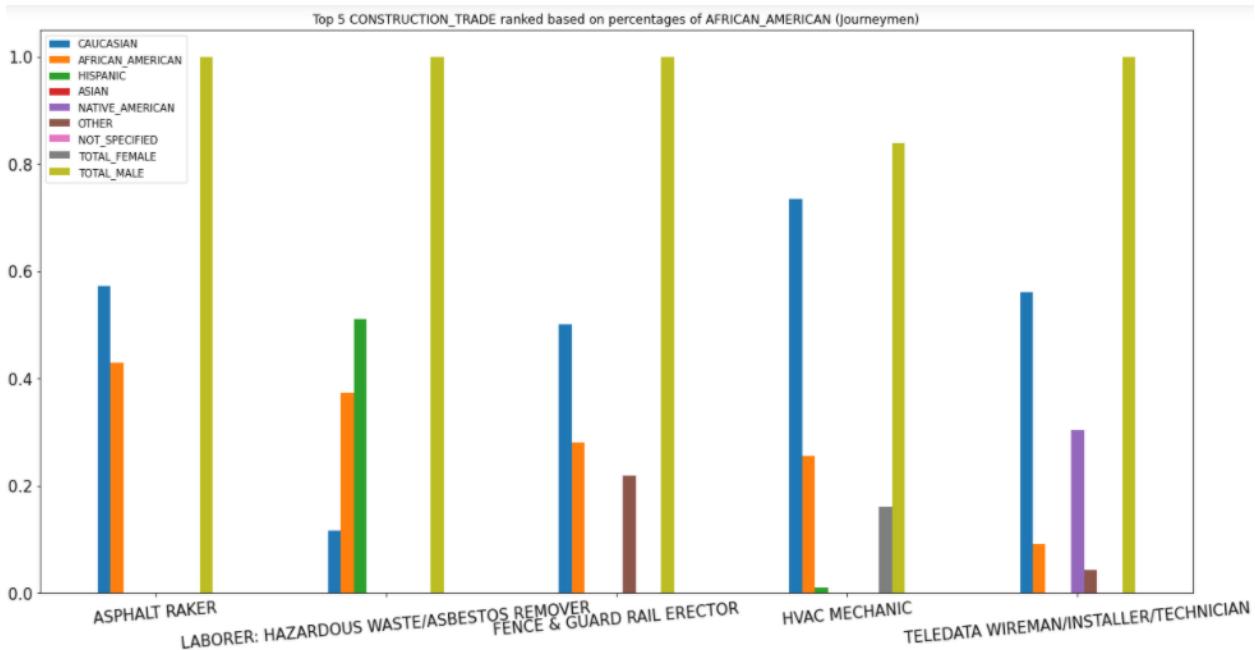
Top 5 trades based on the total number of Apprentice hours worked by African Americans in 2019.



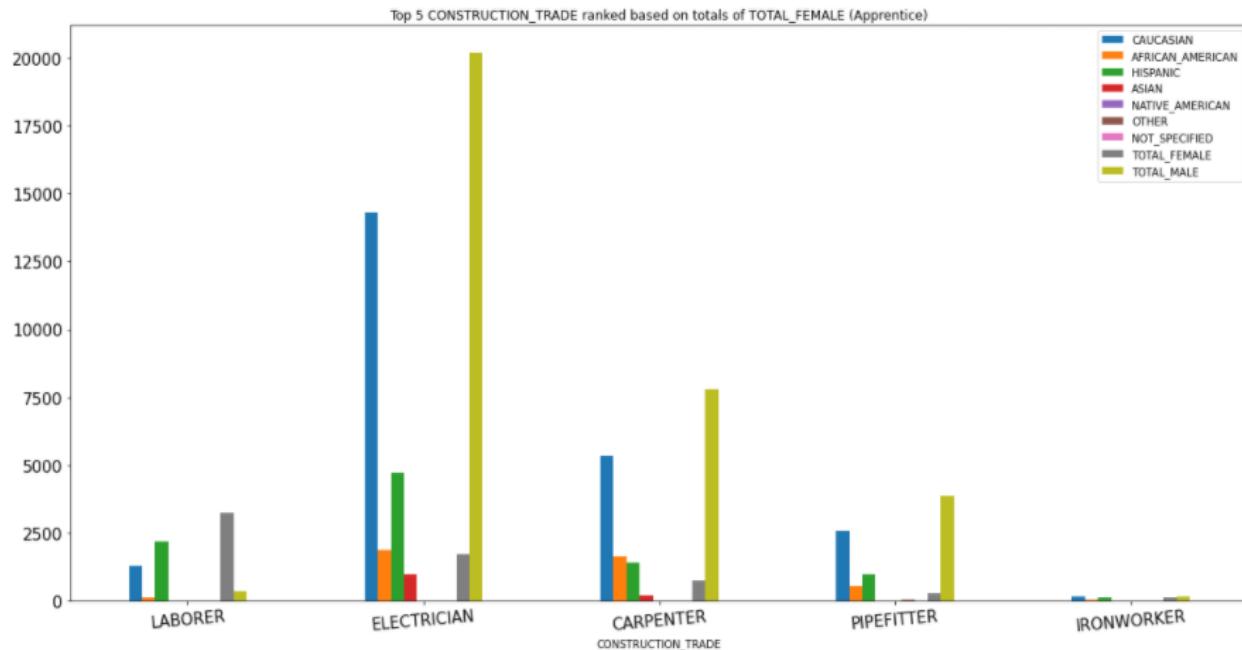
Top 5 trades based on the percentage of hours of Apprentice work that belong to African Americans for 2019



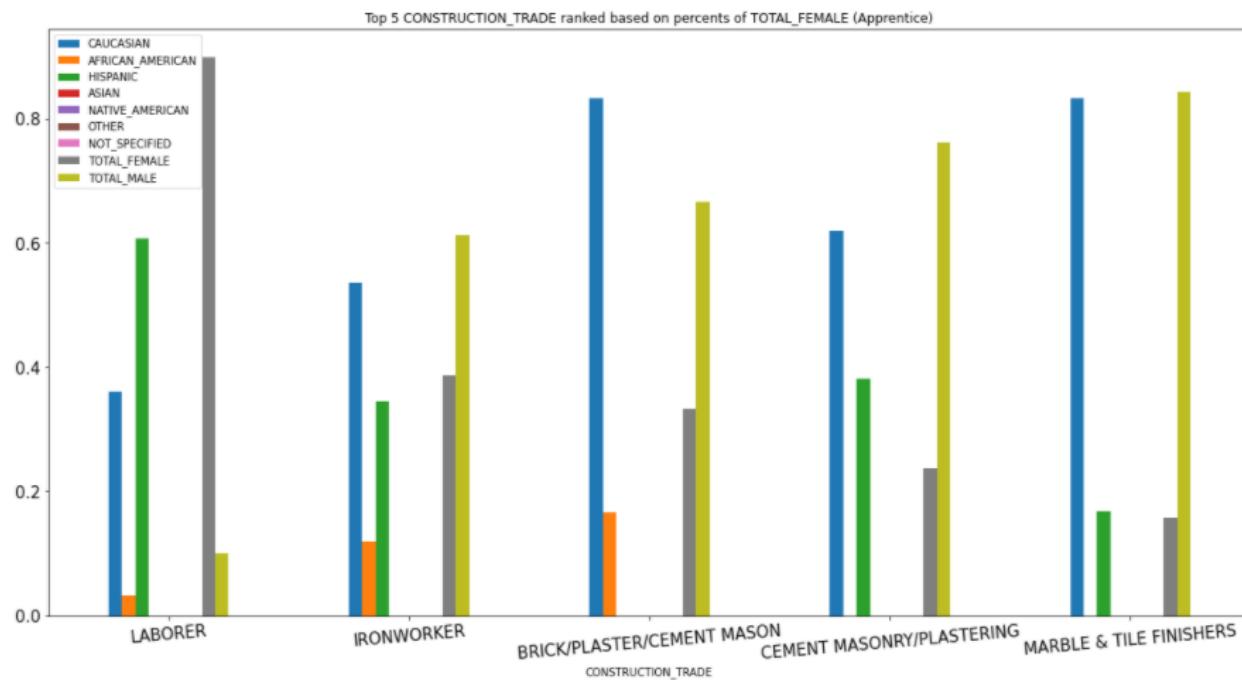
Top 5 trades based on the total number of Journeymen hours worked by African Americans in 2019



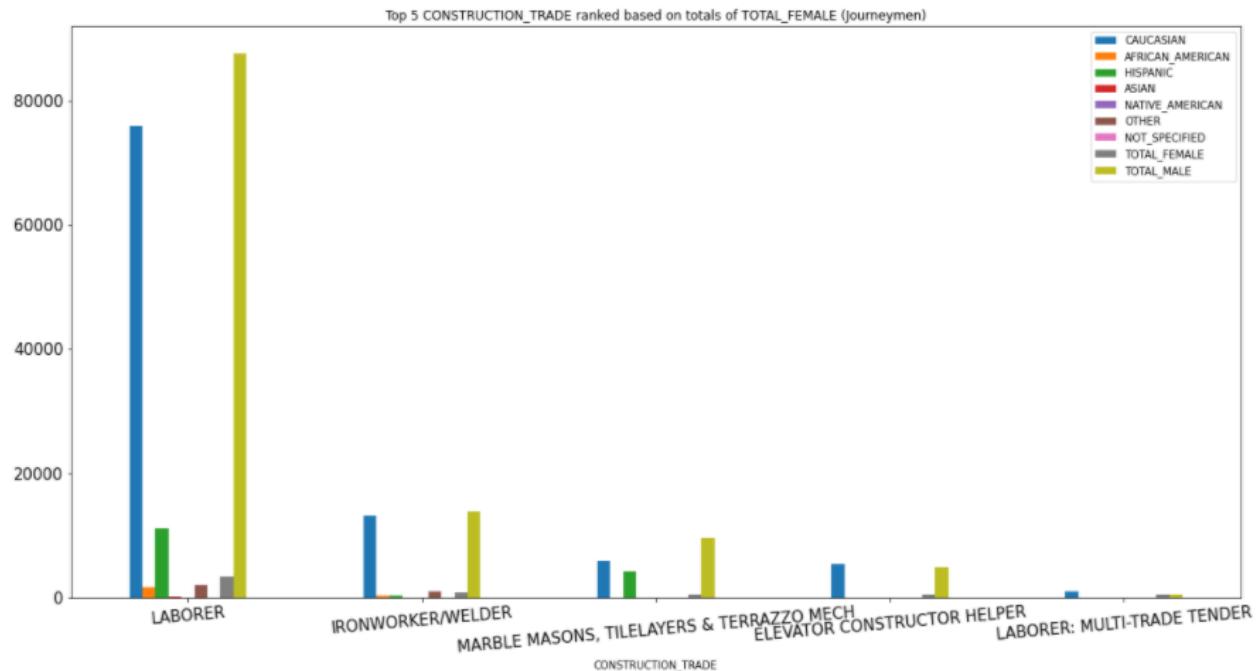
Top 5 trades based on the percentage of hours of Journeymen work that belong to African Americans for 2019



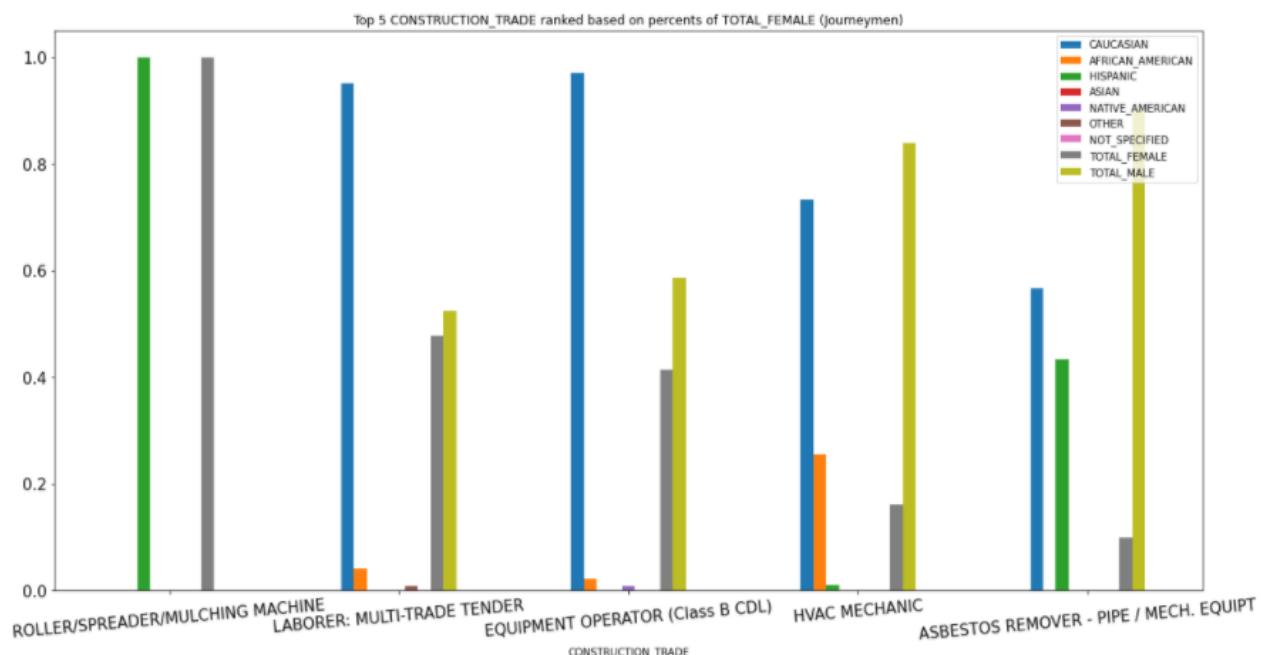
Top 5 trades based on the total number of hours completed by Women Apprentice in 2019



