

H-1B Applications Data Exploration

Sachin Mohan Sujir

2020-10-07

Contents

INTRODUCTION	1
INITIALIZATION	1
DATA PRE-PROCESSING	2
EXPLORATION	4
H1B Visa exploration	5
CONCLUSION	15

INTRODUCTION

In this report, I have performed an exploration of the H1B application data. The dataset size is around 528K, where each record contains information about the visa application filed by the employer for non-immigrant workers. In the data, there are about four types of VISA (H1B, E3 Australian, H1B1 Singapore, and H1B1 Chile) filed during the years from 2011 to 2017. H-1B visas are work authorization visas required by internationals to work in the USA (temporarily).

INITIALIZATION

Here, the required packages and the H1B dataset is loaded and have replaced the empty cells with an NA. Pander is designed to provide a minimal and easy tool for rendering R objects into Pandoc's markdown

```
library(tidyverse)

## -- Attaching packages ----- tidyverse 1.3.0 --
## v ggplot2 3.2.1      v purrr   0.3.3
## v tibble  2.1.3      v dplyr  0.8.3
## v tidyr   1.0.0      v stringr 1.4.0
## v readr   1.3.1      v forcats 0.4.0

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()

library(data.table)

## data.table 1.13.0 using 4 threads (see ?getDTthreads). Latest news: r-datatable.com
##
## Attaching package: 'data.table'

## The following objects are masked from 'package:dplyr':
##
##   between, first, last
```

```
## The following object is masked from 'package:purrr':
##
##      transpose
library(pander)
library(scales)

##
## Attaching package: 'scales'

## The following object is masked from 'package:purrr':
##
##      discard

## The following object is masked from 'package:readr':
##
##      col_factor
library(ggplot2)
library(xtable)
h1bData <- fread("h1bdata.csv", na.strings = c("", "NA", "N/A"))
```

DATA PRE-PROCESSING

Here, I have performed a few data pre-processing steps by selecting only required/relevant columns and removing duplicates from the dataset. Following shows the relevant column names and head of the dataset.

```
names(h1bData)

## [1] "CASE_SUBMITTED_DAY"      "CASE_SUBMITTED_MONTH"
## [3] "CASE_SUBMITTED_YEAR"    "DECISION_DAY"
## [5] "DECISION_MONTH"         "DECISION_YEAR"
## [7] "VISA_CLASS"             "EMPLOYER_NAME"
## [9] "EMPLOYER_STATE"         "EMPLOYER_COUNTRY"
## [11] "SOC_NAME"               "NAICS_CODE"
## [13] "TOTAL_WORKERS"          "FULL_TIME_POSITION"
## [15] "PREVAILING_WAGE"         "PW_UNIT_OF_PAY"
## [17] "PW_SOURCE"              "PW_SOURCE_YEAR"
## [19] "PW_SOURCE_OTHER"        "WAGE_RATE_OF_PAY_FROM"
## [21] "WAGE_RATE_OF_PAY_TO"    "WAGE_UNIT_OF_PAY"
## [23] "H-1B_DEPENDENT"         "WILLFUL_VIOLATOR"
## [25] "WORKSITE_STATE"         "WORKSITE_POSTAL_CODE"
## [27] "CASE_STATUS"

options(xtable.comment=FALSE)
options(xtable.booktabs=TRUE)
options(xtable.result=axis)
h1bData<-h1bData %>%
  select(CASE_SUBMITTED_DAY,
         CASE_SUBMITTED_MONTH,
         CASE_SUBMITTED_YEAR,
         DECISION_DAY,
         DECISION_MONTH,
         DECISION_YEAR,
         VISA_CLASS,
         EMPLOYER_NAME,
         SOC_NAME,
```

```

TOTAL_WORKERS,
FULL_TIME_POSITION,
PREVAILING_WAGE,
PW_UNIT_OF_PAY,
WAGE_RATE_OF_PAY_FROM,
WAGE_RATE_OF_PAY_TO,
WAGE_UNIT_OF_PAY,
'H-1B_DEPENDENT',
WILLFUL_VIOLATOR,
WORKSITE_STATE,
CASE_STATUS)

```

```

h1bData <- distinct(h1bData)
dim(h1bData)

```

```
## [1] 456549      20
```

```
names(h1bData)
```

```

## [1] "CASE_SUBMITTED_DAY"      "CASE_SUBMITTED_MONTH"
## [3] "CASE_SUBMITTED_YEAR"    "DECISION_DAY"
## [5] "DECISION_MONTH"         "DECISION_YEAR"
## [7] "VISA_CLASS"              "EMPLOYER_NAME"
## [9] "SOC_NAME"                "TOTAL_WORKERS"
## [11] "FULL_TIME_POSITION"     "PREVAILING_WAGE"
## [13] "PW_UNIT_OF_PAY"         "WAGE_RATE_OF_PAY_FROM"
## [15] "WAGE_RATE_OF_PAY_TO"    "WAGE_UNIT_OF_PAY"
## [17] "H-1B_DEPENDENT"         "WILLFUL_VIOLATOR"
## [19] "WORKSITE_STATE"         "CASE_STATUS"

```

```
pander(head(h1bData))
```

Table 1: Table continues below

CASE_SUBMITTED_DAY	CASE_SUBMITTED_MONTH	CASE_SUBMITTED_YEAR	DECISION_DAY
24	2	2016	1
4	3	2016	1
10	3	2016	1
28	9	2016	1
22	2	2015	2
12	3	2015	2

Table 2: Table continues below

DECISION_MONTH	DECISION_YEAR	VISA_CLASS	EMPLOYER_NAME
10	2016	H1B	DISCOVER PRODUCTS INC
10	2016	H1B	DFS SERVICES LLC
10	2016	H1B	EASTBANC TECHNOLOGIES LLC
10	2016	H1B	INFO SERVICES LLC
10	2016	H1B	BBandT CORPORATION
10	2016	H1B	SUNTRUST BANKS INC

Table 3: Table continues below

SOC_NAME	TOTAL_WORKERS	FULL_TIME_POSITION	PREVAILING_WAGE
ANALYSTS	1	Y	59197
ANALYSTS	1	Y	49800
ANALYSTS	2	Y	76502
COMPUTER OCCUPATION	1	Y	90376
ANALYSTS	1	Y	116605
ANALYSTS	1	Y	59405

Table 4: Table continues below

PW_UNIT_OF_PAY	WAGE_RATE_OF_PAY_FROM	WAGE_RATE_OF_PAY_TO
Year	65811	67320
Year	53000	57200
Year	77000	0
Year	102000	0
Year	132500	0
Year	71750	0

