

# DS311 Data Explore

Gabrielle Salamanca

March 28, 2023

## We-R-Finished Project

### 1. Loading necessary libraries & reading dataset into R

```
## Loading required package: ggplot2

## Warning: package 'ggplot2' was built under R version 4.2.3

## Loading required package: lattice

##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
## 
##     filter, lag

## The following objects are masked from 'package:base':
## 
##     intersect, setdiff, setequal, union

## Warning: Coercing text to numeric in Y146963 / R146963C25: '45870'

## Warning: Coercing text to numeric in Y164631 / R164631C25: '76700'

## [1] 167278      27
```

### 2. Renaming columns for sake of ease

```
## [1] "Case Number"           "Case Status"          "Received Date"
## [4] "Decision Date"        "Employer Name"        "Submitted Prevail Wage"
## [7] "SPrW Unit"            "Submitted Paid Wage" "SPaW Unit"
## [10] "Job Title"             "Work City"            "Required Edu"
## [13] "Required College Major" "Exp Req"              "Exp Req (Months)"
## [16] "Citizenship"          "Prevail Wage SOC Code" "PWSOC Title"
## [19] "Work State"           "WS Abb"                "WPostal Code"
## [22] "Full Time"            "Visa Class"           "Prevail Wage/Yr"
## [25] "Paid Wage/Yr"         "Job Title Sub"        "Order"
```

### 3. What are the job title subcategories & their numbers

```

##      assistant professor          attorney      business analyst
##                  18866             1488           27811
##      data analyst            data scientist management consultant
##                  3840              1227             770
##      software engineer          teacher
##                  99364             13912
##      Sub_Job Frequency
## 1 Software Engineer      99364
## 2 Business Analyst       27811
## 3 Assistant Prof        18866
## 4 Teacher                 13912
## 5 Data Scientist          1227
## 6 Data Analyst            3840
## 7 Attorney                 1488
## 8 Management Consultant    770

```

Was expecting higher number of data scientist and analyst Not expecting assistant professor & teacher in top 5

	NEW YORK	MOUNTAIN VIEW	SAN FRANCISCO	SAN DIEGO
##	6714	6712	6458	2983
##	SAN JOSE	HOUSTON	SANTA CLARA	ATLANTA
##	2925	2857	2684	2388
##	CHICAGO	DALLAS	SUNNYVALE	PALO ALTO
##	2351	2277	2198	2145
##	AUSTIN	IRVING	BOSTON	JERSEY CITY
##	1554	1410	1401	1257
##	MENLO PARK	SEATTLE	WASHINGTON	PITTSBURGH
##	1188	1146	1134	1132
##	REDMOND	BELLEVUE	LOS ANGELES	NORTH KANSAS CITY
##	1127	1116	986	932
##	PISCATAWAY	COLUMBUS	ALPHARETTA	BALTIMORE
##	919	911	904	879
##	PLANO	San Jose	SAN MATEO	PHOENIX
##	855	829	816	812
##	FREMONT	IRVINE	RICHMOND	REDWOOD CITY
##	800	768	757	755
##	CHARLOTTE	SAN ANTONIO	EDISON	Mountain View
##	749	728	716	697
##	PHILADELPHIA	CAMBRIDGE	BROOKLYN	DURHAM
##	697	692	678	662
##	NEWARK	PRINCETON	TAMPA	RICHARDSON
##	656	631	627	624
##	CUPERTINO	INDIANAPOLIS	WILMINGTON	FOSTER CITY
##	606	599	587	570
##	SCHAUMBURG	HERNDON	HILLSBORO	BLOOMINGTON
##	562	557	557	542
##	MILPITAS	RESTON	COLUMBIA	MINNEAPOLIS
##	538	536	528	527
##	WALTHAM	SOUTH PLAINFIELD	FAIRFAX	MIAMI

