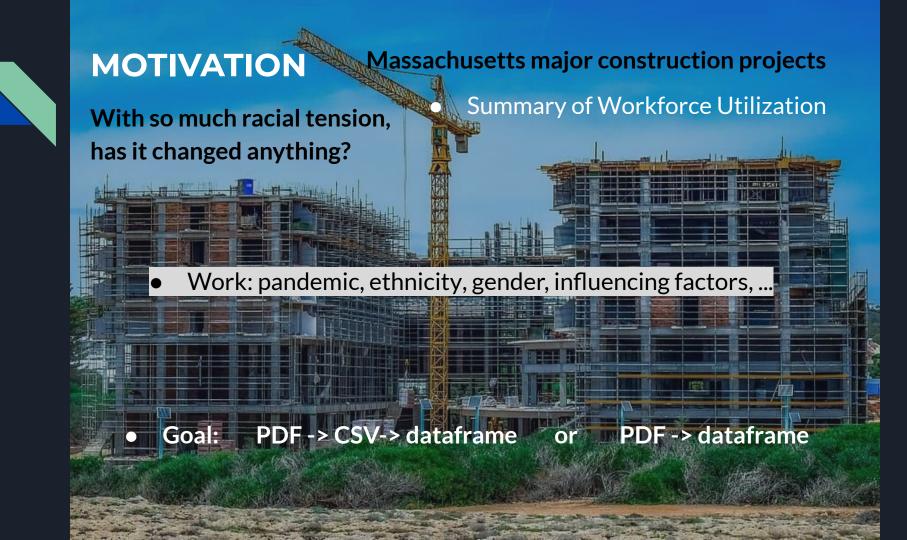
WGBH - Color of Money

Jena Jordahl, Murt Bahrani, Carmen Araújo, Richard Lee, Elisa Cordeiro Lopes

CS506 - Spring 2021



Method: configure bounding box parameters

Current Progress

+ tabula PDF reader -> pandas

+ PyPDF2 page count for tabula

+ custom parser groups report lines per pd.df and per page

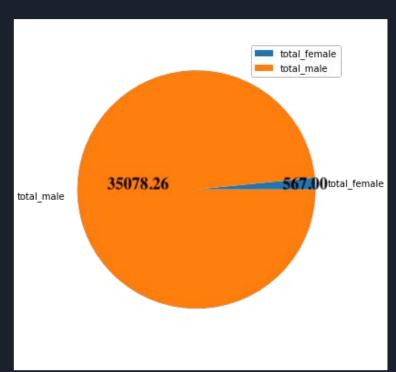
Project Code: CHE1604 DC1 CM						Hours Work	ed				
Construction Trade	Craft Level	Total Employee	Caucasian	African American	Hispanic	Asian	Native American	Other	Not Specified	Total Female	Total Male
E. Amanti & Sons, Inc											
	Journey	19.50	19.50	0.00	0.00	0.00	0.00	0.00	0.00	0.00	19.50
	Apprentice	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	00.0	0.00
LABORER	AU Ratio	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	New Hire	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	Subtotal	19.50	19.50	0.00	0.00	0.00	0.00	0.00	0.00	0.00	19.5
	Journey	20.00	20.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	20.0
	Apprentice	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.0
PIPERITTER & STEAMFITTER	AU Ratio	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.0
	New Hire	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.0
	Subtotal	20.00	20.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	20.0
	Journey	40.00	40.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	40.0
	Apprentice	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.0
PLUMBERS & GASFITTERS	AU Ratio	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.0
	New Hire	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.0
	Subtotal	40.00	40.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	40.0
	Journey	79.50	79.50	0.00	0.00	0.00	0.00	0.00	0.00	0.00	79.5
	Apprentice	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.0
Total for Contractor	AU Ratio	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.0
	New Hire	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.0
	Subtotal	79.50	79.50	0.00	0.00	0.00	0.00	0.00	0.00	0.00	79.5
arnum Industries LTD											
	Journey	11.00	11.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	11.0
	Accrentice	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.0
LABORER	AU Ratio	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	New Hire	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