Table 5: Table continues below

WAGE_UNIT_OF_PAY	H- 1B_DEPENDENT	WILLFUL_VIOLATOR	WORKSITE_STATE
Year	N	N	IL
Year	N	N	IL
Year	Y	N	DC
Year	Y	N	NJ
Year	N	N	NY
Year	N	N	GA

CASE_STATUS
CERTIFIEDWITHDRAWN
CERTIFIEDWITHDRAWN
CERTIFIEDWITHDRAWN
WITHDRAWN
CERTIFIEDWITHDRAWN
CERTIFIEDWITHDRAWN

EXPLORATION

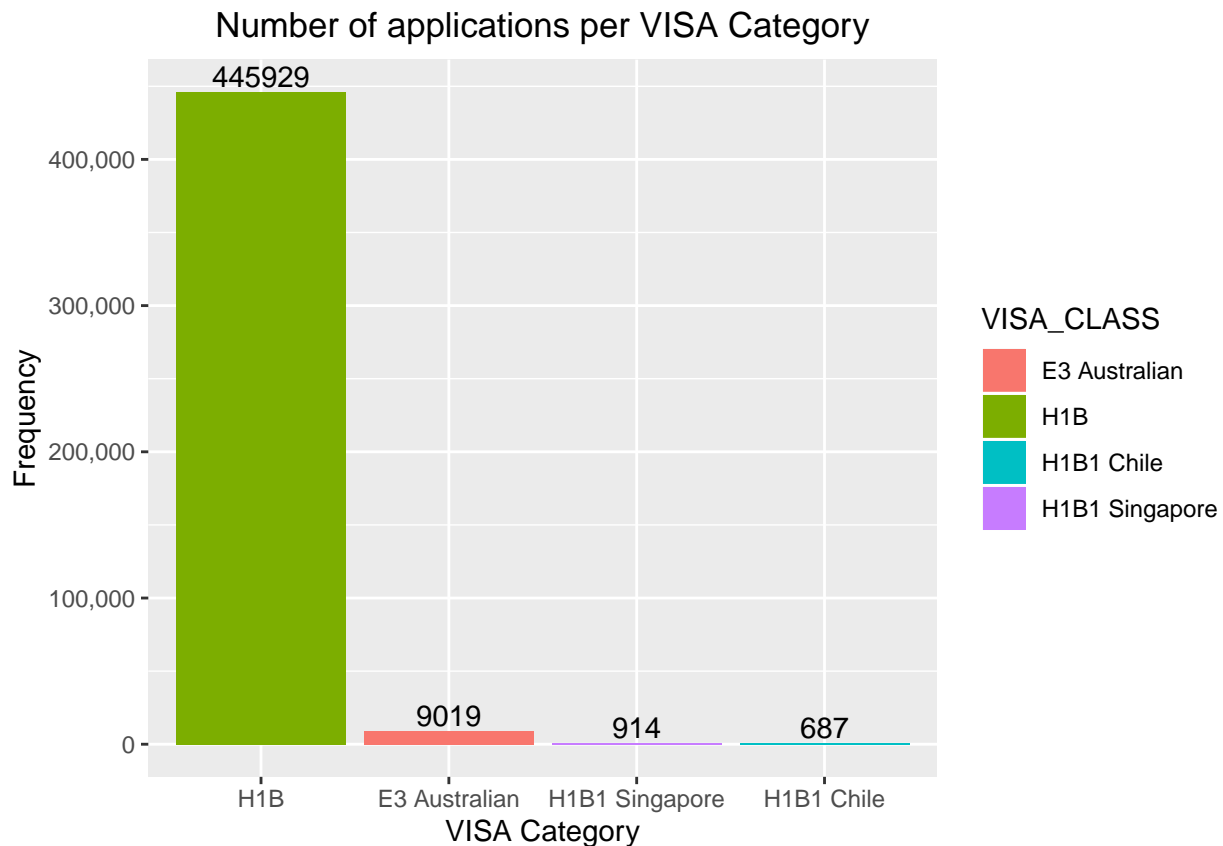
Initially, I have explored the frequency of applications per VISA category. From the below bar graph, it looks like more than 95% of the applications were for H1B visa category with approximately 44K records that belong to the category.

```
visaCategory <- h1bData %>%
  group_by(VISA_CLASS) %>%
```

```

summarize(frequency=n())
(ggplot(visaCategory, aes(x=reorder(VISA_CLASS, -frequency),
                                y=frequency, fill=VISA_CLASS)) +
  geom_bar(stat="identity") +
  scale_y_continuous(breaks = seq(0, 500000, by = 100000), labels = comma) +
  geom_text(aes(label=frequency), position=position_dodge(width=0.9),
            vjust=-0.25) +
  xlab("VISA Category") +
  ylab("Frequency") +
  ggtitle("Number of applications per VISA Category") +
  theme(plot.title = element_text(hjust = 0.5)))

```



H1B Visa exploration

The following shows the top 15 states that had the most H1B applicants. Looks like California had the maximum number of applicants. California is one of the hubs that provide a lot of employment to internationals. It is not a wonder that it is on the top of the list.