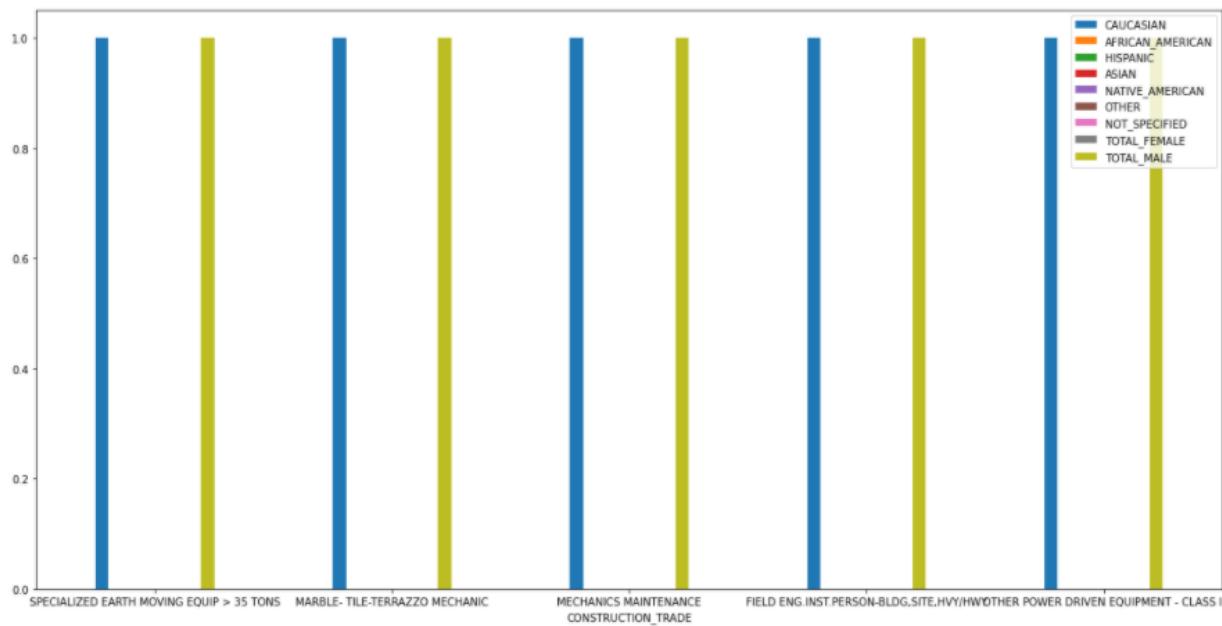
Top 5 trades based on the percentage of hours Women Apprentices worked in 2019



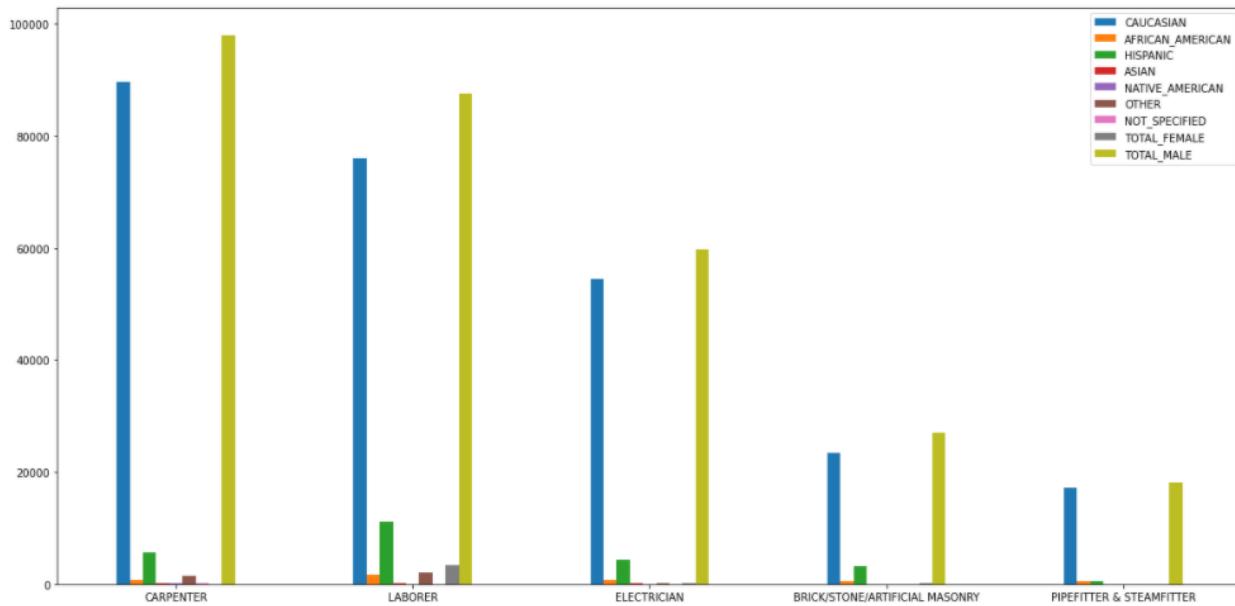
Top 5 trades based on the total number of hours completed by Women Journeymen in 2019



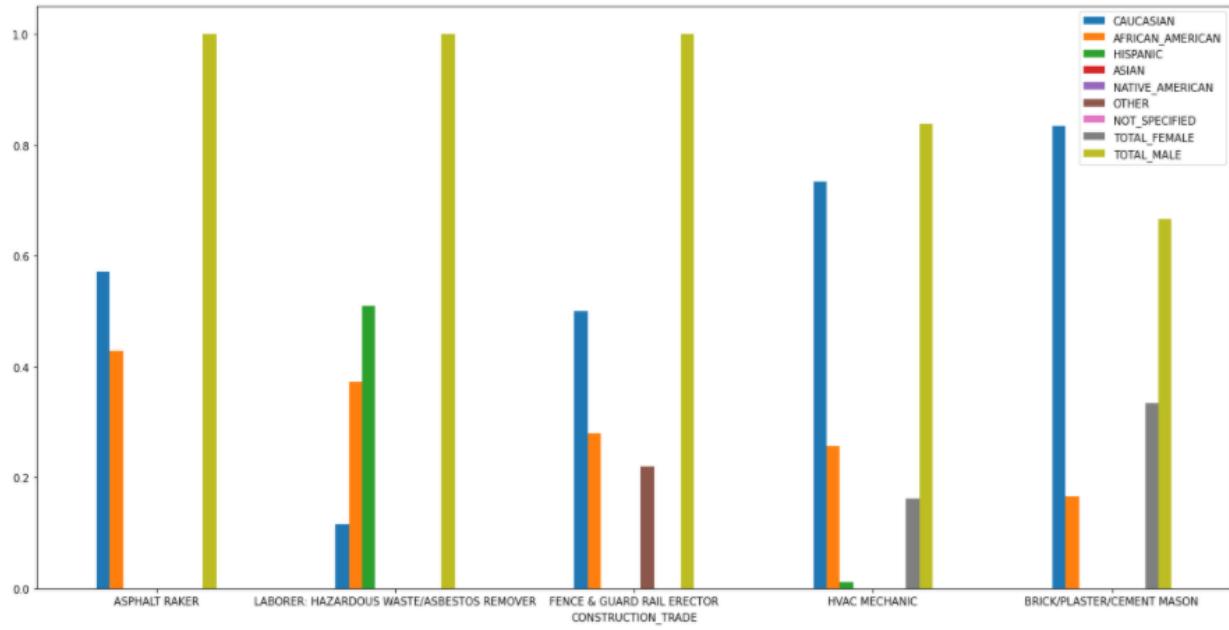
Top 5 trades based on the percentage of hours worked by Women Journeymen in 2019



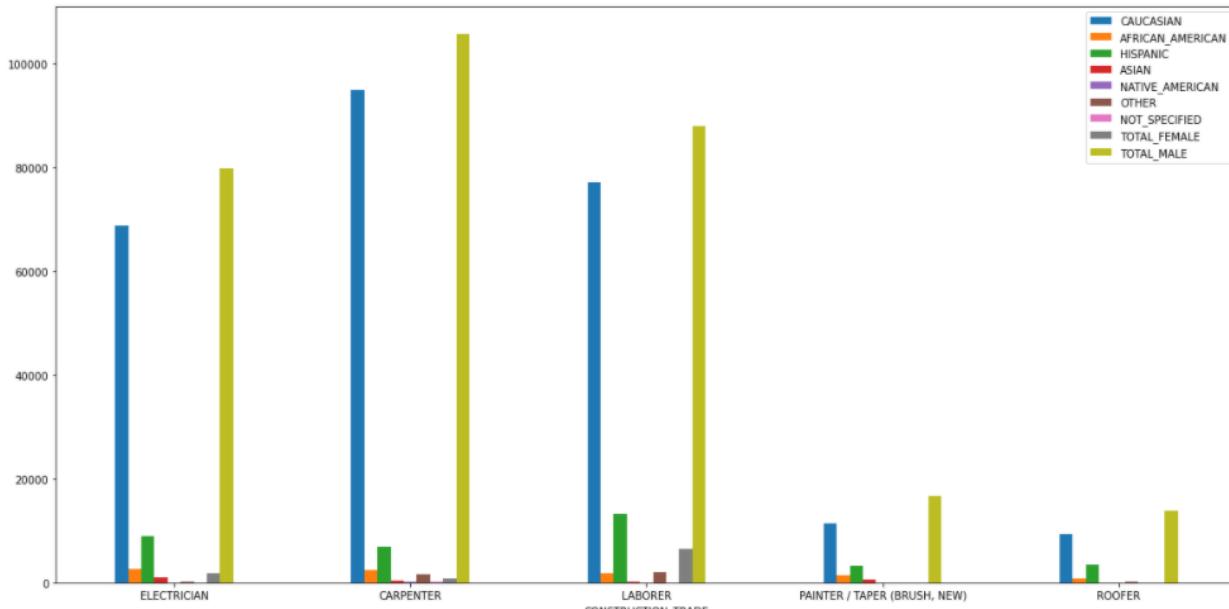
Top 5 trades ranked based on percentage of Caucasians (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



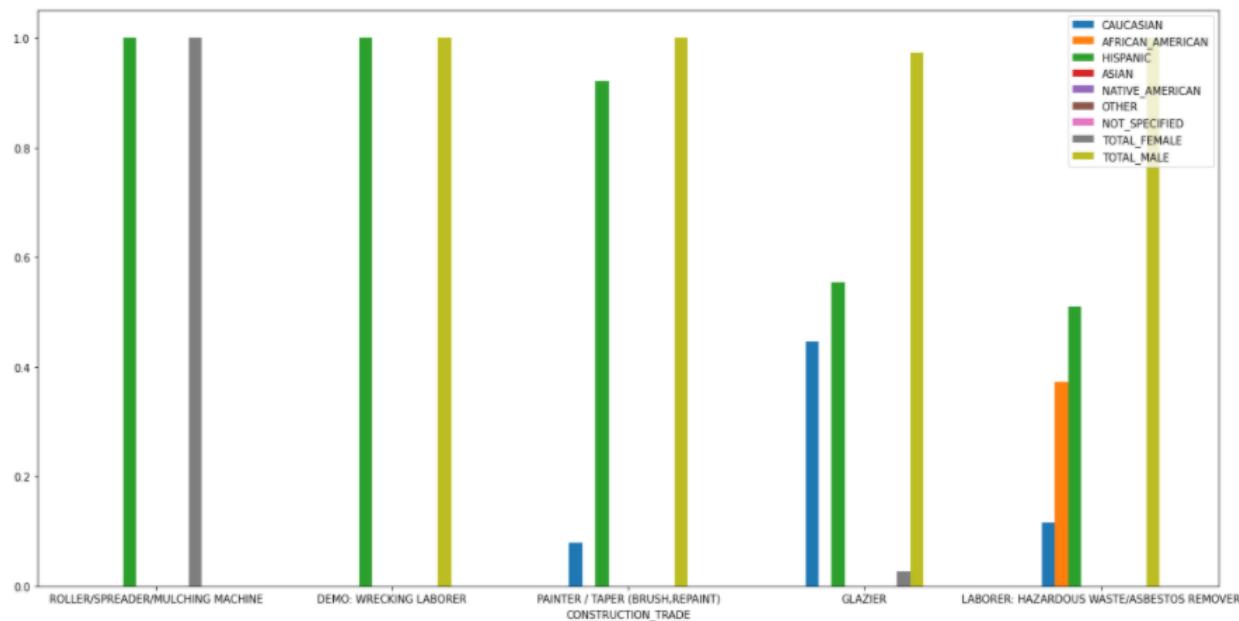
Top 5 trades ranked based on total hours of Caucasians (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



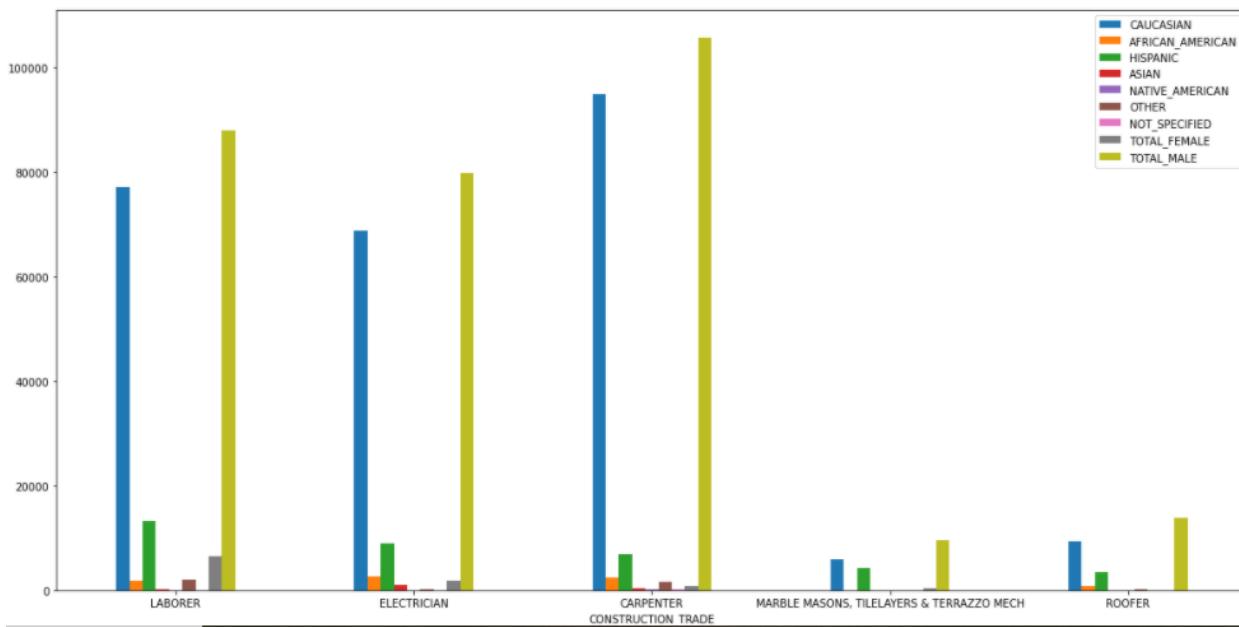
Top 5 trades based on percentage of African Americans (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