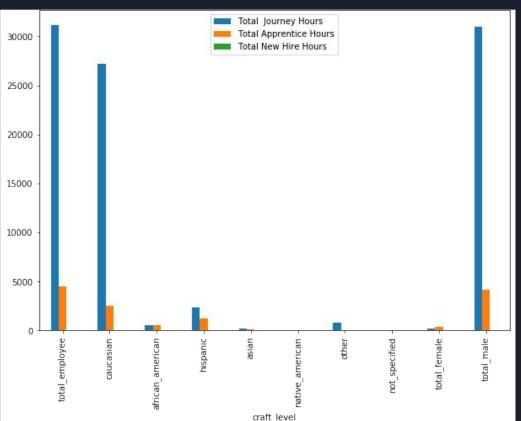
	OR	CONSTRUCTION_TRADE	CRAFT_LEVEL	TOTAL_EMPL	CAUCASIAN	AFRICAN	HISPANIC	ASIAN	NATIVE_A	OTHER	NOT_SPEC	TOTAL_FE	TOTAL_MALE	HOUR
197		GLAZIER	Apprentice	64.00000	0.00000	0.00000	64.00000	0.00000	0.00000	0.00000	0.00000	0.00000	64.00000	HOUR
198		GLAZIER (GLASS PLANK/AIR BA	Journeymen	264.00000	264.00000	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000	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	0.00000	HOUR
200		IRONWORKER	Journeymen	0.00000	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	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

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

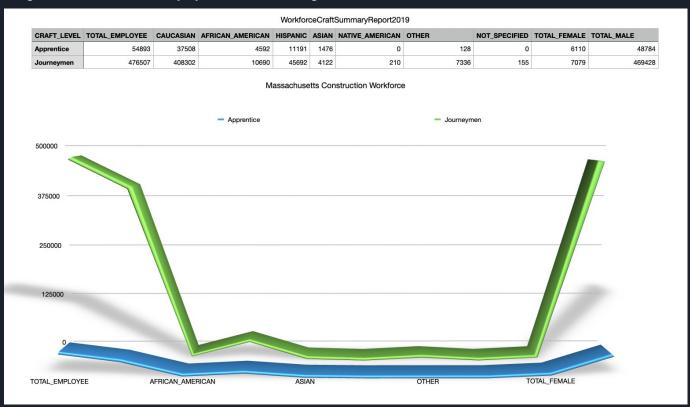
Early Results

Total number of hours by ethnicity split up by Journey, Apprentice, and New Hires





Findings Apprentice vs Journeymen If you don't apprentice you cannot become a Journeymen



KEY FINDINGS

Journeymen 86% Caucasian Americans
Apprentice 68% Caucasion Americans
Journeymen

98% Male Americans88% Male Americans2% African Americans8% African Americans

The color of money is significant in 2019, and so is its gender.

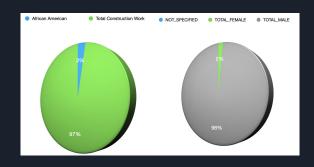
Apprentice

CRAFT_LEVEL	TOTAL_EMPLOYEE	CAUCASIAN	AFRICAN_AMERICAN	HISPANIC	ASIAN	NATIVE_AMERICAN	OTHER	NOT_SPECIFIED	TOTAL_FEMALE	TOTAL_MALE
Apprentice	54893.25	37507.75	4592.0	11190.5	1475.5	0.0	127.5	0.0	6109.75	48783.5
Journeymen	476507.39	408302.44	10689.83	45692.34	4121.5	209.5	7336.44	155.34	7079.44	469427.95
	• A	Apprentice	Journeyr		500000	100 - NOT 0	DEGISIED - TOTAL	N SEMALE -	TOTAL MALE	
	0.00	125000.00	250000.00 375000	0.00	500000	0.00 — NOT_S	PECIFIED - TOTA	AL_FEMALE —	TOTAL_MALE	
FOTAL_EMPLOYE	EE					-	en 1.5% Wo		50	0.0000
CAUCASIA	N				_	Apprentice	11% Wo	men_		
RICAN_AMERICA	AN .								37	5000.0
HISPAN	IIC								250	0.000.0
ASIA	AN								230	0000.0
NATIVE_AMERIC	AN _					Apprentice			125	5000.0
отн	ER					Apprentice			0.0	

Next Steps

• Discussing on possible analysis to do

ONE YEAR IS NOT AS IMPORTANT AS A 10 YR TREND



SOME PEOPLE FEEL BLOCKED IN EARNING
ENOUGH HOURS TO BECOME A JOURNEYMEN

How many different companies get contracts?

How many distinct companies are all one race?

How much bigger are the big companies?