Following the horizontal bar graph, the table shows the frequency of applications across the years (2011 to 2017) in the top 15 states. It is clear that the number of applications filed has increased over the years and California has the maximum number of applicants compared to all the states. The increase in number is drastic and it has been increasing over the years.

```

h1bAppln <- h1bData %>%
  filter(VISA_CLASS=="H1B")

h1bTopState <- h1bAppln %>%

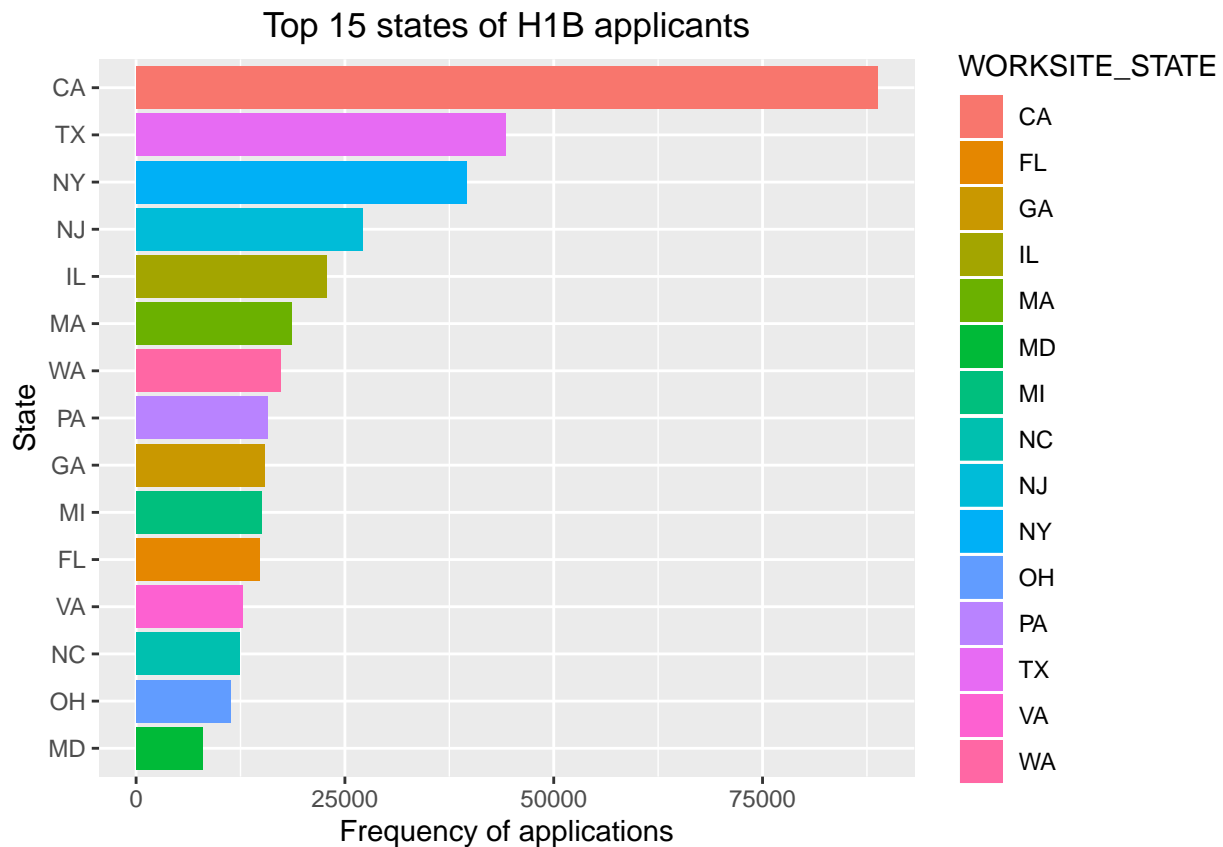
```

```
group_by(WORKSITE_STATE) %>%
summarize(frequency= n()) %>%
arrange(desc(frequency)) %>%
top_n(15)
```

Selecting by frequency

```
(ggplot(h1bTopState,aes(x=reorder(WORKSITE_STATE, frequency),
                             y=frequency, fill=WORKSITE_STATE)) +
  geom_bar(stat="identity") +

  coord_flip() +
  xlab("State") +
  ylab("Frequency of applications") +
  ggtitle("Top 15 states of H1B applicants")+
  theme(plot.title = element_text(hjust = 0.5)))
```



```
h1bTopYear <- h1bAppln %>%
  filter(WORKSITE_STATE %in% h1bTopState$WORKSITE_STATE) %>%
  group_by(WORKSITE_STATE, CASE_SUBMITTED_YEAR) %>%
  summarize(frequency=n())

# Year-wise spread of h1b application with respect to top 15 states
h1bYearSpread <- h1bTopYear %>%
  spread(key=CASE_SUBMITTED_YEAR, value = frequency)
colnames(h1bYearSpread)[1] <- "STATE"
```

```
h1bYearSpread[is.na(h1bYearSpread)] <- 0
h1bYearSpread
```

```
## # A tibble: 15 x 8
## # Groups:   WORKSITE_STATE [15]
##   STATE `2011` `2012` `2013` `2014` `2015` `2016` `2017`
##   <chr> <dbl> <dbl> <int> <int> <int> <int> <int>
## 1 CA      1      6     45   1012   1565  16994  69110
## 2 FL      0      0      5    125    161   2830  11664
## 3 GA      0      0      6    168    215   3098  11867
## 4 IL      0      0     13    189    247   4710  17689
## 5 MA      0      0     14    223    330   3602  14495
## 6 MD      1      0      5    111    183   1697   5928
## 7 MI      0      0      3    109    200   2746  11966
## 8 NC      0      0      5     97    171   2902   9220
## 9 NJ      0      0     14    160    320   5664  21018
## 10 NY     0      0     35   357    483   7021  31619
## 11 OH     0      0      5     90    141   2354   8766
## 12 PA     0      0     13    137    194   3316  12140
## 13 TX     0      2     30   508    742   8550  34363
## 14 VA     0      0      9    131    218   2617   9724
## 15 WA     0      0     31   276    324   4222  12437
```

Now, I am determining the decision status of the applications across the top states. From the vertically stacked bar graph, looks like all the states have more certified cases compared to other decision statuses. After this, I have also determined the acceptance rate of the H1B applications for states, shown in the form of a table. The maximum acceptance rate is for NY state which is around 89.8%. But almost all the top states have an acceptance rate on an average of around 88.5%.