##	510	506	498	472
##	MCLEAN	SAN BRUNO	ST. LOUIS	FARMINGTON HILLS
##	470	461	445	438
##	KANSAS CITY	FRISCO	STERLING	JACKSONVILLE
##	432	431	422	408
##	PORLTAND	PLEASANTON	RALEIGH	DENVER
##	408	383	381	372
##	NASHVILLE	CHANDLER	MADISON	TROY
##	372	370	363	359
##	WARREN	BURLINGTON	San Diego	SOMERSET
##	356	349	348	348
##	DETROIT	GAITHERSBURG	CINCINNATI	Santa Clara
##	340	336	334	329
##	SANTA MONICA	OMAHA	DUBLIN	MILWAUKEE
##	326	324	314	312
##	SAN RAMON	ASHBURN	FORT WORTH	ROCHESTER
##	308	306	306	305
##	ALBANY	GLENDALE	ARLINGTON	(Other)
##	304	299	298	70689

```

##          top10 freq
## 1      New York 6714
## 2  Mountain View 6712
## 3  San Francisco 6458
## 4      San Diego 2983
## 5      San Jose 2925
## 6      Houston 2957
## 7  Santa Clara 2684
## 8      Atlanta 2388
## 9      Chicago 2351
## 10     Dallas 2277

```

Not surprised NY and CA areas are in top 5 cities Expected LA to be in top Wait i didn't expect to see my city, but it sure is a lot of CA areas

#### 4. Who are the top 10 employers

##	GOOGLE INC.
##	6213
##	FUJITSU AMERICA, INC.
##	1814
##	INTEL CORPORATION
##	1781
##	MICROSOFT CORPORATION
##	1364
##	QUALCOMM TECHNOLOGIES, INC.
##	1277
##	FACEBOOK, INC.
##	1210
##	CERNER CORPORATION
##	1038
##	CISCO SYSTEMS, INC.
##	1011

## HITACHI CONSULTING CORPORATION  
## 812  
## DALLAS INDEPENDENT SCHOOL DISTRICT  
## 808  
## EMC CORPORATION  
## 789  
## IBM CORPORATION  
## 762  
## VMWARE, INC.  
## 701  
## IBM INDIA PRIVATE LIMITED  
## 694  
## LARSEN & TOUBRO INFOTECH LIMITED  
## 679  
## WAL-MART ASSOCIATES, INC.  
## 674  
## LINKEDIN CORPORATION  
## 637  
## TWITTER, INC.  
## 632  
## YASH & LUJAN CONSULTING, INC.  
## 629  
## CAPITAL ONE SERVICES, LLC  
## 526  
## SYMANTEC CORPORATION  
## 520  
## APPLE INC.  
## 482  
## ACCENTURE LLP  
## 472  
## LOGISTIC SOLUTIONS, INC.  
## 453  
## ORACLE AMERICA, INC.  
## 448  
## DELOITTE CONSULTING LLP  
## 436  
## ITECH US, INC.  
## 417  
## BALTIMORE CITY PUBLIC SCHOOLS  
## 411  
## MOTOROLA MOBILITY LLC  
## 408  
## QUALCOMM INNOVATION CENTER INC.  
## 372  
## CAPGEMINI FINANCIAL SERVICES USA INC  
## 350  
## AKVARR INC  
## 336  
## IGATE TECHNOLOGIES INC.  
## 335  
## SOFTWARE PARADIGMS INTERNATIONAL GROUP, LLC  
## 334  
## INTUIT INC.  
## 333

## INTRAEDGE, INC. 331  
## ## QUALCOMM AHEROS, INC. 328  
## ## NEW YORK CITY DEPARTMENT OF EDUCATION 311  
## ## ## QUALCOMM INCORPORATED 302  
## ## ## ALINDUS, INC. 295  
## ## ## MCKINSEY & COMPANY, INC. UNITED STATES 294  
## ## ## KPIT INFOSYSTEMS INC. 286  
## ## ## SCM DATA, INC 284  
## ## ## MINDTREE LIMITED 282  
## ## ## OPENLOGIX CORPORATION 274  
## ## ## QUALCOMM TECHNOLOGIES INC. 267  
## ## ## THE MATHWORKS, INC. 267  
## ## ## ARGHA SERVICES, INC 260  
## ## ## GLOBAL TEACHERS RESEARCH & RESOURCES, INC 250  
## ## ## SYNECHRON, INC. 250  
## ## ## CYMA SYSTEMS INC 248  
## ## ## SALESFORCE.COM, INC. 247  
## ## ## VISION IT SERVICES USA INC 243  
## ## ## PHOTON INFOTECH, INC. 239  
## ## ## FRONTIER TECHNOLOGIES, LLC 232  
## ## ## SAMSUNG TELECOMMUNICATIONS AMERICA, LLC 229  
## ## ## FIDELITY TECHNOLOGY GROUP LLC 227  
## ## ## VISIONET SYSTEMS, INC. 227  
## ## ## BAHAI INDUSTRIES CORP. 225  
## ## ## CAPITAL ONE SERVICES II LLC 224  
## ## ## SUDHI INFOMATICS INC 219  
## ## ## CLIENT NETWORK SERVICES, INC. 217