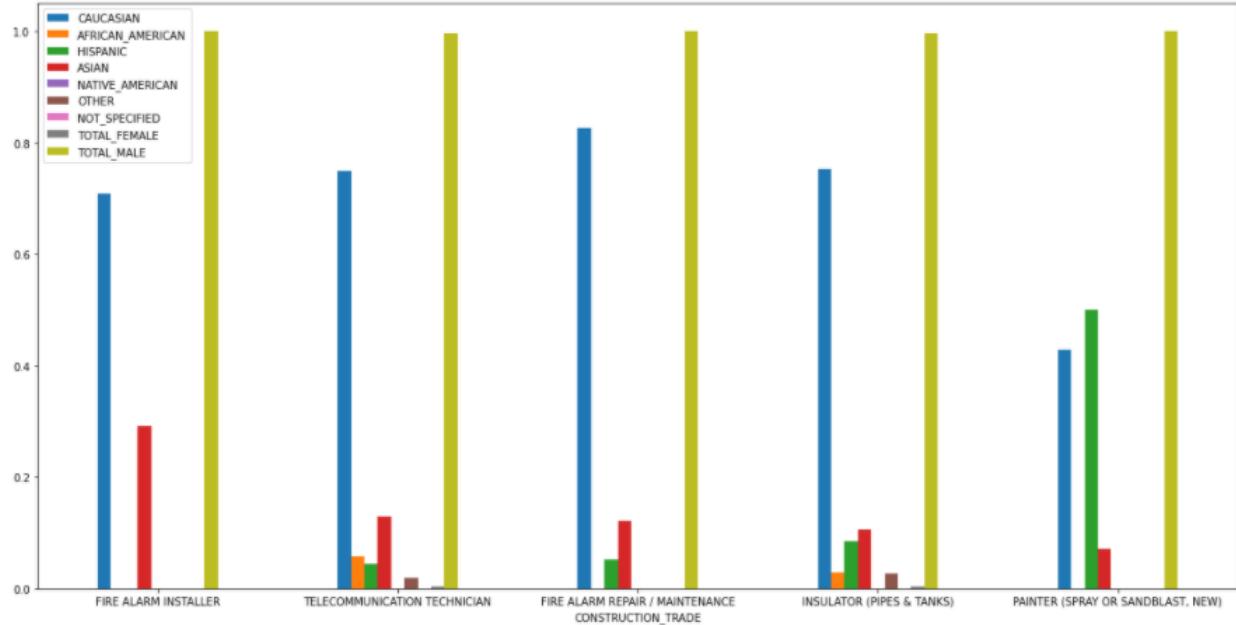
Top 5 trades based on total hours of African Americans (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



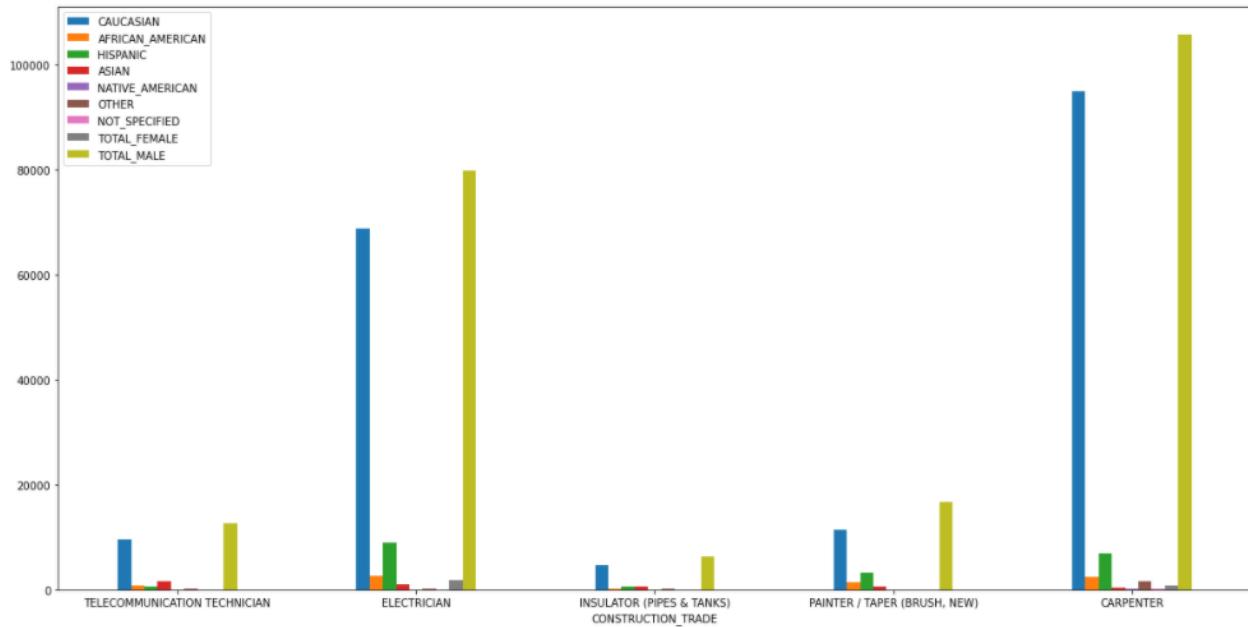
Top 5 trades based on percentage of Hispanics (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



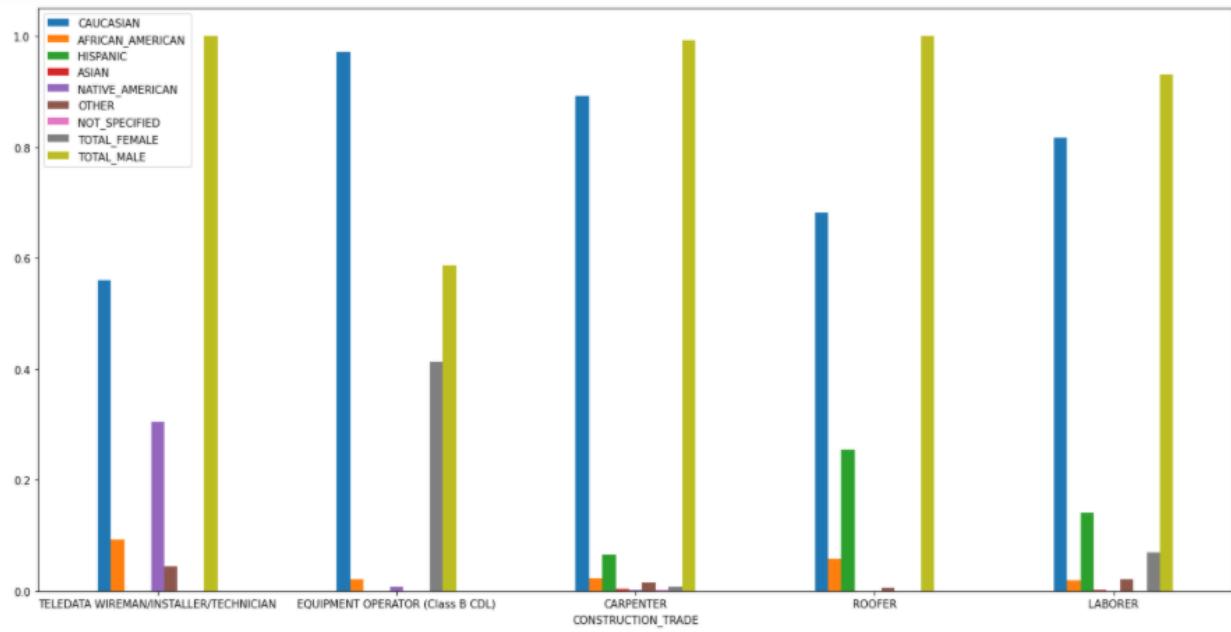
Top 5 trades based on total hours of Hispanics (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



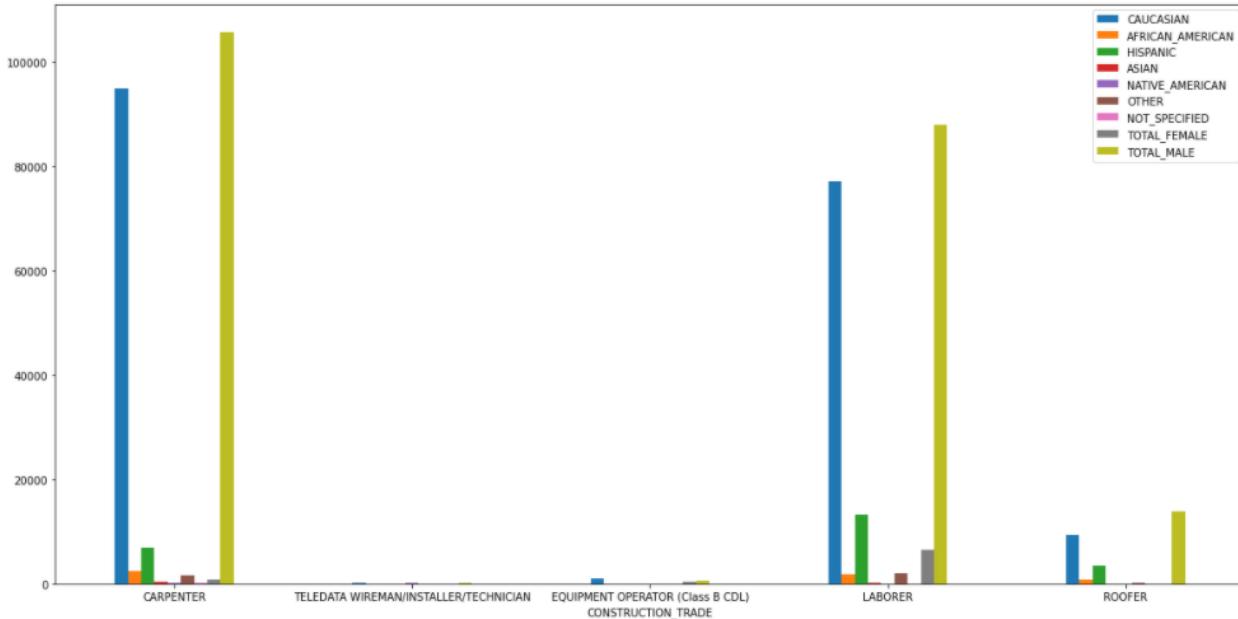
Top 5 trades based on percentage of Asians (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



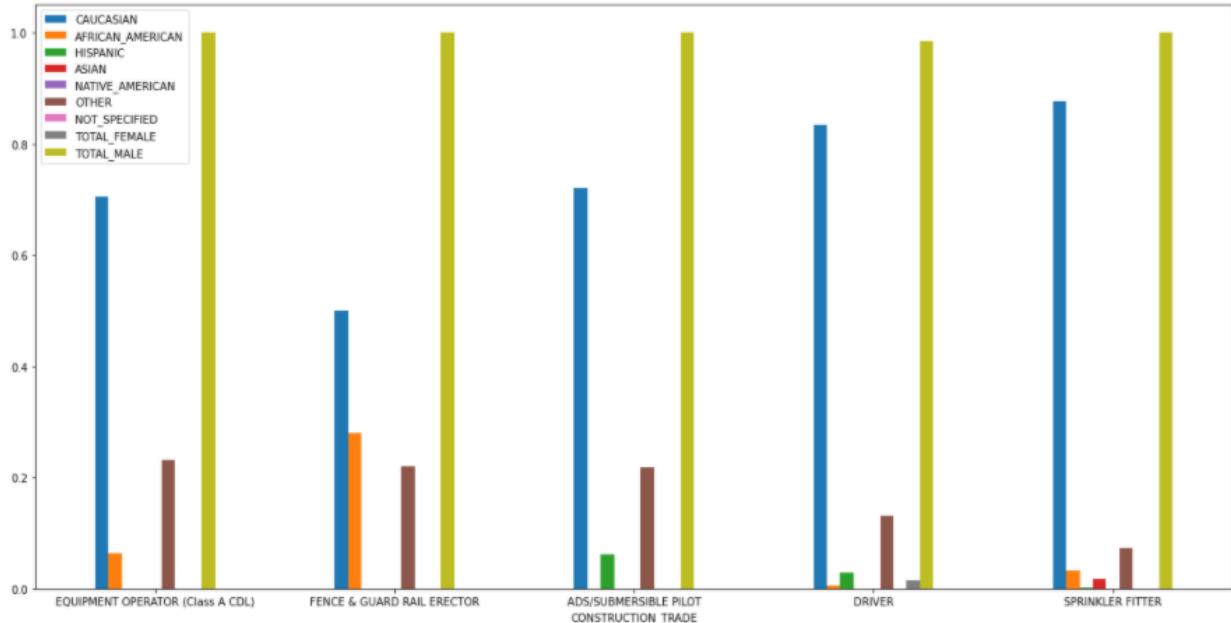
Top 5 trades based on total hours of Asians (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



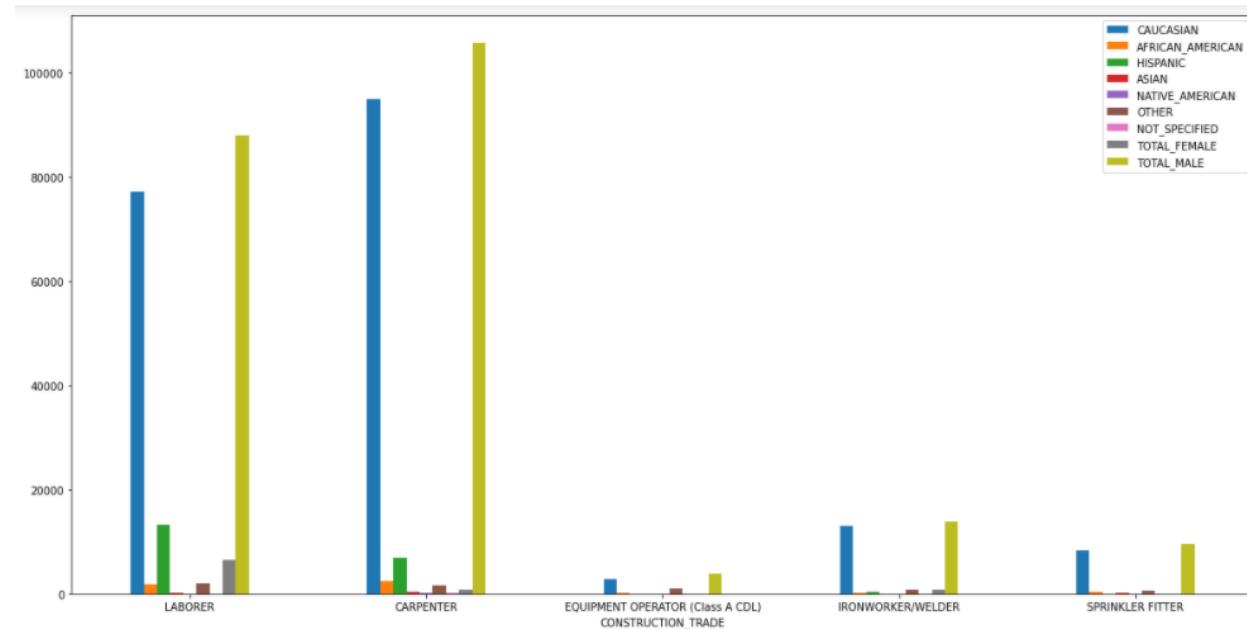
Top 5 trades based on percentage of Native Americans (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