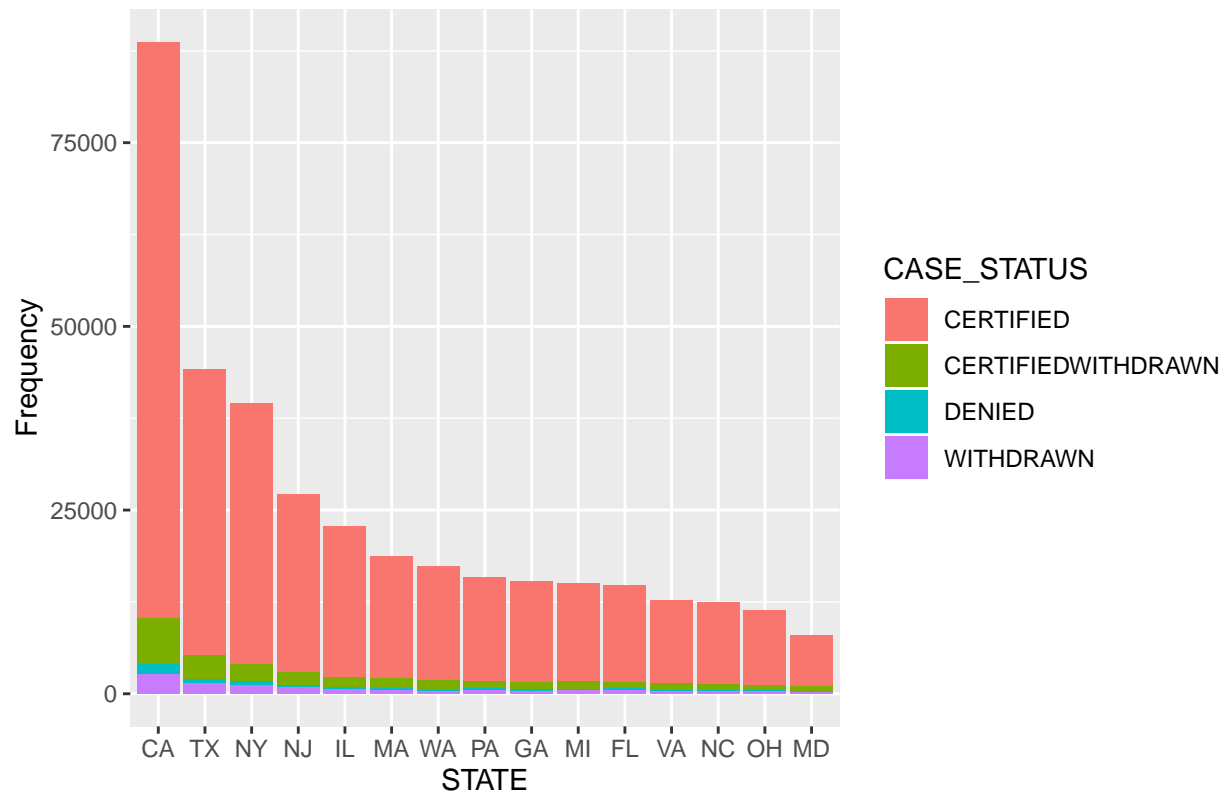
Even though California has the most number of H1B applications, NY has a better acceptance rate than California.

```
# decision with respect to top 15 states
h1bStatus <- h1bAppln %>%
  filter(WORKSITE_STATE %in% h1bTopState$WORKSITE_STATE) %>%
  group_by(WORKSITE_STATE, CASE_STATUS) %>%
  summarize(frequency=n())

(ggplot(h1bStatus, aes(x=reorder(WORKSITE_STATE, -frequency),
                             y=frequency, fill=CASE_STATUS, label=frequency)) +
  geom_bar(stat = "identity") +

  xlab("STATE") +
  ylab("Frequency") +
  ggtitle("Status of H1B applications of top 15 states") +
  theme(plot.title = element_text(hjust = 0.5)))
```

Status of H1B applications of top 15 states



Certified acceptance rate for the top 15 states

```
h1bStateCertified <- h1bAppln %>%
  filter(WORKSITE_STATE %in% h1bTopState$WORKSITE_STATE &
    CASE_STATUS=="CERTIFIED") %>%
  group_by(WORKSITE_STATE) %>%
  summarize(certifiedCases = n())

h1bCertifiedRate <- merge(h1bTopState, h1bStateCertified, by="WORKSITE_STATE")
h1bCertifiedRate$acceptanceRate <-
  h1bCertifiedRate$certifiedCases/h1bCertifiedRate$frequency
h1bCertifiedRate
```

##	WORKSITE_STATE	frequency	certifiedCases	acceptanceRate
## 1	CA	88733	78348	0.8829635
## 2	FL	14785	13097	0.8858302
## 3	GA	15354	13692	0.8917546
## 4	IL	22848	20490	0.8967962
## 5	MA	18664	16476	0.8827690
## 6	MD	7925	6869	0.8667508
## 7	MI	15024	13273	0.8834531
## 8	NC	12395	11073	0.8933441
## 9	NJ	27176	24113	0.8872903
## 10	NY	39515	35481	0.8979122
## 11	OH	11356	10172	0.8957379
## 12	PA	15800	14019	0.8872785
## 13	TX	44195	38900	0.8801901

## 14	VA	12699	11197	0.8817230
## 15	WA	17290	15419	0.8917872

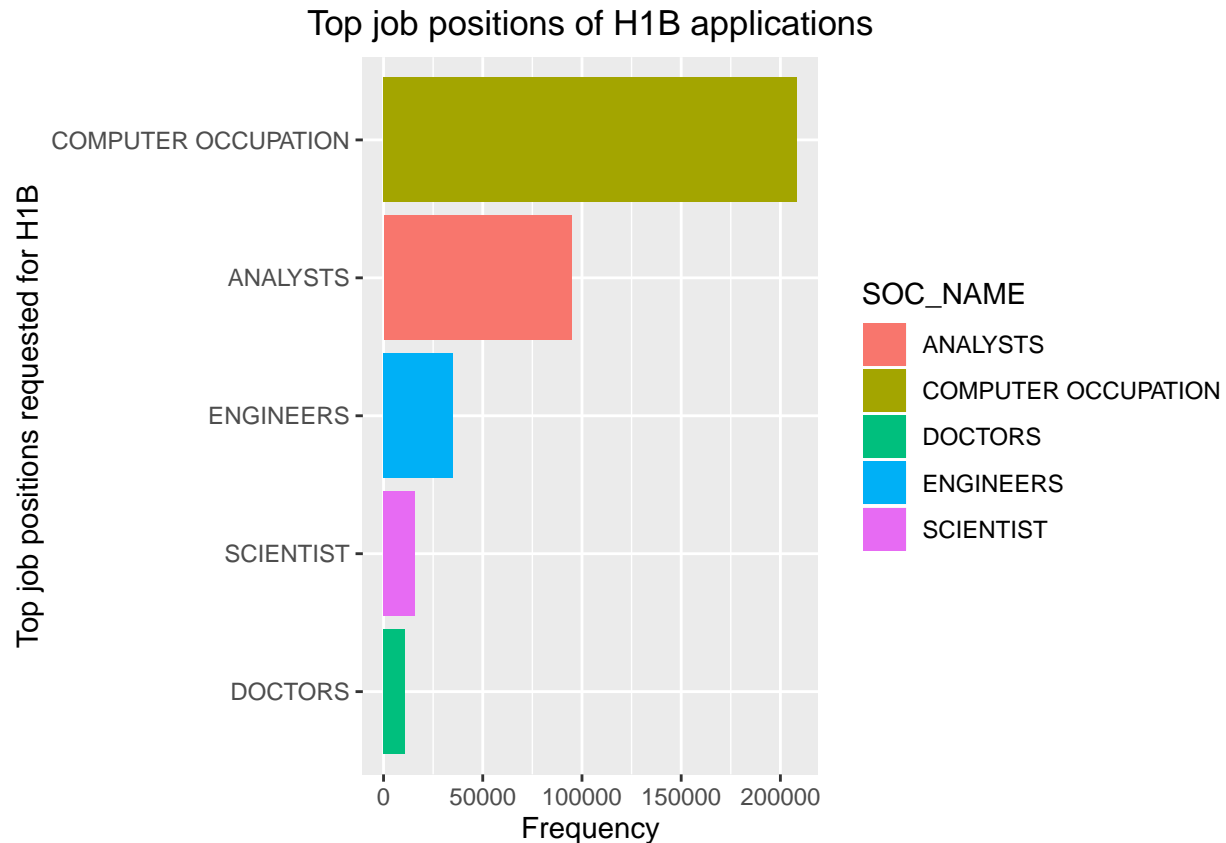
Next look into job positions- initially, I have determined the top five job titles. Looks like more than 200K applications are requested for Computer occupation jobs and the top five jobs are Computer occupation, analysts, engineers, scientists, and doctors.