Data Format Challenges



Report Header
Separate DF:
(3 rows & two categories)
Project Information,
Column Header Info

Trade Totals (5 lines):
Journeymen, Apprentice
A/J ratio, New Hire,
Trade SubTotal

Contractor Totals (5 lines): Journeymen, Apprentice A/J ratio, New Hire, Contractor SubTotal

Project Name: CHE1604 DC1 CM Construc Project Code: CHE1604 DC1 CM	tion Management	Services Che	elsea Soldiers	' Home Comu	ınity Living Ce	nter					
Troject Gode. CHE1004 DOT ON						Hours Work	ed				
Construction Trade	Craft Level	Total Employee	Caucasian	African American	Hispanic	Asian	Native American	Other	Not Specified	Total Female	Total Male
E. Amanti & Sons, Inc				.,					,		
	oouney								0.00		
	Apprentice	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
LABORER	A/J Ratio	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	New Hire	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	0	10.50	10.50	0.00	0.00	0.00	0.00	0.00	0.00	0.00	10.50
	Journey	20.00	20.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	20.00
	Apprentice	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
PIPEFITTER & STEAMFITTER	A/J Ratio	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	New Hire	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	Subtotal	20.00	20.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	20.00
	journe,	10.00	40.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	10.00
	Apprentice	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
PLUMBERS & GASFITTERS	A/J Ratio	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	New Hire	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	0.11.11	40.00	40.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	40.00
	Journey	79.50	79.50	0.00	0.00	0.00	0.00	0.00	0.00	0.00	79.50
	Apprentice	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Total for Contractor	A/J Ratio	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	New Hire	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	Subtotal	79.50	79.50	0.00	0.00	0.00	0.00	0.00	0.00	0.00	79.50
Harnum Industries LTD											
	Journey	11.00	11.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	11.00
	Apprentice	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
LABORER	A/J Ratio	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	New Hire	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
	Subtotal	11.00	11.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	11.00

Contractor
Header
(1 lines):
Name of
Contractor
without any
labeling



Data Format Challenges

33,000 lines of text and 1,097 pages

2019 Data

• 5,540 Detail lines

DF2

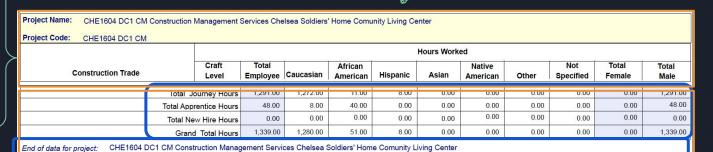
LIMITATIONS

We need more years for trend analysis

Report Header
Separate DF:
(3 rows & two categories)
Project Information
After 1st, throw away the
column heading info

Project Totals (4 lines): Journeymen, Apprentice New Hire, Contractor SubTotal

Project End (1 lines): Project Description



Jan2019 pages 61

Feb2019 pages 63
Mar2019 pages 75
Apr2019 pages 82
PG 2 of 82
May2019 pages 89
Jun2019 pages 90
Jul2019 pages 81
Aug2019 pages 81
Aug2019 pages 81
Aug2019 pages 103
Oct2019 pages 108
Nov2019 pages 129
Dec2019 pages 131
Total pages processed by the parser: 1097

DF3

CHALLENGES

- Directly read the PDF file into a Pandas DataFrame, each pg .. df
- Manually transformed into CSV
- PyPDF2, Tabula, Camelot, Pandas, PDFplumber, and batch Tabula-CSV

Hindsight:

The core process of a parser is identify and process, "IP".