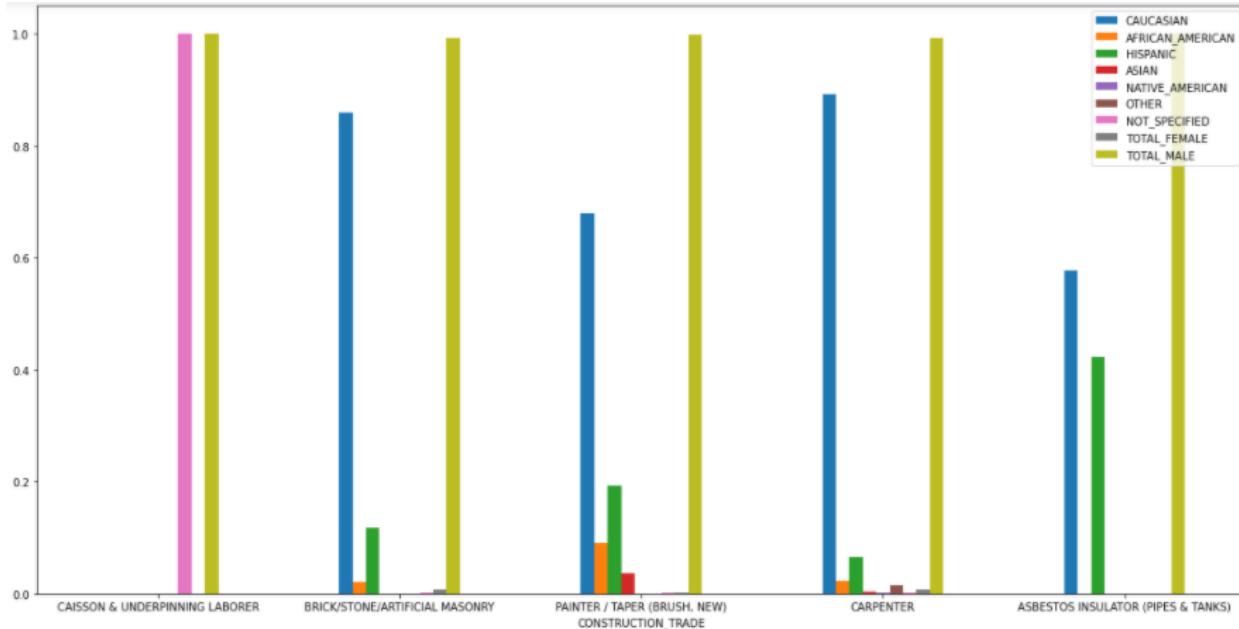
Top 5 trades based on total hours of Native Americans (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



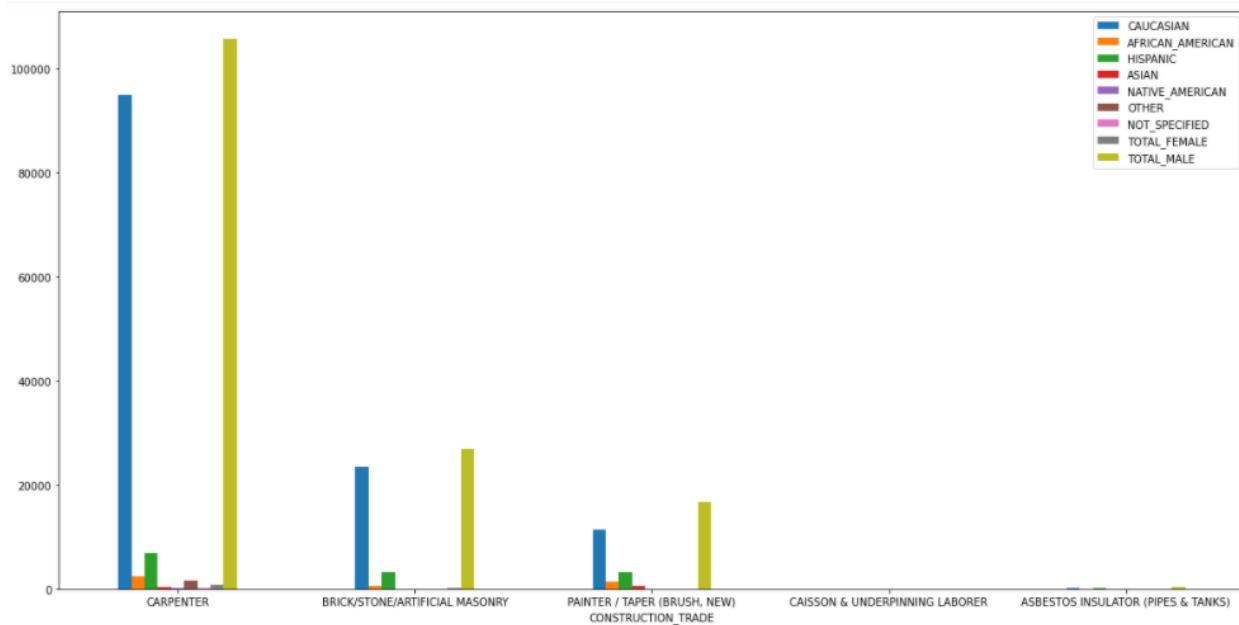
Top 5 trades based on percentage of Other (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



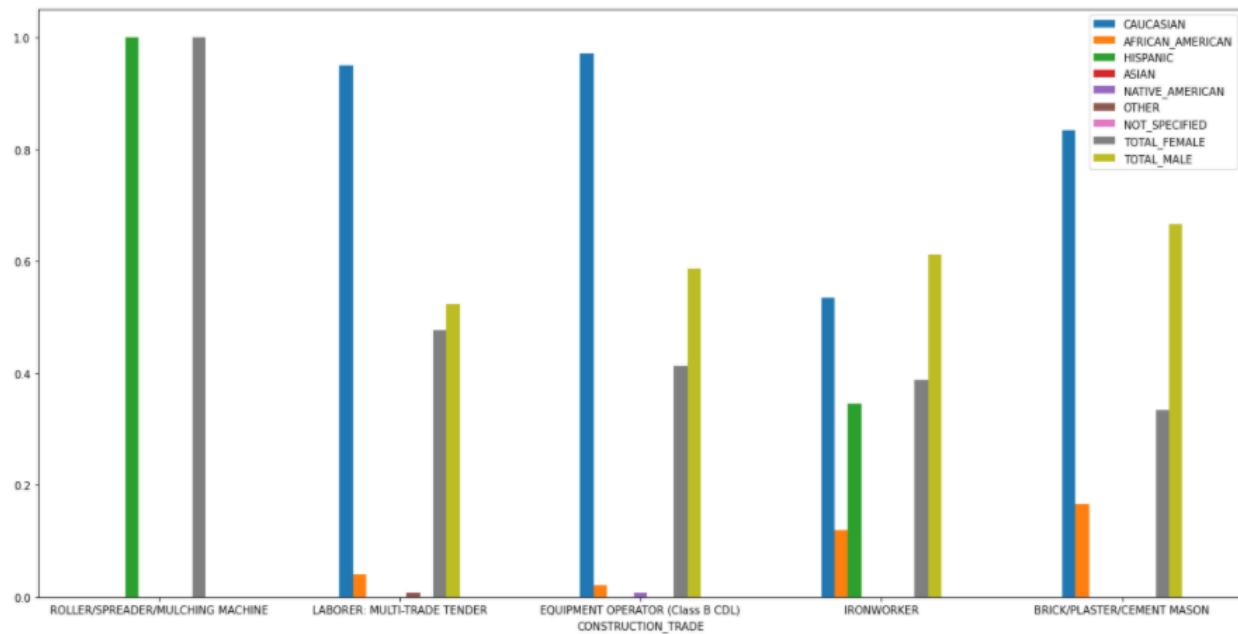
Top 5 trades based on total hours of Other (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



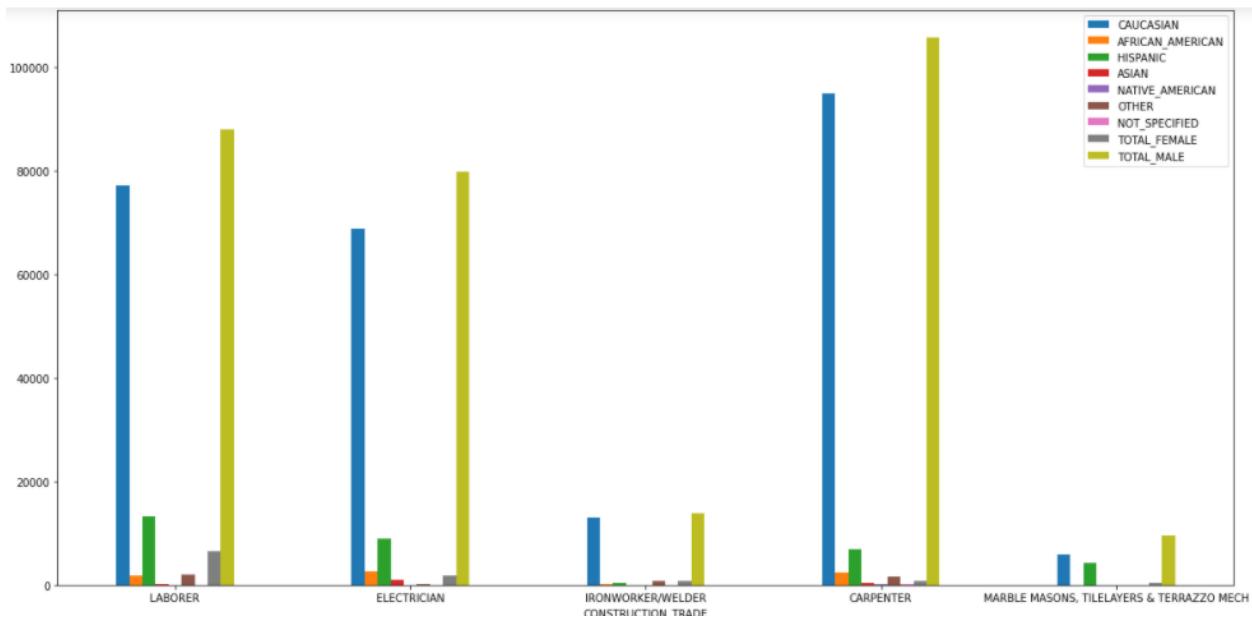
Top 5 trades based on percentage of Not Specified (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



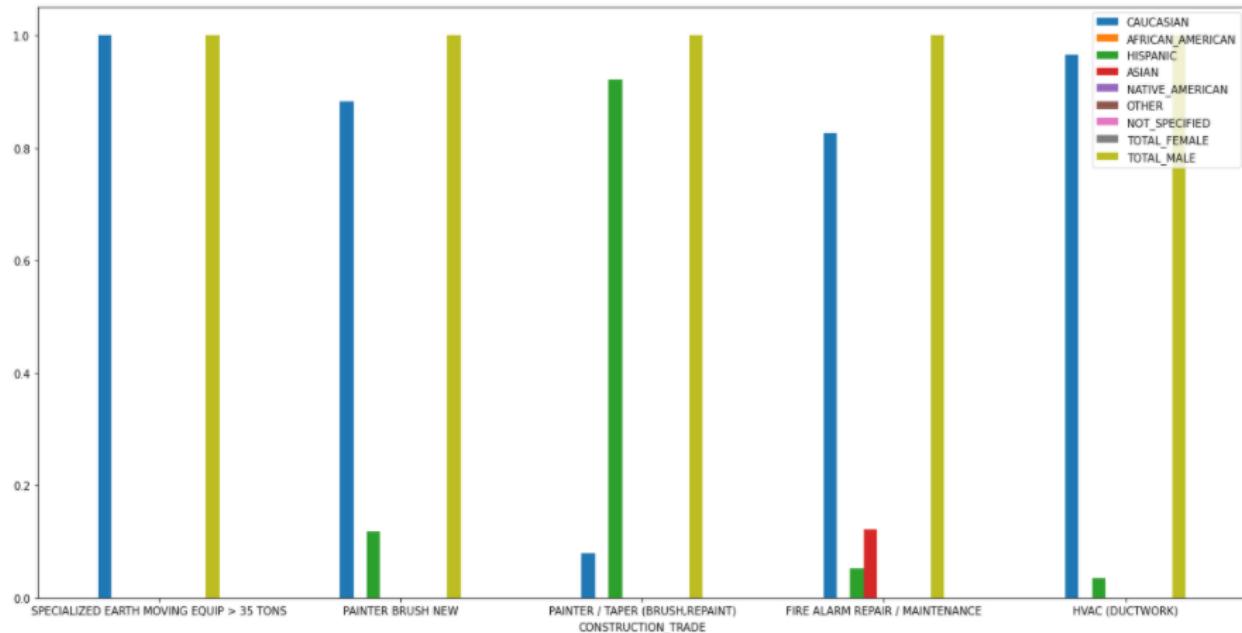
Top 5 trades based on total hours of Not Specified (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



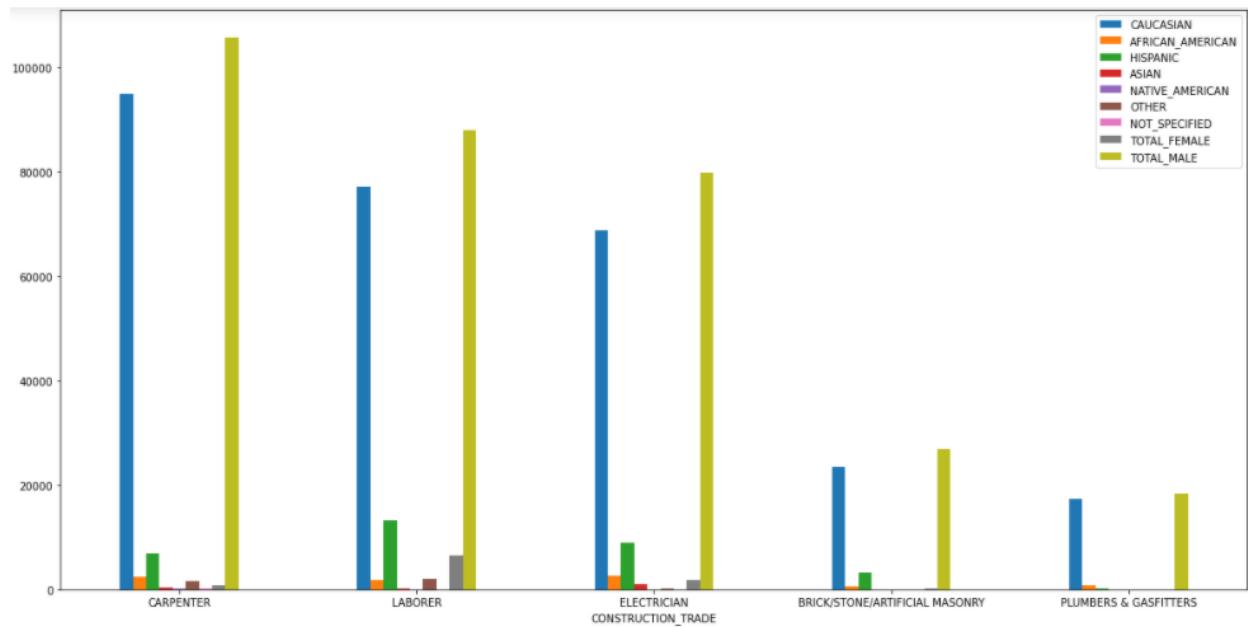
Top 5 trades based on percentage of Females (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