Now, let's explore how many of these top job positions are requested in the top 15 states. The line graph shows the applicants across the states specific to the top 5 job titles. California, being the top state, has the maximum number of applications with respect to all the job titles as depicted. California and NY are IT hubs in the USA and it is clear that the most number of applications are in California and the most number of applications accepted is in NY and it is also clear from the visualization that California and NY top in Computer occupation jobs. Also, the topmost job title which is Computer Occupation has been leading with respect to all the states, thus showing that computer occupation has the highest demand for all other job titles.

```
# top job positions
h1bTopPositions <- h1bAppln %>%
  group_by(SOC_NAME) %>%
  summarize(frequency=n()) %>%
  arrange(desc(frequency)) %>%
  top_n(5)
```

Selecting by frequency

```
(ggplot(h1bTopPositions, aes(x=reorder(SOC_NAME, frequency),
                                y=frequency, fill=SOC_NAME)) +
  geom_bar(stat="identity") +
  coord_flip() +
  xlab("Top job positions requested for H1B") +
  ylab("Frequency") +
  ggtitle("Top job positions of H1B applications") +
  theme(plot.title = element_text(hjust = 0.5)))
```



Exploring the trends in frequency of the top 5 job titles across the top 15 states

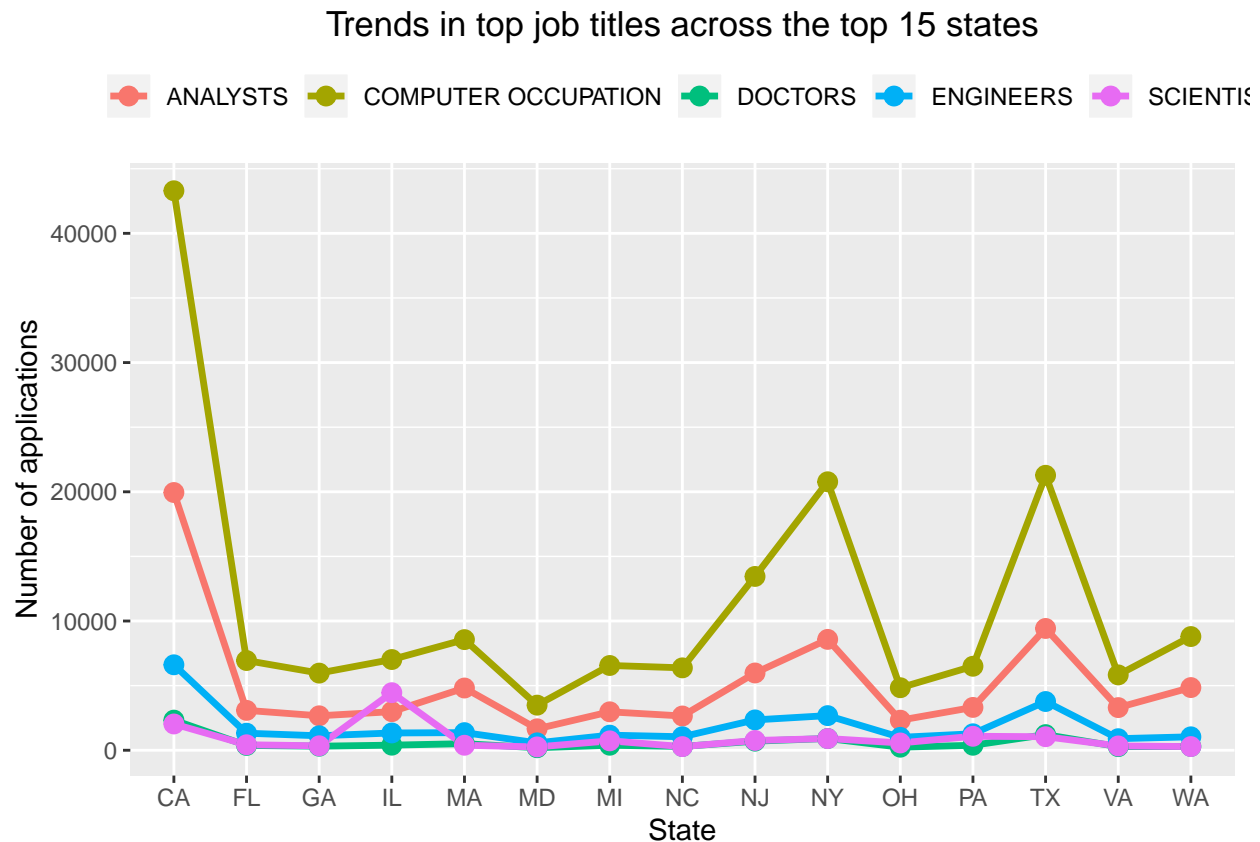
```
h1bStatePosition <- h1bAppln %>%
  filter(SOC_NAME %in% h1bTopPositions$SOC_NAME &
         WORKSITE_STATE %in% h1bTopState$WORKSITE_STATE) %>%
  group_by(WORKSITE_STATE, SOC_NAME) %>%
  summarize(frequency = n())
```

```
h1bJobSpread <- h1bStatePosition %>%
  spread(key=WORKSITE_STATE, value=frequency)
h1bJobSpread
```

```
## # A tibble: 5 x 16
##   SOC_NAME    CA    FL    GA    IL    MA    MD    MI    NC    NJ    NY
##   <chr>    <int> <int> <int> <int> <int> <int> <int> <int> <int> <int>
## 1 ANALYSTS 19944 3087 2666 2986 4816 1660 2974 2652 5979 8579
## 2 COMPUTE~ 43295 6944 5970 7020 8561 3490 6560 6376 13451 20777
## 3 DOCTORS  2332  373  313  394  510  195  389  293  704  910
## 4 ENGINEE~ 6614 1319 1119 1338 1364  572 1169 1045 2348 2680
## 5 SCIENTI~ 2035  425  343 4455  387  259  728  283  756  906
## # ... with 5 more variables: OH <int>, PA <int>, TX <int>, VA <int>,
## #   WA <int>
```

```
(ggplot(data=h1bStatePosition, aes(x=WORKSITE_STATE, y=frequency, group=SOC_NAME)) +
  geom_line(linetype="solid", size=1.2, aes(color=SOC_NAME)) +
  geom_point(aes(color=SOC_NAME), size=3) +
```

```
ggtitle("Trends in top job titles across the top 15 states") +
  xlab("State") +
  ylab("Number of applications") +
  theme(plot.title = element_text(hjust = 0.5),
        legend.position = "top", legend.title = element_blank()))
```



Now, I am exploring the yearly starting salary(wage) of the majoring job titles. The following histogram shows the applicants falling into each of the wage ranges from the lowest to highest wage, across the job titles as depicted by the vertically stacked histogram.