## NORTHSTAR GROUP INC  
## 214  
## FACTSET RESEARCH SYSTEMS, INC.  
## 213  
## RANDSTAD TECHNOLOGIES, LP  
## 212  
## CITRIX SYSTEMS, INC.  
## 210  
## GOOGLE, INC.  
## 210  
## HEWLETT-PACKARD COMPANY  
## 209  
## YAHOO! INC.  
## 208  
## COOLSOFT, LLC  
## 207  
## ICONSOFT INC.  
## 207  
## CONSULTADD INC  
## 206  
## AKAMAI TECHNOLOGIES, INC.  
## 202  
## APEX TECHNOLOGY SYSTEMS, INC  
## 201  
## BROCADE COMMUNICATIONS SYSTEMS, INC.  
## 200  
## MOTOROLA SOLUTIONS, INC.  
## 198  
## EPAM SYSTEMS, INC  
## 194  
## UNIVERSITY OF ILLINOIS AT CHICAGO  
## 193  
## STAPLES, INC.  
## 189  
## AIM BIG, INC.  
## 188  
## CGI TECHNOLOGIES AND SOLUTIONS INC.  
## 187  
## MASTECH, INC., A MASTECH HOLDINGS, INC. COMPANY  
## 186  
## TABNER, INC.  
## 185  
## CSC COVANSYS CORPORATION  
## 183  
## JPMORGAN CHASE & CO.  
## 183  
## SASKEN COMMUNICATION TECHNOLOGIES, LTD.  
## 182  
## Groupon, INC.  
## 181  
## RIVERBED TECHNOLOGY, INC.  
## 180  
## UNIVERSITY OF MICHIGAN  
## 179

```

##          THE OHIO STATE UNIVERSITY
##                                177
##          MMC SYSTEMS INC.
##                                176
##          PURDUE UNIVERSITY
##                                176
##          THE UNIVERSITY OF ARIZONA
##                                175
##          TECH MAHINDRA (AMERICAS), INC.
##                                174
##          UNIVERSITY OF FLORIDA
##                                173
##          RJT COMPUQUEST, INC.
##                                171
##          ASTA CRS, INC.
##                                168
##          LOS ANGELES UNIFIED SCHOOL DISTRICT
##                                167
##          MIRACLE SOFTWARE SYSYTEM INC
##                                164
##          (Other)
##                                123034

```

```

##          top10 freq
## 1          Google Inc. 6213
## 2          Fujitsu America, Inc. 1814
## 3          Intel Corporation 1781
## 4          Microsoft Corporation 1364
## 5          Qualcomm Technologies, Inc. 1277
## 6          Facebook, Inc. 1210
## 7          Cerner Corporation 1038
## 8          Cisco Systems, Inc. 1011
## 9          Hitachi Consulting Corporation 812
## 10 Dallas Independent School District 2277

```

We know Google, Intel, Microsoft, Facebook, Cisco. The others I don't recognize. Also Dallas Independent School District is one of the top 10 employers?

