

Project Deliverable 1

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. 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
- ❑ Second meeting with the client on Thurs March 4th
- ❑ 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

Data

The data is collected weekly by DCAMM. They sort it by months and keep it in PDF form. DCAMM already provided WGBH the work force from 2019 and will provide in March the data from 2020. The data is organized as tables of projects (such as bridges, buildings, etc) containing the companies included, their types of workers, and the hour rate separated 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 analysis. An example of a file is April 2019:

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

We have been given a file folder with files for each month Jan - Dec 2019, e.g. WorkforceUtilizationSummaryReportApril2018.pdf.

Methods

Since the main goal of this project is to get the data out of the PDF and transform it into a csv or .xlsx file in an organized way, our team is exploring multiple methods. Preliminary, we extracted the data using the online Tabula application, as well as programmatically extracting the data using the Tabula library (<https://pypi.org/project/tabula-py/>) and the PyPDF2 library. We also used acrobat to save the file in CSV format. Each method produced the same misalignments between the hourly data rows and the company/trades header data. The issue stems from the PDF merging cells to pretty print the data for human readability.

Our methodology will include a means to read through the data file and realign the company and the trade data with the hourly time spent rows. We may leverage Pandas and NumPy libraries when realigning the data. We will need to write a custom python script to complete our task.

When the data is aligned, we will run various machine learning algorithms to predict the number of hours assigned per project for each group. We can use the Scikit-learn library for the predictive analysis step. For visualization we will use Matplotlib and Plotly Dash.

Specific questions that will be answered:

1. How will we extract data from our PDF files?
2. Is there a difference between state-paid contractual hours based on color and/or sex?
3. What are the factors, e.g. location of the project, that fair in hiring working crews?
4. How do state-wise elections affect hiring decisions across projects? Did construction companies hire fewer minorities, people of color and/or women, during the pandemic?
5. Did the companies' work workload change because of COVID-19?

Our project work will have a larger portion of total time devoted to data cleanup vs data exploration due to the condition of the data sources. Subsequently, we will start to look for any patterns or trends to answers questions revolving around race, sex, and opportunities for state contracts.

Discussion and Limitations

Please refer to our screenshots 1-5 below to see our data import. On a preliminary analysis, we concluded that the data is divided by minorities and project type, as well as jobs per month. In header areas, the data is very sparse containing multiple zero rows and columns with nans inserted by Pandas(screenshot 5). Due to the obvious data misalignment and extraneous values, we will continue to build on our custom python parser for this project.

The client was only interested in the extraction of the data from PDF files and said we would discuss specific questions for analysis after this was completed. Our next step is to align our columns and make our data more organized and accessible to manipulate.

For this deliverable we were able to answer the first question listed in our methods section: How will we extract data from our PDF files? Now we have a new question which is: How to pre-process the data for our later analysis?

For limitations, when using libraries to transform PDFs into csv or xml, misplacements happen, such as elements in the wrong place and generating the wrong number of columns. To handle these issues, a common practice is to read the file in a series of passes. The first pass will extract the company names and be validated for accuracy by ourselves and our client. The second pass will extract the names of the trades and any

Our client has not specified analysis of our data because they understand the difficulty of massaging the data. Due to these limitations, our project will be limited and will not include all four aspects of data science analysis as defined by our instructor.

1. Tabula software example. Note the Company name row is empty for all other columns. This is the cause of the unalignment.

2. Tabula software CSV example. In this format, there is no visual problem, however when transforming into a Pandas Dataframe the row below the company name will merge with it the job name, and have one extra column.

[illegible]

3. Tabula library example showing how Tabula reformats all PDFs in a directory into like named csv files. The raw CSV output is shown on the right side panel:

The screenshot shows a VS Code editor with a project named 'WGBH-DCAMM'. The file explorer on the left shows a directory structure with 'data' and 'docs' folders. The 'data' folder contains numerous CSV files named 'WorkforceUtilizationSummaryReport' followed by a month and year (e.g., 'April2019.csv', 'Oct2019.csv'). The 'docs' folder contains 'Project Deliverable 0.pdf' and 'project_description.md'.