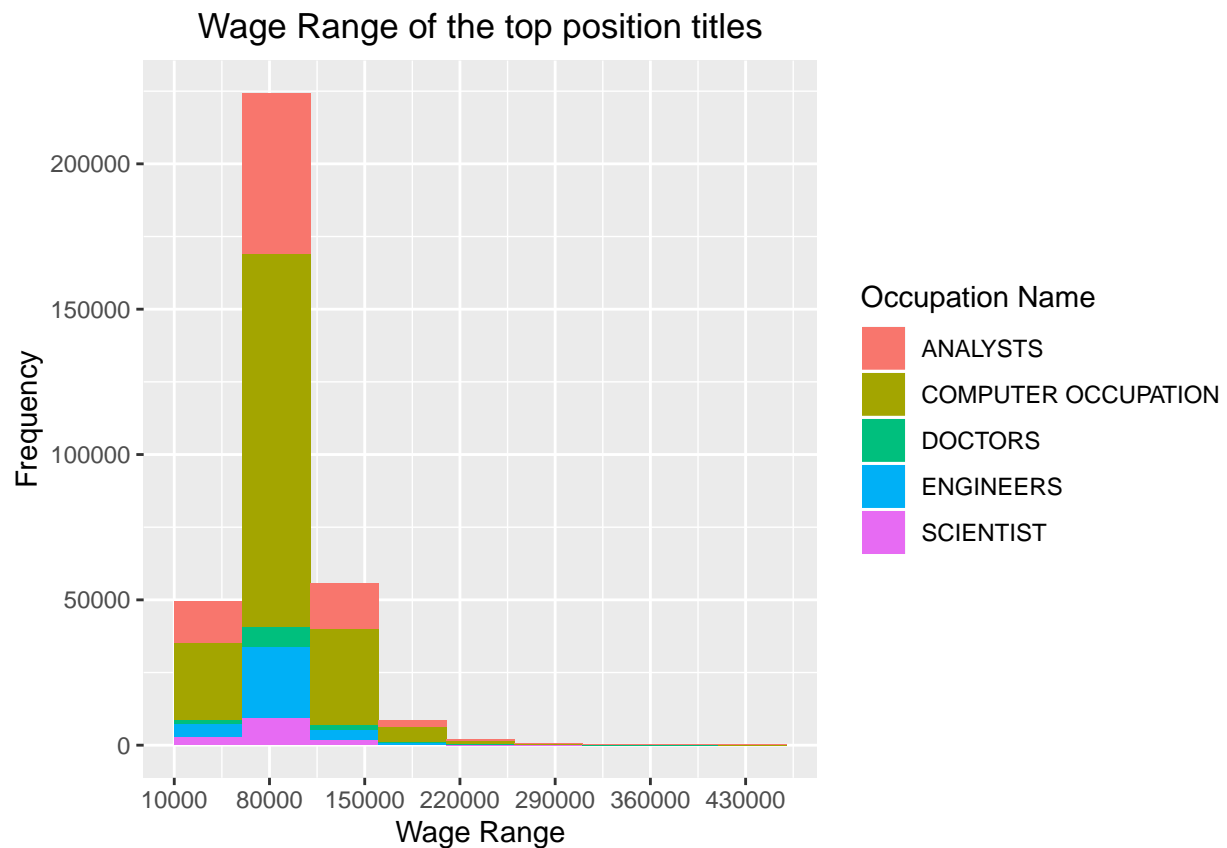
Following that, as salary depends on the state, I have determined the average salary for each of the top job titles across the 15 states. This will give us an idea about the average salary provided by the employers for these jobs with respect to states. Looks like California and Washington has the maximum average salary across all the job titles. The reason for such a pattern could be because the cost of living is expensive in California and Washington. As a resident of NY, we know that the cost of living and the taxes are a little high (and from the graph) but seem like it is not as high as California and Washington.

```
# wageRange of top positions

h1bTopPosAppl <- h1bAppln %>%
  filter(SOC_NAME %in% h1bTopPositions$SOC_NAME & WAGE_UNIT_OF_PAY=="Year")

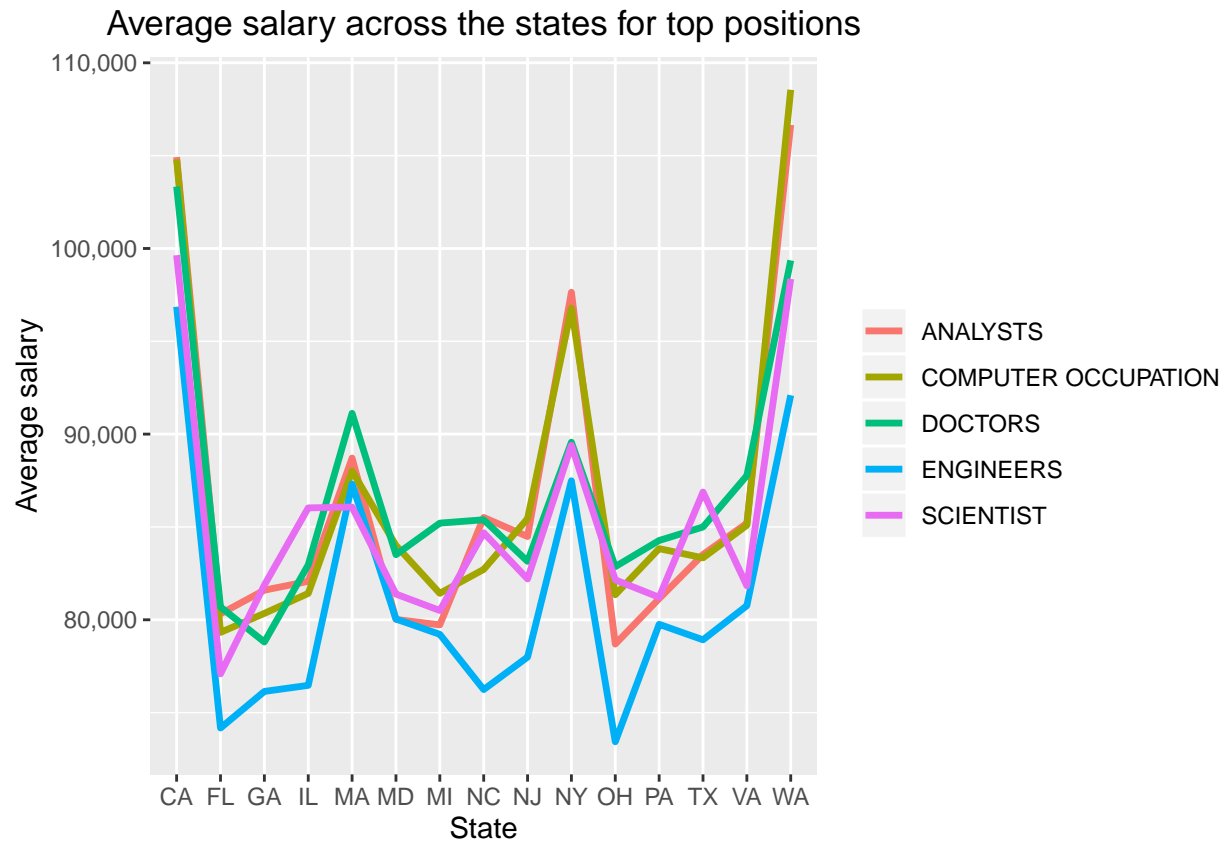
(ggplot(data=h1bTopPosAppl, aes(x=WAGE_RATE_OF_PAY_FROM)) +
  geom_histogram(aes(fill=SOC_NAME), breaks=seq(10000, 500000, by=50000)) +
  scale_x_continuous(breaks = seq(10000, 500000, by=70000)) +