## 5. Top 10 states to work in

##	Alabama	Alaska	Arizona
##	559	66	2570
##	Arkansas	California	Colorado
##	558	46782	1614
##	Connecticut	Delaware	District of Columbia
##	2023	874	1370
##	Florida	Georgia	Guam
##	4064	5615	61
##	Guamam	Hawaii	Idaho
##	1	170	186
##	Illinois	Indiana	Iowa
##	7411	1617	1012
##	Kansas	Kentucky	Louisiana

```

##          746          641          883
##          Maine        Maryland      Massachusetts
##          216          3275         6848
##          Michigan      Minnesota
##          3844         2116         Mississippi
##          Missouri      Montana
##          2773          63           Nebraska
##          Nevada        New Hampshire
##          427           427          New Jersey
##          New Mexico     New York
##          575           11373        North Carolina
##          North Dakota   Northern Mariana Islands
##          138            58           Ohio
##          Oklahoma       Oregon
##          551           1535         Palau
##          Pennsylvania   Puerto Rico
##          4725          109          Rhode Island
##          South Carolina South Dakota
##          851           105          Tennessee
##          Texas          Utah
##          15498          757          Vermont
##          Virgin Islands Virginia
##          109           6031         Washington
##          West Virginia  Wisconsin
##          307           1433         Wyoming
##          45

##          top10  freq
## 1    California  46782
## 2    Texas      15498
## 3    New York   11373
## 4    New Jersey  10198
## 5    Illinois   7411
## 6    Massachussets 6848
## 7    Virginia   6031
## 8    Pennsylvania 4725
## 9    Washington  4610
## 10   Michigan    3844

```

as expected CA is the largest, but WOW did the numbers really vary after. I was expecting TX and NY, but def not NJ.

## 6. Top 10 Jobs

```

##          SOFTWARE ENGINEER
##          54478
##          BUSINESS ANALYST
##          19516
##          SENIOR SOFTWARE ENGINEER
##          11714
##          ASSISTANT PROFESSOR
##          10134
##          Software Engineer
##          3743

```

## SR. SOFTWARE ENGINEER  
## 2133  
## COMPUTER SOFTWARE ENGINEER, APPLICATIONS  
## 1808  
## TEACHER  
## 1566  
## DATA ANALYST  
## 1524  
## SENIOR BUSINESS ANALYST  
## 1437  
## STAFF SOFTWARE ENGINEER  
## 1335  
## ELEMENTARY BILINGUAL TEACHER  
## 1020  
## COMPUTER SOFTWARE ENGINEER  
## 982  
## PRINCIPAL SOFTWARE ENGINEER  
## 962  
## Senior Software Engineer  
## 862  
## RESEARCH ASSISTANT PROFESSOR  
## 809  
## SOFTWARE ENGINEER II  
## 800  
## DATA SCIENTIST  
## 762  
## Assistant Professor  
## 678  
## LEAD SOFTWARE ENGINEER  
## 658  
## ASSOCIATE ATTORNEY  
## 650  
## SPECIAL EDUCATION TEACHER  
## 625  
## ASSOCIATE SOFTWARE ENGINEER  
## 591  
## SOFTWARE ENGINEER III  
## 579  
## SOFTWARE ENGINEER & TESTER  
## 529  
## SOFTWARE ENGINEER IN TEST  
## 496  
## VISITING ASSISTANT PROFESSOR  
## 495  
## CLINICAL ASSISTANT PROFESSOR  
## 485  
## SCIENCE TEACHER  
## 481  
## ELEMENTARY SCHOOL TEACHER  
## 463  
## ATTORNEY  
## 436  
## SR. BUSINESS ANALYST  
## 385

##	MATHEMATICS TEACHER	
##		375
## MEMBER OF TECHNICAL STAFF - SOFTWARE ENGINEERING		
##		369
##	SENIOR DATA ANALYST	
##		368
##	PUBLIC SCHOOL TEACHER	
##		361
##	SOFTWARE ENGINEER I	
##		327
##	MATH TEACHER	
##		290
##	software engineer	
##		285
##	Business Analyst	
##		278
##	IT BUSINESS ANALYST	
##		275
## EMBEDDED SOFTWARE ENGINEER		
##		266
##	KINDERGARTEN TEACHER	
##		237
##	SR SOFTWARE ENGINEER	
##		231
##	Sr. Software Engineer	
##		231
## FOREIGN LANGUAGE TEACHER		
##		230
## SOFTWARE ENGINEER AND TESTER		
##		225
## ELEMENTARY TEACHER		
##		222
## LEAD TEACHER		
##		213
## SENIOR CONSULTANT/BUSINESS ANALYST		
##		202
## BUSINESS DATA ANALYST		
##		192
## GRAPHICS SOFTWARE ENGINEER		
##		190
## BILINGUAL ELEMENTARY TEACHER		
##		189
## SECONDARY SCHOOL TEACHER		
##		186
## CLINICAL DATA ANALYST		
##		183
## SPANISH TEACHER		
##		177
## BILINGUAL TEACHER		
##		172
## BUSINESS ANALYST II		
##		159
## COMPUTER SOFTWARE ENGINEER - APPLICATIONS		
##		157