The right panel shows the raw CSV output of the 'WorkforceUtilizationSummaryReportOct2019.csv' file. The output is as follows:

```

1 "Project Name:
2 AEP1801E UT1 C Utility Vendor Contract- DCR- Statewide
3 Project Code:
4 AEP1801E UT1 C",,,,,,,,,,
5 Construction Trade,Hours Worked,,,,,,,,
6 "Craft
7 Level","Total
8 Employee",Caucasian,"African
9 American",Hispanic,Asian,"Native
10 American",Other,"Not
11 Specified","Total
12 Female","Total
13 Male"
14 "Jacqueline Electric and Contracting, Inc.",,,,,,,,,,
15 ELECTRICIAN,Journey,5.00,5.00,0.00,0.00,0.00,0.00,0.00,0.00,5.00
16 Apprentice,5.00,5.00,0.00,0.00,0.00,0.00,0.00,0.00,5.00,
17 A/J Ratio,1.00,1.00,0.00,0.00,0.00,0.00,0.00,0.00,1.00,
18 New Hire,0.00,0.00,0.00,0.00,0.00,0.00,0.00,0.00,0.00,
19 Subtotal,10.00,10.00,0.00,0.00,0.00,0.00,0.00,0.00,10.00,
20 "Journey
21 Apprentice
22 Total for Contractor
23 A/J Ratio
24 New Hire
25 Subtotal",5.00,"5.00
26 0.000.000.000.000.00","0.000.00
27 5.00",,,,,,,,,
28

```

The terminal at the bottom shows an error message: 'ImportError: cannot import name 'import_pdf' from 'src.parsers.pdf_parser' (/Users/jenajordahl/workspaces/CS-506-Homeworks/CS506Spring2021Repository/WGBH-DCAMM)'. Below the error, there is a command prompt showing the execution of a script: 'python3 color-of-money.py'.

4. An example of the first try at creating a custom python script to read the csv output is shown on the right side panel. The output is printed in the terminal window at the bottom.:

```
# WorkforceUtilizationSummaryReportDec2019.csv
# WorkforceUtilizationSummaryReportFeb2019.csv
# WorkforceUtilizationSummaryReportFeb2019.csv
# WorkforceUtilizationSummaryReportJan2019.csv
# WorkforceUtilizationSummaryReportJan2019.csv
# WorkforceUtilizationSummaryReportJuly2019.csv
# WorkforceUtilizationSummaryReportJuly2019.csv
# WorkforceUtilizationSummaryReportJune2019.csv
# WorkforceUtilizationSummaryReportJune2019.csv
# WorkforceUtilizationSummaryReportMarch2019.csv
# WorkforceUtilizationSummaryReportMarch2019.csv
# WorkforceUtilizationSummaryReportMay2019.csv
# WorkforceUtilizationSummaryReportMay2019.csv
# WorkforceUtilizationSummaryReportNov2019.csv
```

The screenshot shows a Python script in VS Code. The left sidebar displays a file explorer with various CSV files named 'WorkforceUtilizationSummaryReport' followed by months from Jan to Nov 2019. The main editor area contains the following code:

```
res = []
with open("data/WorkforceUtilizationSummaryReportDec2019.csv", 'r') as csv_file:
    lines = csv_file.readlines()
    for line in lines:
        row = []
        values = line.strip('\n').split(',')
        for val in values:
            row.append(convert_str(val))
        res.append(row)
print(res)
return res
```

Below the code editor, there's a terminal window showing the command prompt. It indicates the current directory is /workspaces/CS-506-Homeworks/CS506Spring2021Repository/WGBH-DCAMM and shows the execution of python3 color-of-money.py.

5. Pandas Dataframe example from Tabula. Here in row 14 it is possible to see one of the merging problems. Also, rows 7 to 12 were supposed to carry the race and gender, but they also have a misalignment.

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
0	Project Name:	None	None	None	None	None	None	None	None	None	None	None	None	None	None	None
1	AEP1802E UT1 C Utility Simple Fix	None	None	None	None	None	None	None	None	None	None	None	None	None	None	None
2	Project Code:	None	None	None	None	None	None	None	None	None	None	None	None	None	None	None
3	AEP1802E UT1 C											None	None	None	None	None
4	Construction Trade	Hours Worked										None	None	None	None	None
5	Craft	None	None	None	None	None	None	None	None	None	None	None	None	None	None	None
6	Level	Total	None	None	None	None	None	None	None	None	None	None	None	None	None	None
7	Employee	Caucasian	African	None	None	None	None	None	None	None	None	None	None	None	None	None
8	American	Hispanic	Asian	Native	None	None	None	None	None	None	None	None	None	None	None	None
9	American	Other	Not	None	None	None	None	None	None	None	None	None	None	None	None	None
10	Specified	Total	None	None	None	None	None	None	None	None	None	None	None	None	None	None
11	Female	Total	None	None	None	None	None	None	None	None	None	None	None	None	None	None
12	Male	None	None	None	None	None	None	None	None	None	None	None	None	None	None	None
13	Rise Engineering												None	None	None	None
14	ELECTRICIAN	Journey	16.0	8.0	0.0	0.0	0.0	0.0	8.0	0.0	0.0	16.0	None	None	None	None
15	Apprentice	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0		None	None	None	None
16	A/J Ratio	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0		None	None	None	None
17	New Hire	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0		None	None	None	None