Around the "IP", Is logic to get the Pages and Data Frames

```
Project Code
                                                     Project Name
                                                                               Contractor Name
                                                                                                                  Construction Trade
0 AFP1802F UT1 C
                                                    Utility Simple Fix
                                                                             Batallas Electric Inc.
                                                                                                                        FLECTRICIAN
    AEP1802E UT1 C
                                                    Utility Simple Fix
                                                                             Batallas Electric Inc.
                                                                                                                           LABORER
   AEP1802E UT1 C
                                                    Utility Simple Fix
                                                                                Rise Engineering
                                                                                                          INSULATOR (PIPES & TANKS)
3 CHE1604 DC1 CM Construction Management Services Chelsea Soldi...
                                                                             E. Amanti & Sons Inc
                                                                                                                           LABORER
4 CHE1604 DC1 CM Construction Management Services Chelsea Soldi.
                                                                             F Amanti & Sons Inc.
                                                                                                          PIPEFITTER & STEAMFITTER
5 CHE1604 DC1 CM Construction Management Services Chelsea Soldi...
                                                                             F Amanti & Sons Inc.
                                                                                                           PLUMBERS & GASFITTERS
6 CHE1604 DC1 CM Construction Management Services Chelsea Soldi.
                                                                           Harnum Industries LTD
                                                                                                                           LABORER
7 CHE1604 DC1 CM Construction Management Services Chelsea Soldi.
                                                                     S&F Concrete Contractors Inc. EQUIPMENT OPERATOR (Class B CDL)
8 CHE1604 DC1 CM Construction Management Services Chelsea Soldi.
                                                                    S&F Concrete Contractors Inc.
                                                                                                                           LABORER
9 CHE1604 DC1 CM Construction Management Services Chelsea Soldi... W.L. French Excavating Corp.
                                                                                                            ADS/SUBMERSIBLE PILOT
```

CUSTOM Parser

To test for a smaller number of pages set m to a small number while proc_pages <= m: print(f'starting page {proc_pages}') df0 = tabula.read_pdf(filename, pages=proc_pages, lattice=True, area=(11, 26, 582, 829), pandas_options={'header': None}) rows = 0 proc_dfs = 0 while proc_dfs < len(df0): print(f' Processing DF number: {proc_dfs}') print(df0[proc_dfs]) print() proc_rows = 0 while proc_rows < len(df0[proc_dfs]):</pre> print(f' Processing row number: {proc_rows}') if beginning: rows, COL, PROJCT, PROJCT_CD= start_process(df0[proc_dfs]) beginning = False else: grp = identify_grp(df0[proc_dfs],proc_rows) rows, remaining, STOREIT, NEWHIRE, PROJCT, PROJCT_CD, CONTRACTOR, CONSTR_TRAD process_grp(remaining, df0[proc_dfs],proc_rows,grp, STOREIT, NEWHIRE, CO proc_rows += rows proc_dfs += 1 proc_pages += 1 if debug: "./data/WorkforceUtilizationSummaryReportApril2019.pdf") outfile1 = outfile + MONTHYR[3:7] + MONTH + '.csv' STOREIT.to_csv(outfile1, index=False)

CHALLENGES: News stories need the facts to be in evidence

If WGBH wants to use the information we identify then they need proof that we parsed all that data according to the PDFs they received.

For use in a random sample test plan, the parser creates a proof.txt file to confirm that each row in the received PDF file was accurately categorized and processed in the final pandas DataFrame. On the right, please see excerpts from proof.txt:

FYI, proof.txt is 115,304 lines long.

```
starting page 1
Processing DF number: 0
  Project Name:\rAEP1802E UT1 C Utility Simple F...
                                         Craft\rLevel ... Total\rMale
[3 rows x 11 columns]
Processing row number: 0
['MONTH', 'YEAR', 'PROJECT', 'PROJECT_CODE', 'CONTRACTOR', 'CONSTRUCTION_TRADE',
 'CRAFT_LEVEL', 'TOTAL_EMPLOYEE', 'CAUCASIAN', 'AFRICAN_AMERICAN', 'HISPANIC', 'ASIAN',
 'NATIVE_AMERICAN', 'OTHER', 'NOT_SPECIFIED', 'TOTAL_FEMALE', 'TOTAL_MALE',
 'HOURS WORKED PER MONTH']
 Processing DF number: 1
             Rise Engineering
    INSULATOR (PIPES & TANKS)
                                                   0.0 ... 0.0
                          complete grp cj WSC1902 DB1 C WSC-Campus-wide-Redundant Steam and Condensate Loop :
                           complete arp project
10
                              Total Journey Hours 1,056.00 904.5 0.0 ...
                                                            7.0 0.0 ... 0.0 0.0
                                                            0.0 0.0 ... 0.0 0.0
                                Grand Total Hours 1,063.00 911.5 0.0 ... 0.0 0.0 143.5 919.5
                          [4 rows x 11 columns]
                          Processing row number: 0
                          grand total 4
                                                  16.0
                                                 393.0
```

Thank you!

Reference

1min motivation / background (why does your project matter?)

1min30s current progress, results, findings (what are the key findings of your project?)

30s next steps, limitations, challenges