##	MANAGEMENT CONSULTANT	
##		152
##	MIDDLE SCHOOL TEACHER	
##		149
##	PRESCHOOL TEACHER	
##		149
##	ASSOCIATE (MANAGEMENT CONSULTANT)	
##		145
##	ADVISORY SOFTWARE ENGINEER	
##		144
##	SENIOR DATA SCIENTIST	
##		142
##	ASSISTANT PROFESSOR OF ECONOMICS	
##		141
##	LEAD BUSINESS ANALYST	
##		141
##	SYSTEMS/SOFTWARE ENGINEER	
##		141
##	SAP BUSINESS ANALYST	
##		134
##	COMPUTER TEACHER	
##		132
##	SOFTWARE ENGINEER/DEVELOPER	
##		132
##	SENIOR JAVA SOFTWARE ENGINEER	
##		131
##	JAVA SOFTWARE ENGINEER	
##		130
##	MEMBER OF TECHNICAL STAFF - SOFTWARE ENGINEER	
##		129
##	MUSIC TEACHER	
##		128
##	SOFTWARE ENGINEER APPLICATIONS	
##		128
##	ASSISTANT PROFESSOR OF MATHEMATICS	
##		126
##	ASSOCIATE BUSINESS ANALYST	
##		122
##	MANAGER, SOFTWARE ENGINEERING	
##		122
##	ASSISTANT PROFESSOR OF CLINICAL MEDICINE	
##		121
##	HIGH SCHOOL MATH TEACHER	
##		119
##	SYSTEM SOFTWARE ENGINEER	
##		115
##	HIGH SCHOOL SCIENCE TEACHER	
##		112
##	ASSISTANT PROFESSOR OF MEDICINE	
##		111
##	Elementary Bilingual Teacher	
##		111
##	BUSINESS ANALYSTS	
##		110

```

##                      ESL TEACHER
##                                108
##          SENIOR STAFF SOFTWARE ENGINEER
##                                104
##          SOFTWARE ENGINEER IV
##                                103
##          TECHNICAL