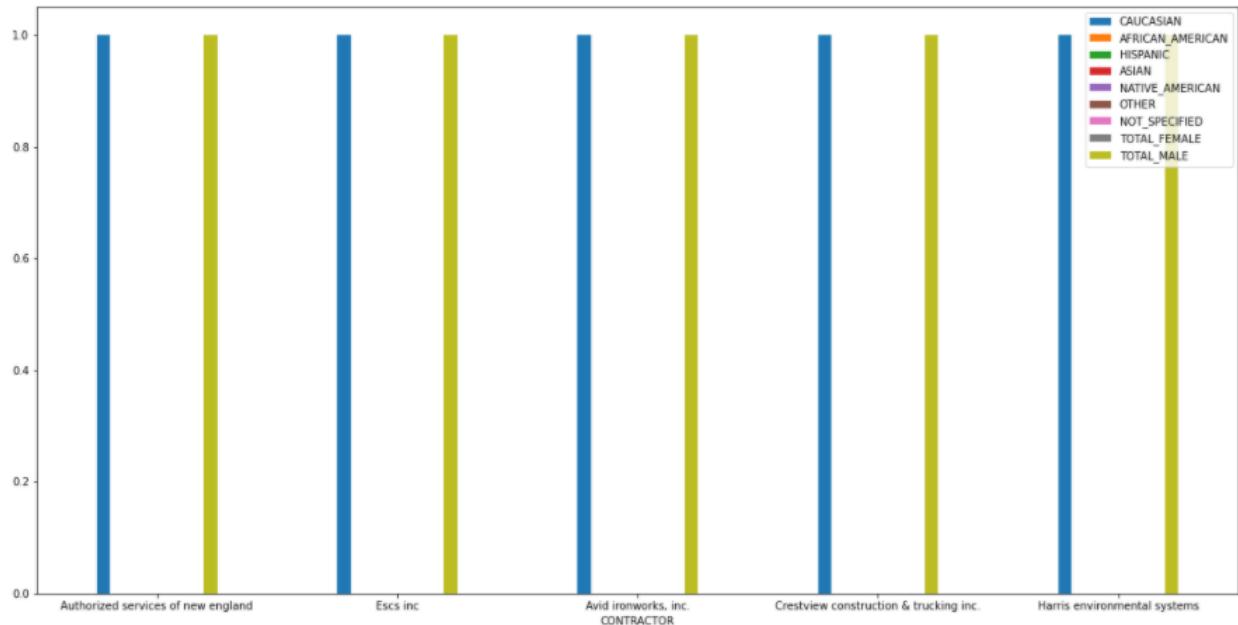
Top 5 trades based on total hours of Females (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



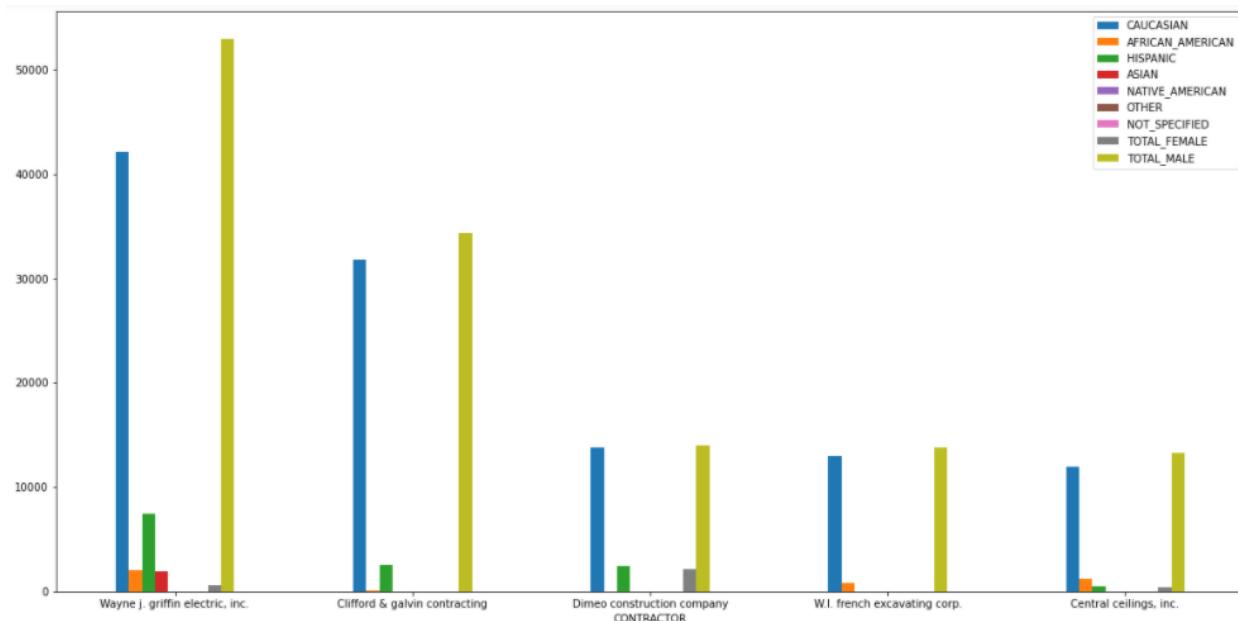
Top 5 trades based on percentage of Males (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



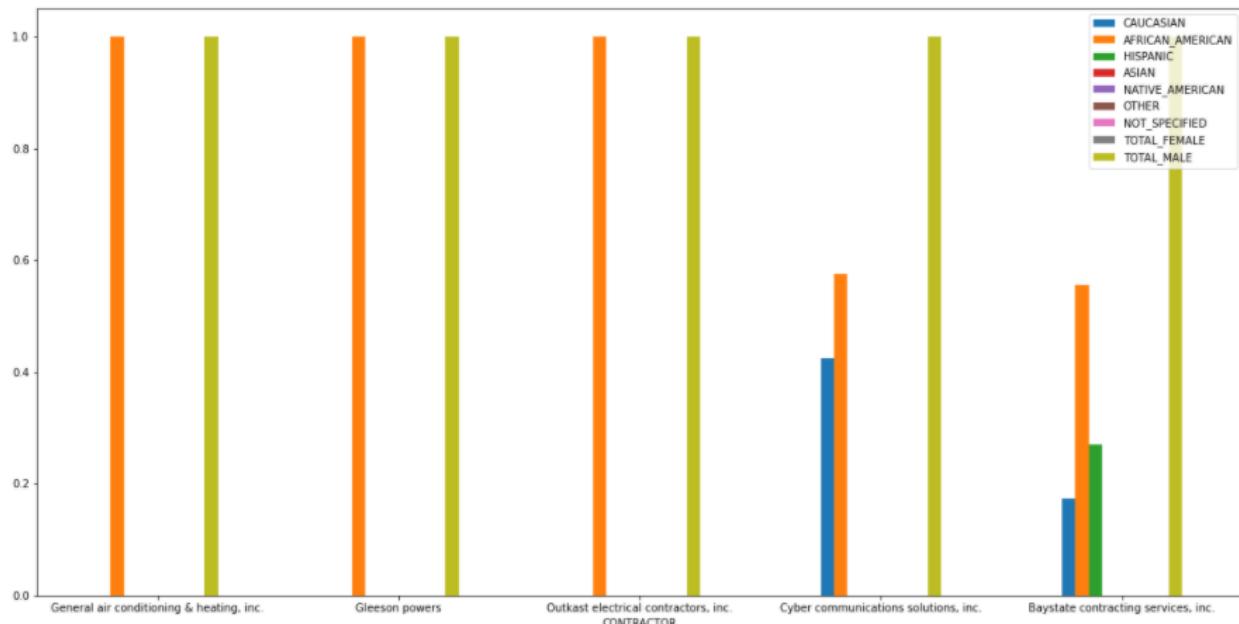
Top 5 trades based on total hours of Males (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



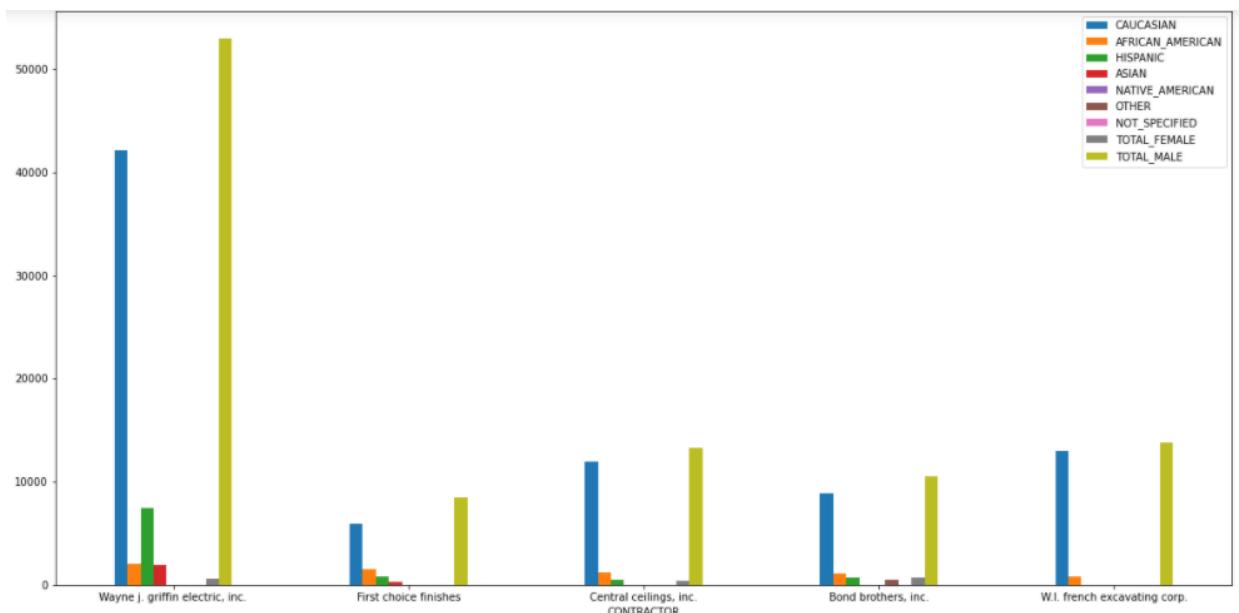
Top 5 companies based on percentage of Caucasians (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



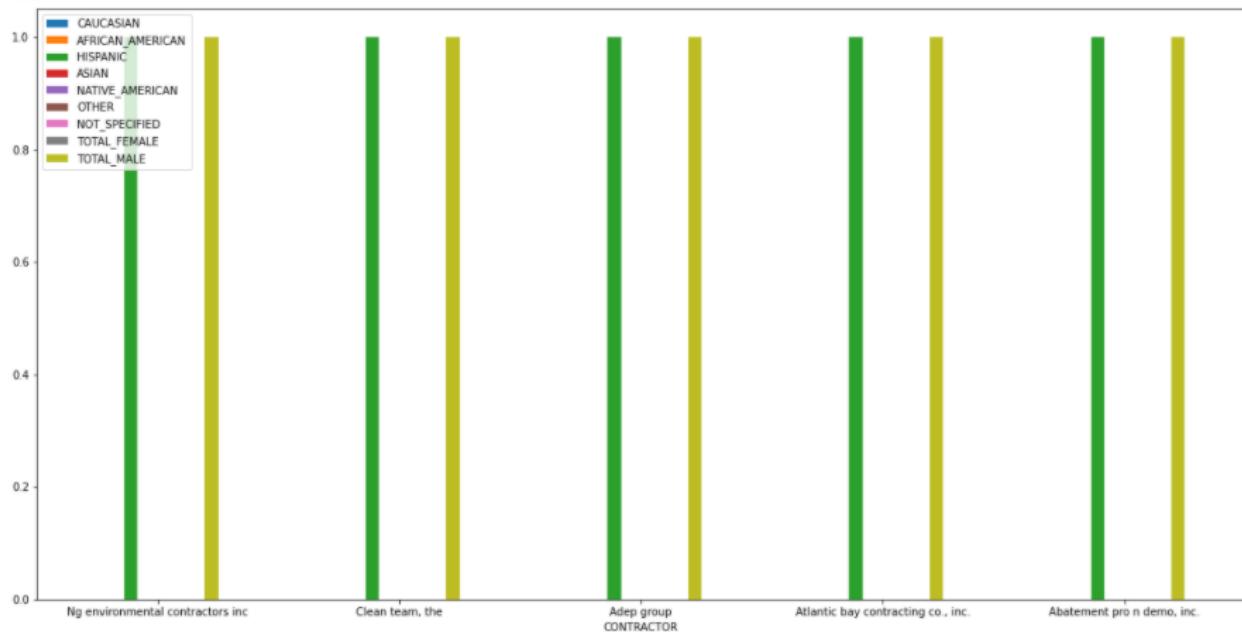
Top 5 companies based on total hours of Caucasians (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