  ggtitle("Wage Range of the top position titles") +
  xlab("Wage Range") +
```

```
ylab("Frequency") +
guides(fill=guide_legend(title="Occupation Name"))+
theme(plot.title = element_text(hjust = 0.5))
```



```
# Average yearly starting salary in the top states with respect to top positions
h1bStateAvgSalary <-h1bAppln %>%
  filter(WAGE_UNIT_OF_PAY=="Year" &
    SOC_NAME %in% h1bTopPositions$SOC_NAME &
    WORKSITE_STATE %in% h1bTopState$WORKSITE_STATE) %>%
  group_by(WORKSITE_STATE,SOC_NAME) %>%
  summarize(`Average Salary` = mean(WAGE_RATE_OF_PAY_FROM))

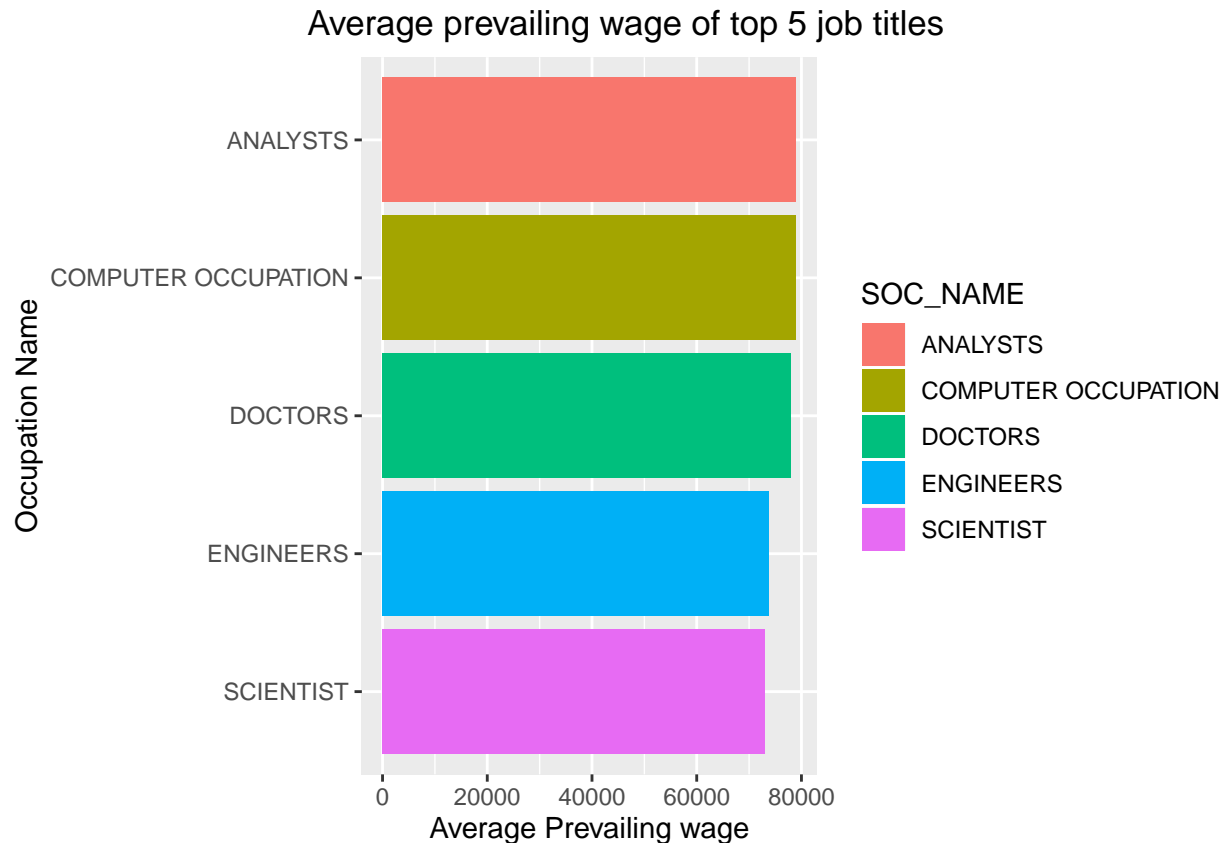
# plot with state and average salary with respect to job title
(ggplot(data=h1bStateAvgSalary, aes(x=WORKSITE_STATE, y=`Average Salary`, group=SOC_NAME))
+ geom_line(linetype="solid", size=1.2, aes(color=SOC_NAME)) +
  ggtitle("Average salary across the states for top positions") +
  xlab("State") +
  ylab("Average salary") +
  scale_y_continuous(labels = comma) +
  theme(plot.title = element_text(hjust = 0.5), legend.title = element_blank()))
```