BUSINESS ANALYST
##                                103
##          SOFTWARE ENGINEER APPS
##                                102
##          SOFTWARE ENGINEER, APPLICATIONS
##                                101
##          ASSOCIATE SOFTWARE ENGINEER DEVELOPER
##                                100
##          ASSISTANT PROFESSOR OF FINANCE
##                                96
##          CLINICAL BUSINESS ANALYST
##                                96
##          INTERNATIONAL BUSINESS ANALYST
##                                95
##          SENIOR .NET SOFTWARE ENGINEER
##                                95
##          SENIOR SOFTWARE ENGINEER DEVELOPER
##                                95
##          ASSISTANT PROFESSOR OF COMPUTER SCIENCE
##                                92
##          (Other)
##                                32745

##                      top10 freq
## 1      Software Engineer 58221
## 2          Business Analyst 19516
## 3      Sr Software Engineer 13847
## 4      Assistant Professor 10134
## 5 Computer Software Engineer, Apps 1808
## 6                  Teacher 1566
## 7          Data Analyst 1524
## 8      Sr Business Analyst 1437
## 9      Staff Software Engineer 1335
## 10     Elementary Bilungual Teacher 1020

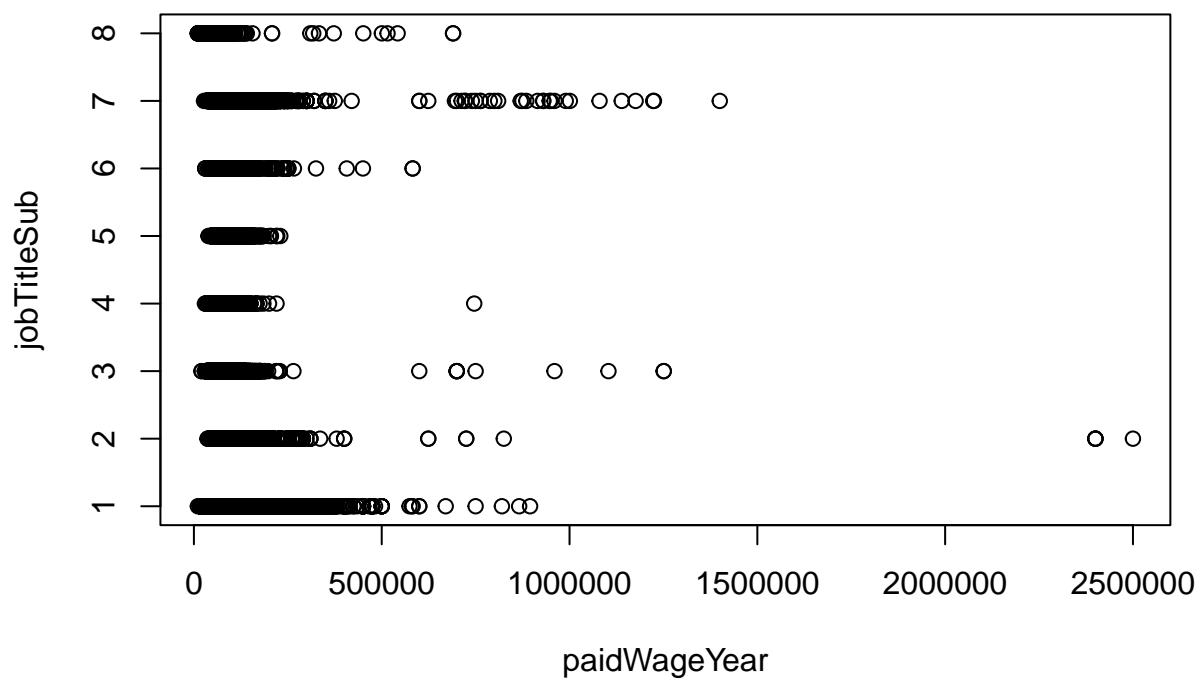
```