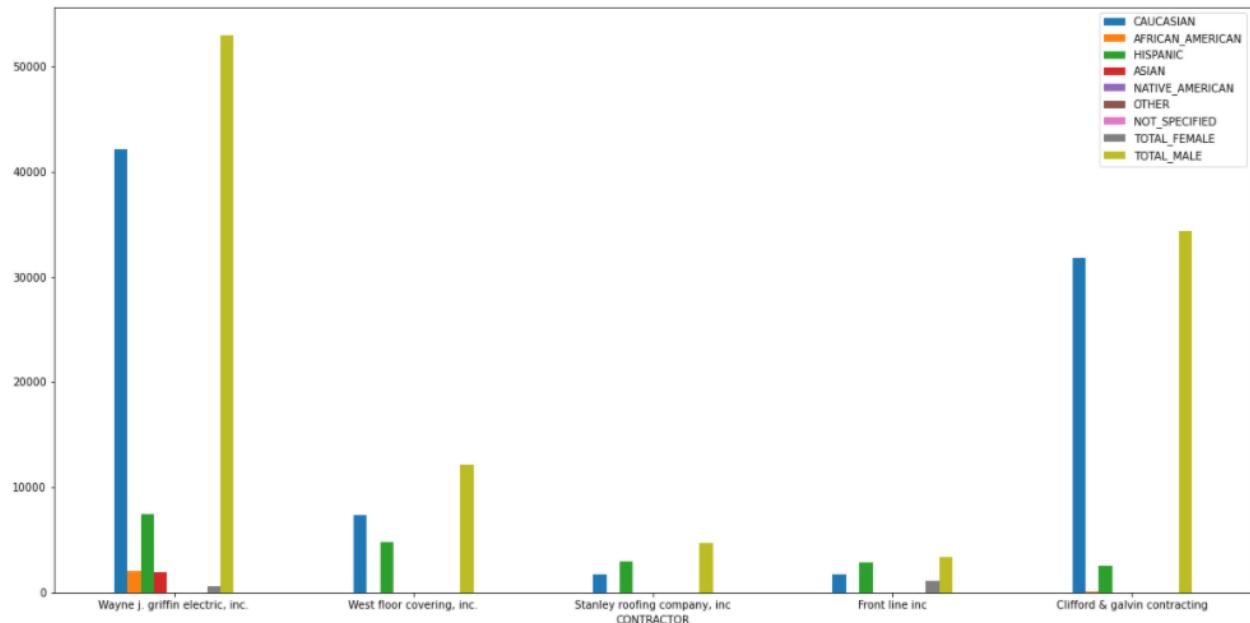
Top 5 companies based on percentage of African Americans (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



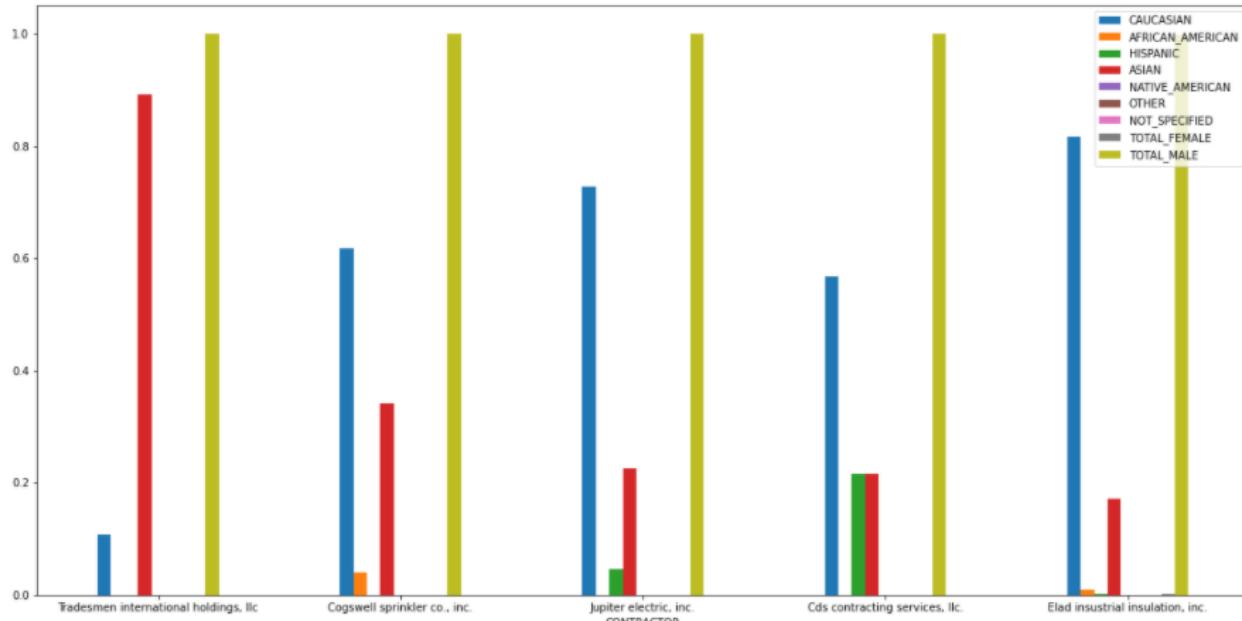
Top 5 companies based on total hours of African Americans (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



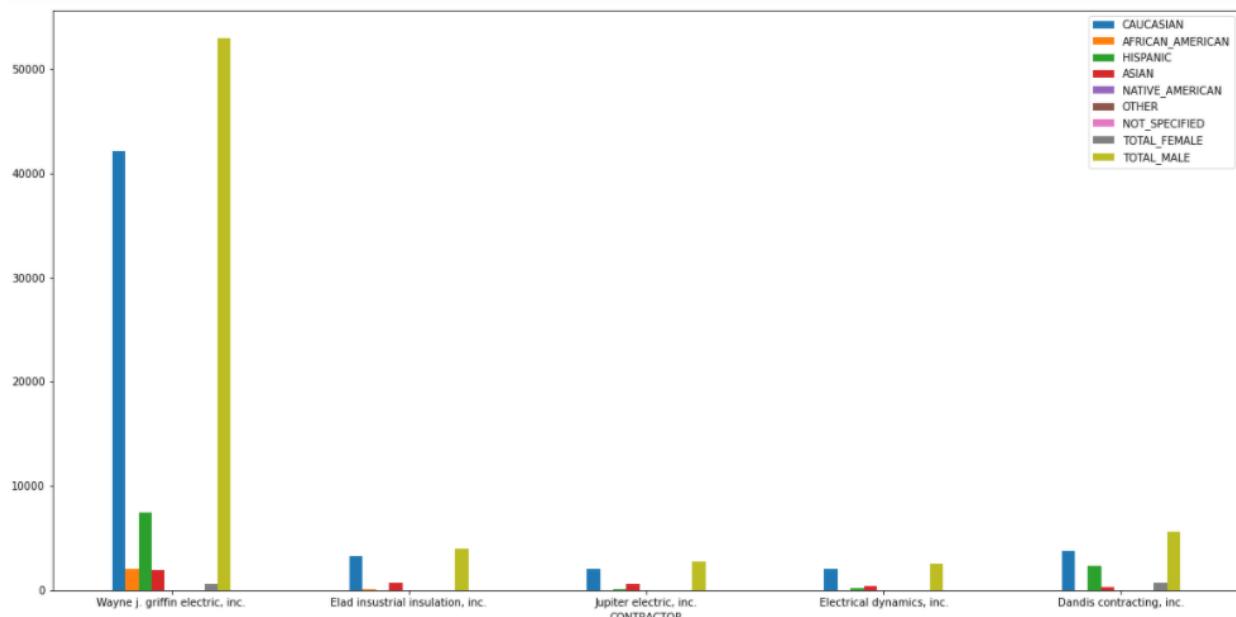
Top 5 companies based on percentage of Hispanics (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



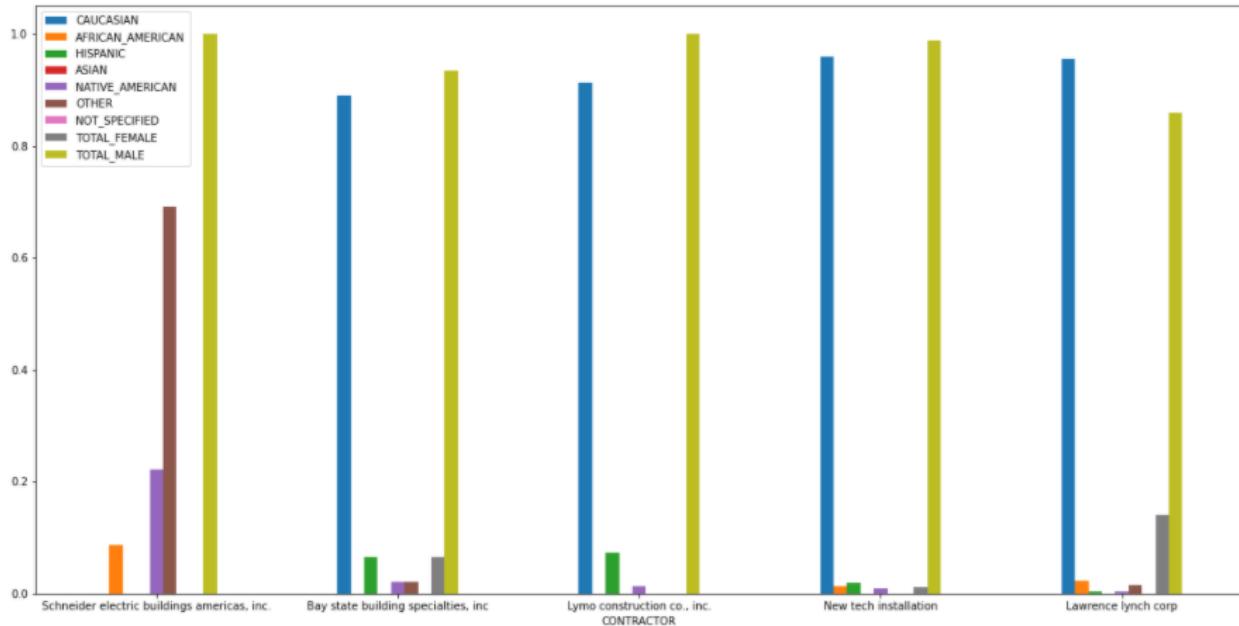
Top 5 companies based on total hours of Hispanics (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



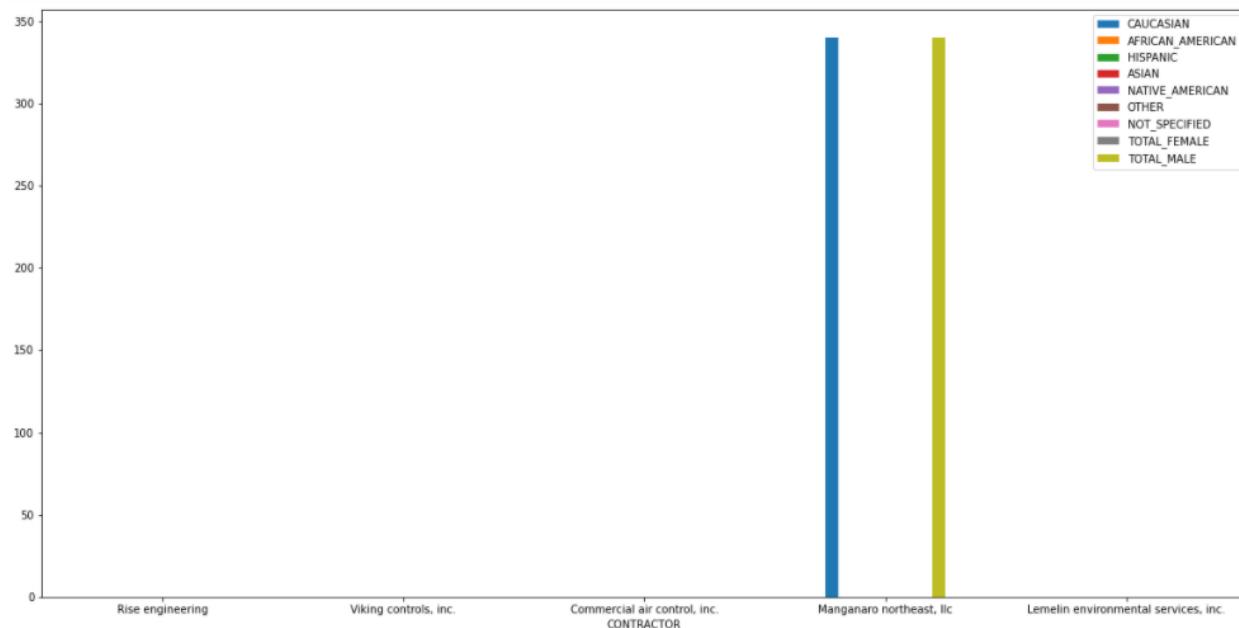
Top 5 companies based on percentage of Asians (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



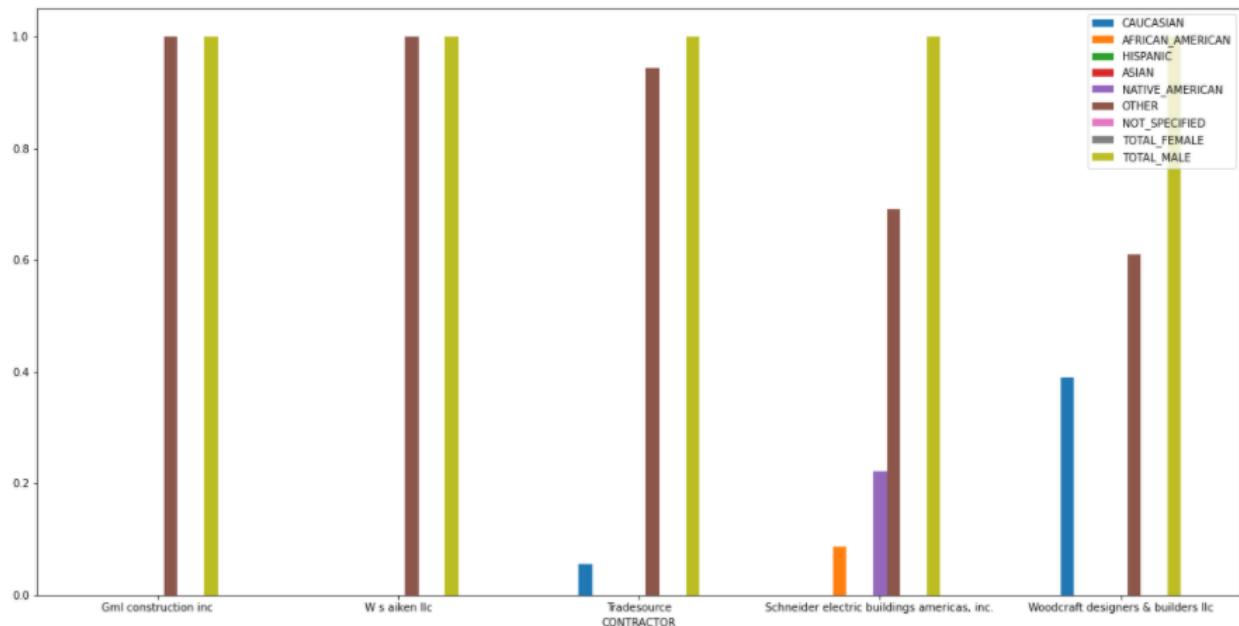
Top 5 companies based on total hours of Asians (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