Having explored the average salary, now I am exploring the average prevailing(current) wage for the top 5 jobs(analysts, computer occupation, doctors, scientists, and engineers). Looks like analysts and computer occupations have almost the similar average prevailing wage. This gives us an idea of what is the current average salary for the top job positions. This graph can also help students to get an idea of their market value and think wisely when negotiation salary when they get a job offer.

```
# prevailing wage for top jobs
h1bPrevailingWage <- h1bAppln %>%
  filter(WAGE_UNIT_OF_PAY=="Year" & SOC_NAME %in% h1bTopPositions$SOC_NAME) %>%
  group_by(SOC_NAME) %>%
  summarize(`Average Prevailing wage`=mean(PREVAILING_WAGE))

(ggplot(h1bPrevailingWage, aes(x=reorder(SOC_NAME, `Average Prevailing wage`),
  y=`Average Prevailing wage`, fill=SOC_NAME)) +
  geom_bar(stat="identity") +
  xlab("Occupation Name") +
  coord_flip() +
  ggtitle("Average prevailing wage of top 5 job titles")+
  theme(plot.title = element_text(hjust = 0.3)))
```



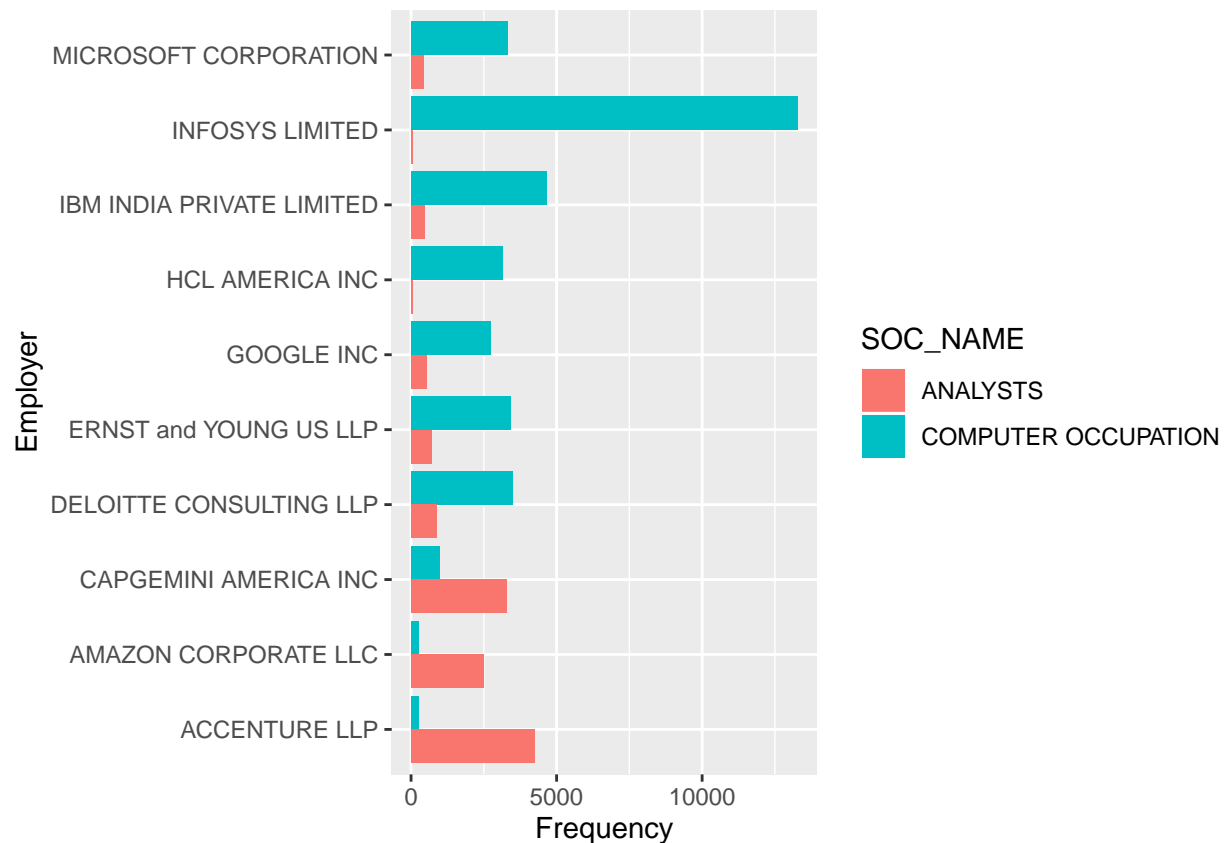
Having explored the wages, now let's find the top 10 employers who have filed H1B for computer programmers and analysts (being the top 2 jobs). This gives us an idea of the top employers sponsoring H1B with a breakdown of both analysts and computer occupation. Looks like, Infosys is majoring in sponsoring computer occupation and Accenture is majoring in sponsoring analysts.

```
# the top employers offering computer occupation and analysts jobs
topEmployers <- h1bAppln %>%
  filter(SOC_NAME=="COMPUTER OCCUPATION" | SOC_NAME=="ANALYSTS") %>%
  group_by(EMPLOYER_NAME) %>%
  summarize(frequency=n()) %>%
  arrange(desc(frequency)) %>%
  top_n(10)
```

Selecting by frequency

```
employerOcc <- h1bAppln %>%
  filter(EMPLOYER_NAME %in% topEmployers$EMPLOYER_NAME &
         (SOC_NAME=="COMPUTER OCCUPATION" | SOC_NAME=="ANALYSTS"))

(ggplot(data = employerOcc) +
  geom_bar(mapping = aes(x=EMPLOYER_NAME,
                        fill=SOC_NAME), position = "dodge") +
  coord_flip() +
  xlab("Employer") +
  ylab("Frequency"))
```



CONCLUSION

In this document, I have made the best use of H1B application data showing various visual explorations using the ggplot2 library. These explorations would be useful for those filing h1b applications and also the current applicants, as it gives us an overall idea of which states have more acceptance rates, the most demanding jobs, and the top employers sponsoring H1B visas for the non-immigrants. To conclude, I found that California is one of the states that has the top-notch tech companies and hence they hire the most. I also could see that Computer Occupation and Analysts have the best average salary. The H1b acceptance rate is high in NY but the number of applications is the highest in California. Similarly, I could also see that California and Washington have the highest paying jobs followed by NY. On the other hand, as the world is turning out to be digital, the most demanding job has become computer software. I feel this trend is likely to be seen in the following years as well with the other jobs been replaced by Computer occupation.