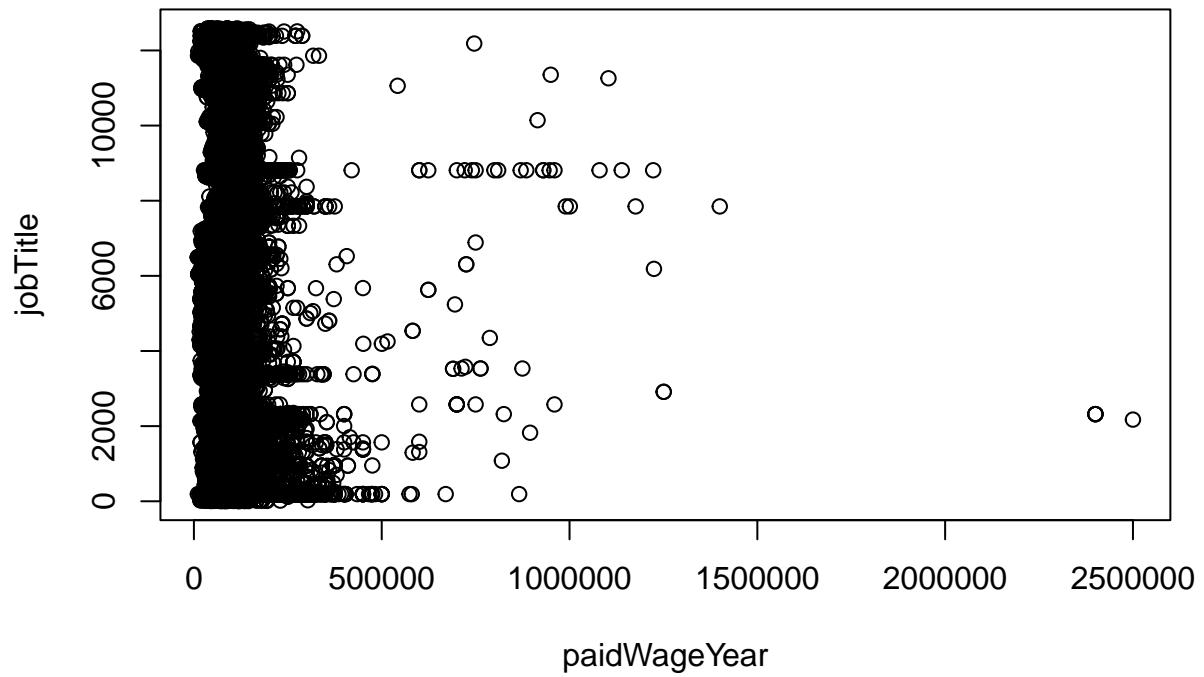
Dataset separated the ones with different cases, had to add up a few. Wasn't expecting teachers, esp the last one. engineer and analyst are popular jobs it seems

```

##    Min. 1st Qu. Median    Mean 3rd Qu.    Max.
##    10500   63000   78600   85533 100006 2500000

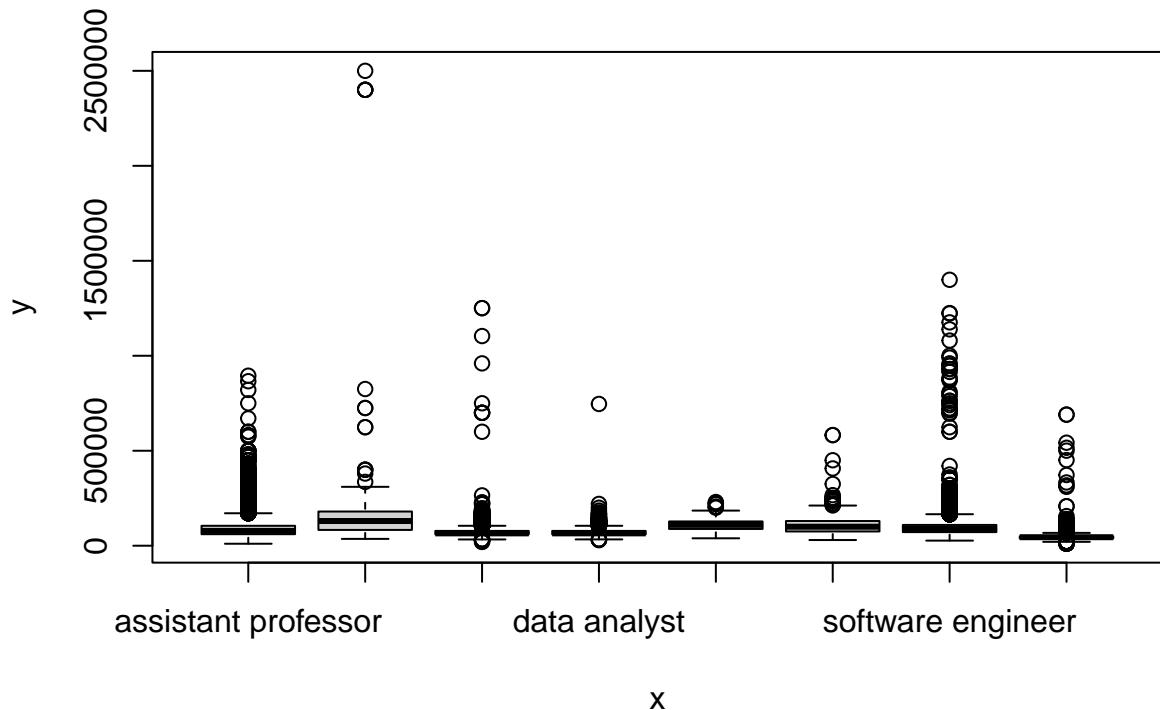
```