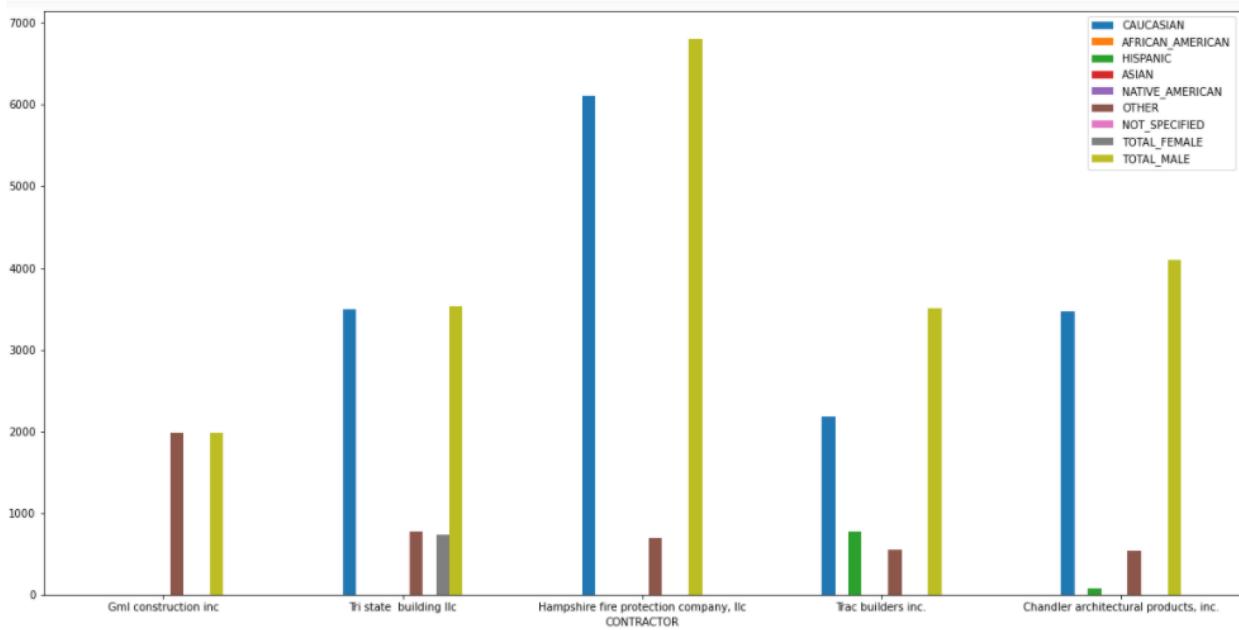
Top 5 companies based on percentage of Native Americans (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



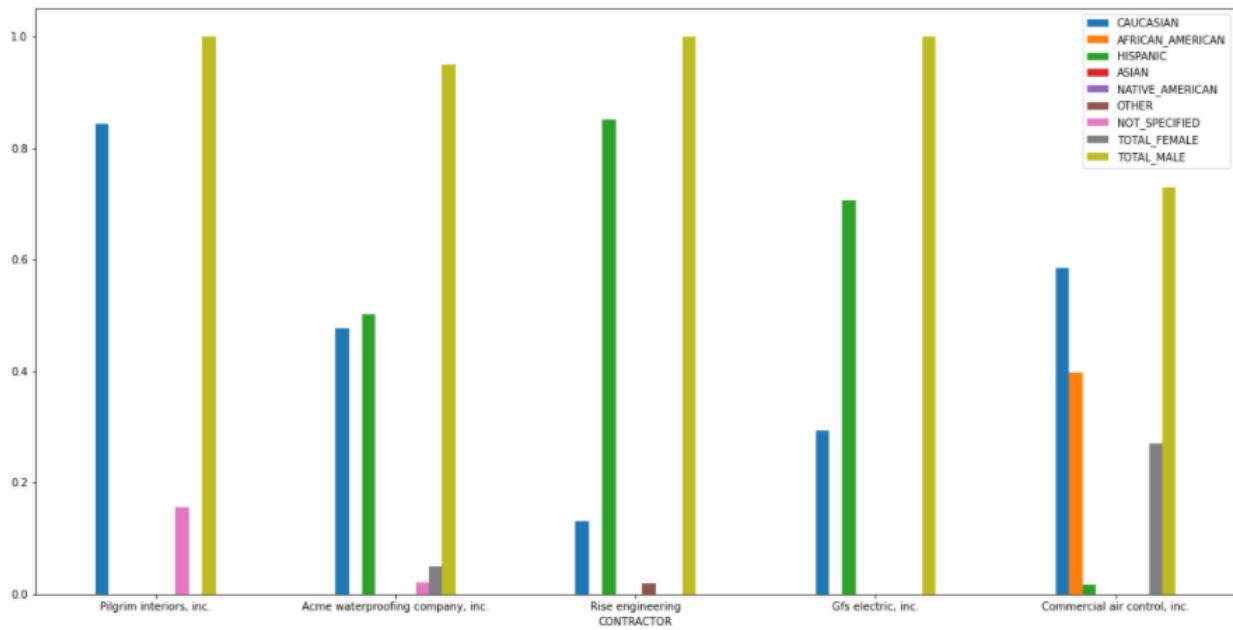
Top 5 companies based on total hours of Native Americans (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



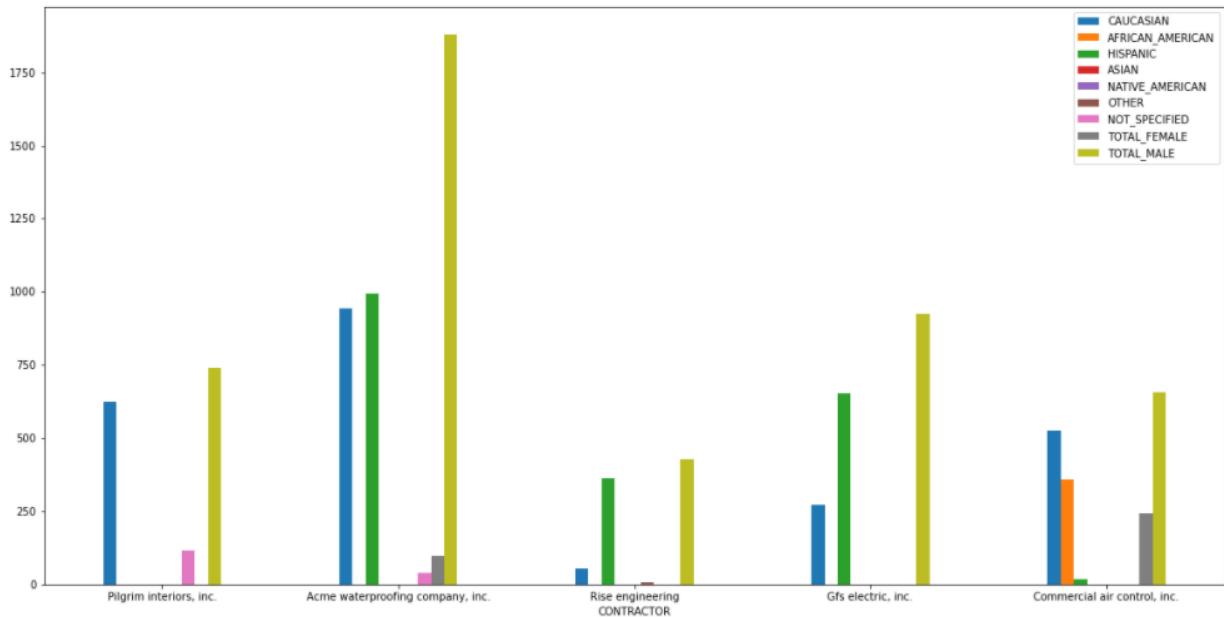
Top 5 companies based on percentage of Other (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



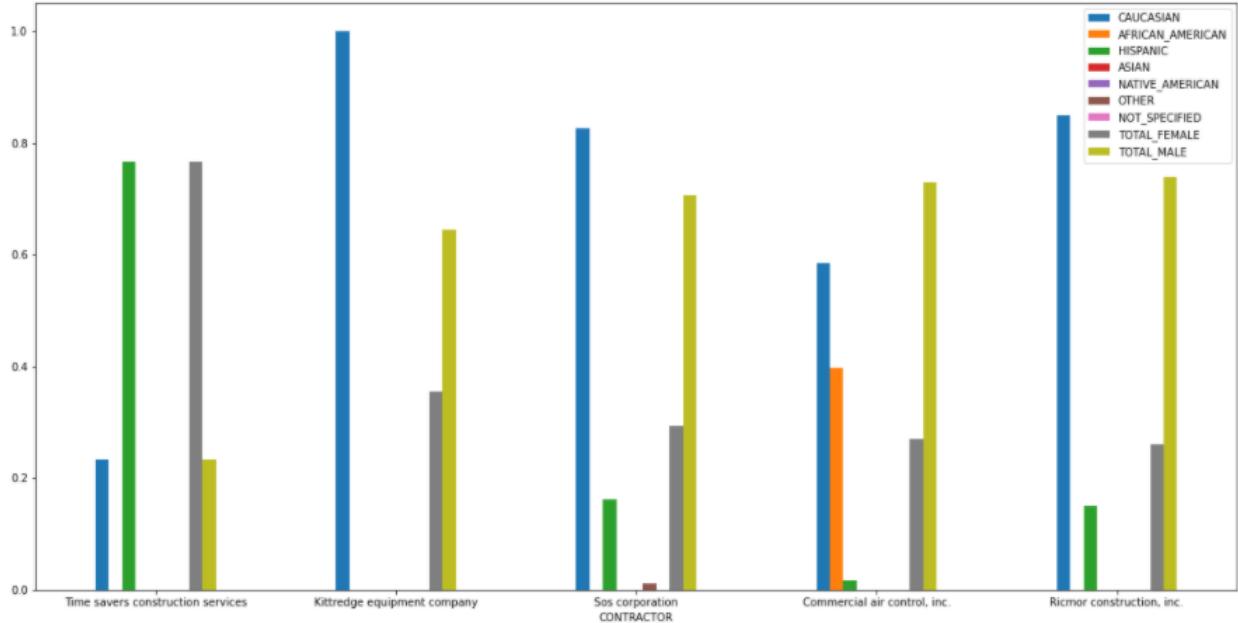
Top 5 companies based on total hours of Other (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



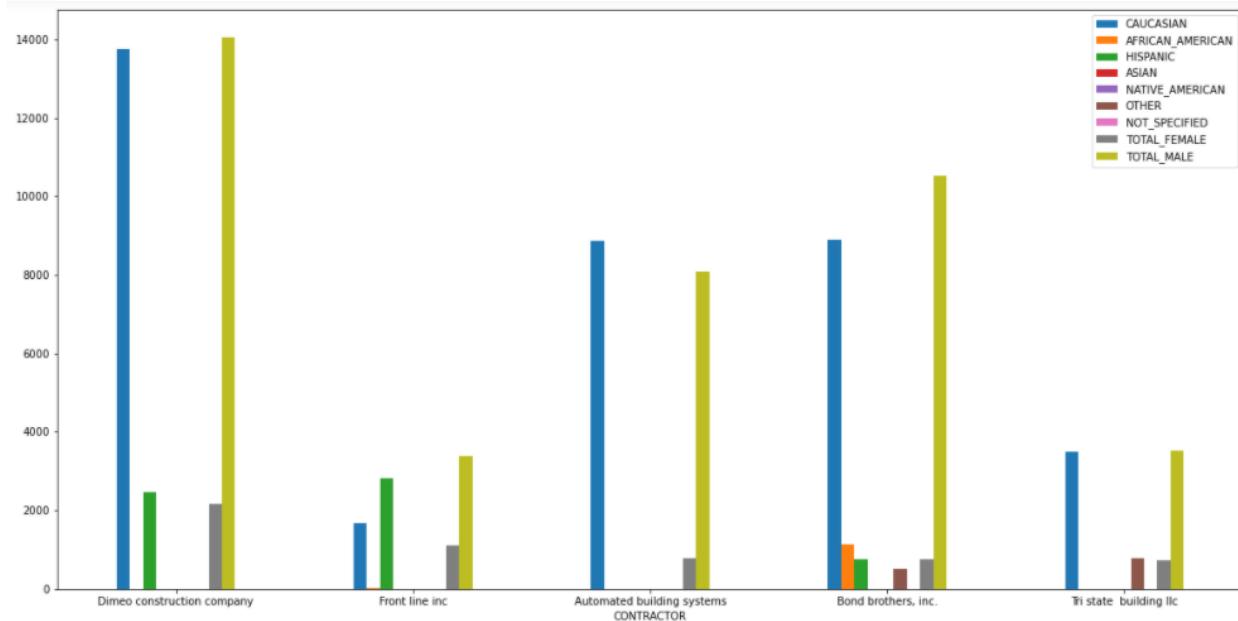
Top 5 companies based on percentage of Not Specified (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



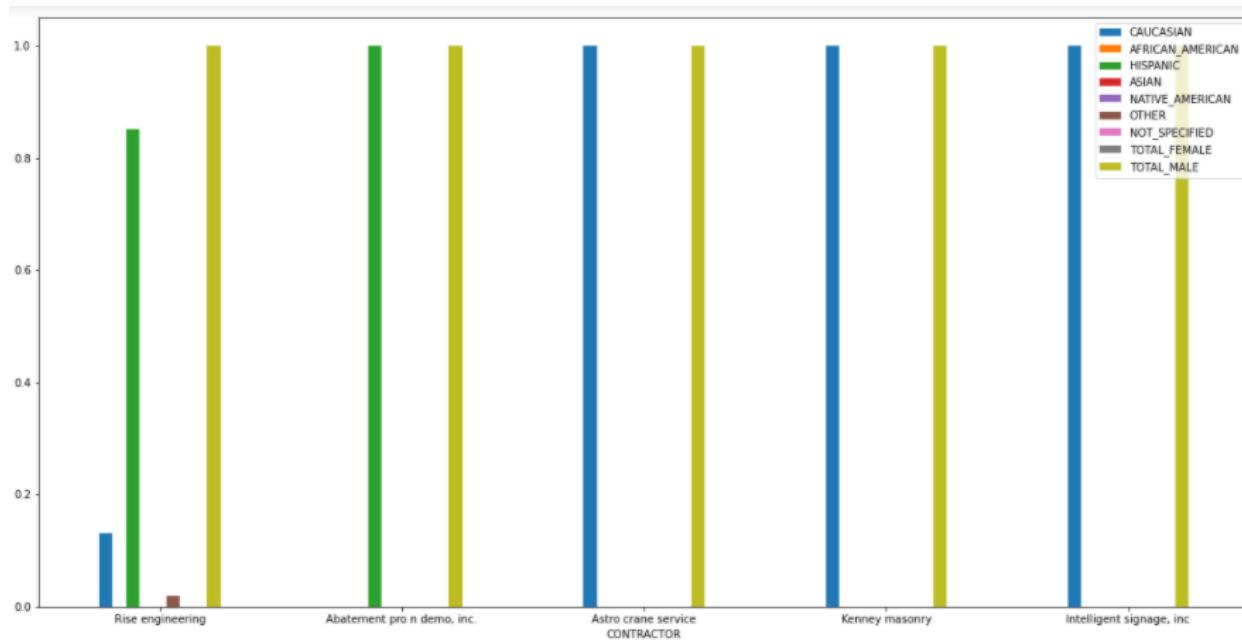
Top 5 companies based on total hours of Not Specified (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