Norm seems to be under \$500,000 wage/year. If only the table can tell me which sub jobs they were... I will need to revisit the jobTitle one, but the norm still holds.

## 7. Trying the plots again



Hard to read

```
knitr::opts_chunk$set(echo = FALSE)
# loading libraries
library(readxl)
library(caret)
library(ggplot2)
library(dplyr)

# dataset overview
salary <- read_excel("C:/Users/knight/OneDrive/Desktop/Github/We-R-Finished/salary/salary_data_states.xlsx")
dim(salary)
# renaming
colnames(salary)[1] = "Case Number"
colnames(salary)[2] = "Case Status"
colnames(salary)[3] = "Received Date"
colnames(salary)[4] = "Decision Date"
colnames(salary)[5] = "Employer Name"
colnames(salary)[6] = "Submitted Prevail Wage"
colnames(salary)[7] = "SPRW Unit"
colnames(salary)[8] = "Submitted Paid Wage"
colnames(salary)[9] = "SPaW Unit"
colnames(salary)[10] = "Job Title"
colnames(salary)[11] = "Work City"
```

```

colnames(salary)[12] = "Required Edu"
colnames(salary)[13] = "Required College Major"
colnames(salary)[14] = "Exp Req"
colnames(salary)[15] = "Exp Req (Months)"
colnames(salary)[16] = "Citizenship"
colnames(salary)[17] = "Prevail Wage SOC Code"
colnames(salary)[18] = "PWSOC Title"
colnames(salary)[19] = "Work State"
colnames(salary)[20] = "WS Abb"
colnames(salary)[21] = "WPostal Code"
colnames(salary)[22] = "Full Time"
colnames(salary)[23] = "Visa Class"
colnames(salary)[24] = "Prevail Wage/Yr"
colnames(salary)[25] = "Paid Wage/Yr"
colnames(salary)[26] = "Job Title Sub"
colnames(salary)[27] = "Order"
names(salary)
jobTitleSub <- as.factor(salary$`Job Title Sub`)
summary(jobTitleSub)
# dataframe
Sub_Job <- c("Software Engineer", "Business Analyst", "Assistant Prof", "Teacher",
             "Data Scientist", "Data Analyst", "Attorney", "Management Consultant")
Frequency <- c(99364, 27811, 18866, 13912, 1227, 3840, 1488, 770)
job2 <- data.frame(Sub_Job, Frequency)
print(job2)
city <- as.factor(salary$`Work City`)
summary(city)
top10 <- c("New York", "Mountain View", "San Francisco", "San Diego", "San Jose", "Houston", "Santa Clara")
freq <- c(6714, 6712, 6458, 2983, 2925, 2957, 2684, 2388, 2351, 2277)
topCity <- data.frame(top10, freq)
print(topCity)
employName <- as.factor(salary$`Employer Name`)
summary(employName)
top10 <- c("Google Inc.", "Fujitsu America, Inc.", "Intel Corporation",
           "Microsoft Corporation", "Qualcomm Technologies, Inc.", "Facebook, Inc.",
           "Cerner Corporation", "Cisco Systems, Inc.", "Hitachi Consulting Corporation", "Dallas Independent School District")
freq <- c(6213, 1814, 1781, 1364, 1277, 1210, 1038, 1011, 812, 2277)
topEmploy <- data.frame(top10, freq)
print(topEmploy)
state <- as.factor(salary$`Work State`)
summary(state)
top10 <- c("California", "Texas", "New York", "New Jersey", "Illinois",
           "Massachusetts", "Virginia", "Pennsylvania", "Washington", "Michigan")
freq <- c(46782, 15498, 11373, 10198, 7411, 6848, 6031, 4725, 4610, 3844)
topState <- data.frame(top10, freq)
print(topState)
jobTitle <- as.factor(salary$`Job Title`)
summary(jobTitle)
top10 <- c("Software Engineer", "Business Analyst", "Sr Software Engineer", "Assistant Professor", "Computer Science", "Teacher", "Data Analyst", "Sr Business Analyst", "Staff Software Engineer", "Elementary Bilingual Education", "Other")
freq <- c(58221, 19516, 13847, 10134, 1808, 1566, 1524, 1437, 1335, 1020)
topJob <- data.frame(top10, freq)
print(topJob)

```

```
paidWageYear <- salary$`Paid Wage/Yr`  
summary(paidWageYear)  
# plots  
plot(paidWageYear, jobTitleSub)  
plot(paidWageYear, jobTitle)  
plot(jobTitleSub, paidWageYear)
```