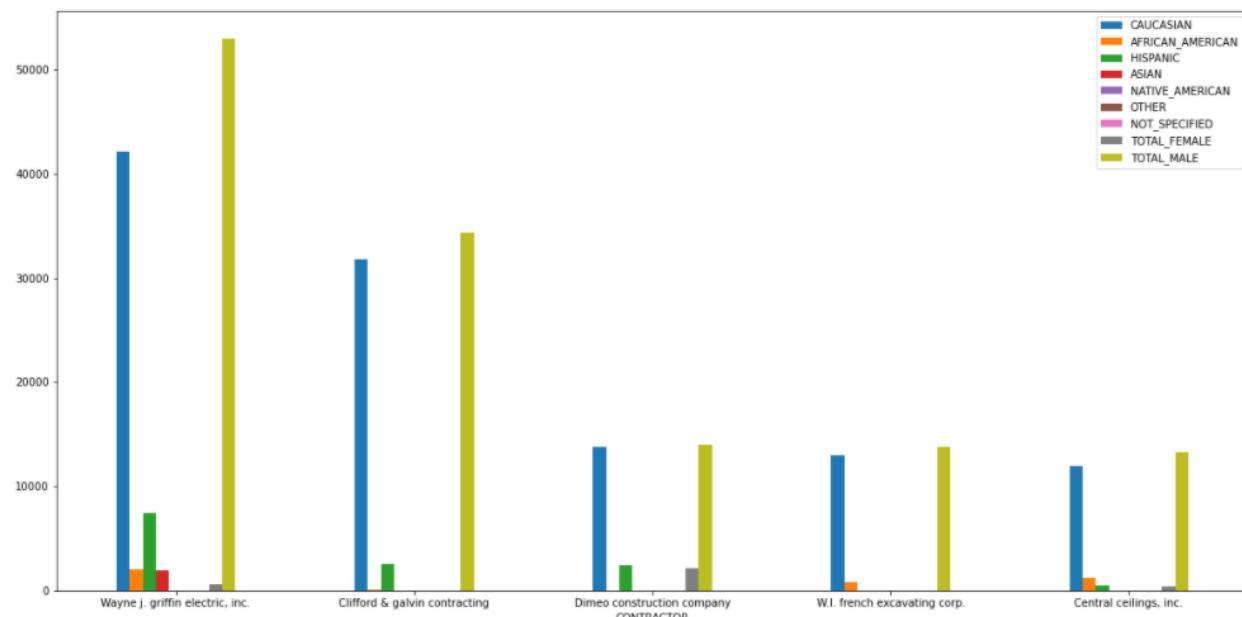
Top 5 companies based on percentage of Females (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



Top 5 companies based on total hours of Females (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



Top 5 companies based on percentage of Males (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)



Top 5 companies based on total hours of Males (Percentage aggregate of total hours worked per month for year for that ethnicity/gender divided by total hours worked for all ethnicities/genders)

Figure 2, Proof.txt file

```

▶ proof.txt x
1 starting page 1
2 Processing DF number: 0
3
4 0 Project Name:\rAEP1802E UT1 C Utility Simple F... ... 0 ... 10
5 1 Construction Trade ... NaN
6 2 Craft\rLevel ... Total\rMale NaN
7
8 [3 rows x 11 columns]
9
10 Processing row number: 0
11 ['MONTH', 'YEAR', 'PROJECT', 'PROJECT_CODE', 'CONTRACTOR', 'CONSTRUCTION_TRADE',
12 'CRAFT_LEVEL', 'TOTAL_EMPLOYEE', 'CAUCASIAN', 'AFRICAN_AMERICAN', 'HISPANIC', 'ASIAN',
13 'NATIVE_AMERICAN', 'OTHER', 'NOT_SPECIFIED', 'TOTAL_FEMALE', 'TOTAL_MALE',
14 'HOURS_WORKED_PER_MONTH']
15 Processing DF number: 1
16
17 0 Rise Engineering NaN 1 NaN 2 NaN 3 ... 8 NaN 9 NaN 10 NaN 11 NaN
18 1 INSULATOR (PIPES & TANKS) Journey 101.0 7.5 ... 0.0 0.0 0.0 0.0 101.0
19 2 Apprentice 0.0 0.0 ... 0.0 0.0 0.0 0.0 0.0 0.0 0.0
20 3 NaN A/J Ratio 0.0 0.0 ... 0.0 0.0 0.0 0.0 0.0 0.0 0.0
21 4 NaN New Hire 0.0 0.0 ... 0.0 0.0 0.0 0.0 0.0 0.0 0.0
22 5 NaN Subtotal 101.0 7.5 ... 0.0 0.0 0.0 0.0 0.0 101.0
23 6 Total for Contractor Journey 101.0 7.5 ... 0.0 0.0 0.0 0.0 0.0 101.0
24 7 Apprentice 0.0 0.0 ... 0.0 0.0 0.0 0.0 0.0 NaN
25 8 A/J Ratio 0.0 0.0 ... 0.0 0.0 0.0 0.0 0.0 NaN
26 9 New Hire 0.0 0.0 ... 0.0 0.0 0.0 0.0 0.0 0.0 NaN
27 10 Subtotal 101.00 7.5 0.0 ... 0.0 0.0 0.0 0.0 101.0 NaN
28 11 Total Journey Hours 101.00 7.5 0.0 ... 0.0 0.0 0.0 0.0 101.0 NaN
29
30
31 115269
32 115270 Processing row number: 0
33 115271 complete grp cj WSC1902 DB1 C WSC-Campus-wide-Redundant Steam and Condensate Loop 1
34 115272 complete grp project
35
36 115273 Processing DF number: 1
37
38 0 Total Journey Hours 1,056.00 904.5 0.0 ... 0.0 0.0 143.5 912.5
39 1 Total Apprentice Hours 7.00 7.0 0.0 ... 0.0 0.0 0.0 7.0
40 2 Total New Hire Hours 0.00 0.0 0.0 ... 0.0 0.0 0.0 0.0
41 3 Grand Total Hours 1,063.00 911.5 0.0 ... 0.0 0.0 143.5 919.5
42
43 115280 [4 rows x 11 columns]
44
45 115281 Processing row number: 0
46 115282 grand_total 4
47 115283 process grp pj
48 115284 process grp pa
49 115285 process grp pnh
50 115286 process grp pgrand
51 115287 complete grp grand_total
52
53 115288 MONTH YEAR ... TOTAL_MALE HOURS_WORKED_PER_MONTH
54 115289 0 12 2019 ... 16.0 HOURS_PER_MONTH
55 115290 0 12 2019 ... 0.0 HOURS_PER_MONTH
56 115291 0 12 2019 ... 0.0 HOURS_PER_MONTH
57 115292 0 12 2019 ... 165.0 HOURS_PER_MONTH
58 115293 0 12 2019 ... 1632.0 HOURS_PER_MONTH
59
60 115294 ... ...
61 115295 0 12 2019 ... ...
62 115296 0 12 2019 ... ...
63 115297 0 12 2019 ... ...
64 115298 0 12 2019 ... 7.0 HOURS_PER_MONTH
65 115299 0 12 2019 ... 393.0 HOURS_PER_MONTH
66 115300 0 12 2019 ... 0.0 HOURS_PER_MONTH
67 115301 0 12 2019 ... 512.5 HOURS_PER_MONTH
68 115302 0 12 2019 ... 0.0 HOURS_PER_MONTH
69
70 115303 [654 rows x 18 columns]
71 115304

```

Figure 3, Final DataMart DataFrame

MONTH	YEAR	PROJECT	PROJECT_CODE	CONTRACTOR	CONSTRUCTION_TRADE	CRAFT_LEVEL	TOTAL_EMPL...	CAUCASIAN	AFRICAN....	HISPANIC					
12	2019	TRC1407 FC1 C Ex...	TRC1407 FC1 C	North shore steel company, inc	IRONWORKER/WELDER	Journeymen	48.00000	0.00000	0.00000	24.00000					
12	2019	TRC1407 FC1 C Ex...	TRC1407 FC1 C	North shore steel company, inc	IRONWORKER/WELDER	Apprentice	0.00000	0.00000	0.00000	0.00000					
12	2019	TRC1407 FC1 C Ex...	TRC1407 FC1 C	North shore steel company, inc	LABORER	Journeymen	48.00000	0.00000	0.00000	24.00000					
12	2019	TRC1407 FC1 C Ex...	TRC1407 FC1 C	North shore steel company, inc	LABORER	Apprentice	0.00000	0.00000	0.00000	0.00000					
12	2019	TRC1407 FC1 C Ex...	TRC1407 FC1 C	Stanley roofing company, inc	ROOFER	Journeymen	620.25000	307.25000	0.00000	313.00000					
12	2019	TRC1407 FC1 C Ex...	TRC1407 FC1 C	Stanley roofing company, inc	ROOFER	Apprentice	0.00000	0.00000	0.00000	0.00000					
12	2019	TRC1407 FC1 C Ex...	TRC1407 FC1 C	Stanley roofing company, inc	SHEETMETAL WORKER	Journeymen	268.50000	100.00000	0.00000	168.50000					
12	2019	TRC1407 FC1 C Ex...	TRC1407 FC1 C	Stanley roofing company, inc	SHEETMETAL WORKER	Apprentice	0.00000	0.00000	0.00000	0.00000					
12	2019	TRC1407 FC1 C Ex...	TRC1407 FC1 C	Zap electric	ELECTRICIAN	Journeymen	523.00000	252.50000	0.00000	270.50000					
12	2019	TRC1407 FC1 C Ex...	TRC1407 FC1 C	Zap electric	ELECTRICIAN	Apprentice	0.00000	0.00000	0.00000	0.00000					
12	2019	TRC1702 HC1 C S...	TRC1702 HC1 C	3 phase elevator	ELEVATOR CONSTRUCTOR	Journeymen	447.50000	447.50000	0.00000	0.00000					
DR		CONSTRUCTION_TRADE	CRAFT_LEVEL	TOTAL_EMPL...	CAUCASIAN	AFRICAN....	HISPANIC	ASIAN	NATIVE_A...	OTHER	NOT_SPEC...	TOTAL_FE...	TOTAL_MALE	HOUR...	
197		GLAZIER	Apprentice	64.00000	0.00000	0.00000	64.00000	0.00000	0.00000	0.00000	0.00000	64.00000		HOUR...	
198		GLAZIER (GLASS PLANK/AIR BA...	Journeymen	264.00000	264.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	264.00000		HOUR...	
199		GLAZIER (GLASS PLANK/AIR BA...	Apprentice	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000		HOUR...	
200		IRONWORKER	Journeymen	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000		HOUR...	
201		IRONWORKER	Apprentice	93.00000	0.00000	0.00000	93.00000	0.00000	0.00000	0.00000	0.00000	0.00000	93.00000		HOUR...
202		IRONWORKER/WELDER	Journeymen	304.00000	304.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	304.00000		HOUR...	
203		IRONWORKER/WELDER	Apprentice	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000		HOUR...
204	Eng	CARPENTER	Journeymen	2344.50000	2160.50000	0.00000	184.00000	0.00000	0.00000	0.00000	0.00000	2344.50000		HOUR...	
205	Eng	CARPENTER	Apprentice	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000		HOUR...
206	Eng	LABORER	Journeymen	608.50000	608.50000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	608.50000		HOUR...	
207	Eng	LABORER	Apprentice	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000		HOUR...
208		ELEVATOR CONSTRUCTOR	Journeymen	350.00000	350.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	350.00000		HOUR...	
209		ELEVATOR CONSTRUCTOR	Apprentice	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000		HOUR...
210		ELEVATOR CONSTRUCTOR HEL...	Journeymen	281.25000	281.25000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	140.00000	141.25000		HOUR...
211		ELEVATOR CONSTRUCTOR HEL...	Apprentice	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000		HOUR...
212	any	CARPENTER	Journeymen	120.00000	120.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	120.00000		HOUR...	
213	any	CARPENTER	Apprentice	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000		HOUR...
214	any	LABORER	Journeymen	575.00000	463.00000	0.00000	112.00000	0.00000	0.00000	0.00000	0.00000	0.00000	575.00000		HOUR...
215	any	LABORER	Apprentice	157.00000	0.00000	0.00000	157.00000	0.00000	0.00000	0.00000	0.00000	0.00000	157.00000		HOUR...

Companies are ethnically specialized

- Excel/sheets behind ALL charts (ex: by trade from email) JENA (put drive)
- Document explaining data, code, and charts SABRINA +MURT
- Something on 2019-2020 months trend or Electricians ELISA
- Wrap up code and organize, visualization RICHARD (verify excel/sheets + run analysis code)
- Presentation (!!!) JENA, ELISA, MURT, RICHARD, SABRINA (2slides)
- Read me (explaining how to use the code to parse + analysis) JENA (verify: ELISA, RICHARD)
- Delv 4, due sunday MURT + RICHARD
- Final Report, 28th JENA, ELISA, MURT, RICHARD, SABRINA
- Paragraphs for client + choose/filter data to printing MURT + SABRINA

- Predictors ('If i wanted to predict how it turned out, what are the predictors?', 'Are your odds better or worse ?') (?)
- Track transfer of ownership (?)
- Doing charts as total, not % (?)
- Get more electricians analysis (?)

Excel behind all charts, bc looking at only %

Paper copies

Electrician -> numbers

Doc of explainer of data

Transfer ownership, from husband to wife, just bc Women on Business. How to identify?

Format: graphics <3, but explanation (!!) and also express totals (for %)

Or change them to totals

Predictive analysis -> no (who knows what is gonna happen tomorrow)

BUT: predictors, yes (trade, location, company...)

'If i wanted to predict how it turned out, what are the predictors?'

'Are your odds better or worse ?'

Wrap up code