

Compare missing versus non-missing Statistics*

Deal and Lawyer Studies

Li Yuan, li.yuan@vanderbilt.edu

26 July, 2021 17:12:46

Abstract

The aim of this paper is to study correlation between missing information and industry sector, year, age, deal value. We will focus on how gender affects missing info using machine learning inference in the second half of this paper.

Contents

Load packages	2
Load data	2
How many lawyers are unknown (NA and “not disclosed”)?	2
Distribution of the missing and non-missing lawyers	2
Chi-Square Test of Independence	3
Two-sample Kolmogorov–Smirnov test	8
How many law firms are unknown (NA and “not disclosed”)?	11
What is the distribution of the missing law firms by the size (in dollars) of the deal, year, and industry? Histogram each.	11
How many deal attorneys are missing biographical information?	13
What is the distribution of the missing data – are these mostly from earlier years, for example?	13
Load the no-matched data set	13
attorneys without biographical info	13
Distribution of Years	14
Distribution of deal type	14
Distribution of attorneys with biograh info	14
Compare them together	14
Two time series by months of counts of attorneys with and without biograh info	15
How many deals are missing attorneys’ biographical information?	16
Counts of deals without attorneys’ biograh info grouped by Year and Deal Type.	16
Create some summary tables	16
Number of deals, gender, age	16
school and law distribution	17
types of deals	19
Buyer and seller of deals	20
Regression Analysis	21
First we used Multiple Logistic Regression with all variables	21
Multiple Logistic Regression with all variables except for Law Schol and Law firm	21
Random Forest with the whole variables used	21
Conclusion	21

*Advisor: Tracey George and Albert

Load packages

```
1 > library(readtext)
2 > library(antiword)
3 > library(tidyverse)
4 > library(ggplot2)
5 > library(textreadr)
6 > library(stringi)
7 > library(textclean)
8 > library(SemNetCleaner)
9 > library(readxl)
10 >
11 > library(patchwork)
12 > library(ggrepel)
13 > library(gghighlight)
14 > library(paletteer)
15 > library(ggExtra)
16 > library(ggbeeswarm)
17 > library(kableExtra)
18 > library(caret)
19 > library(randomForest)
20 > library(corrplot)
21 > library(lubridate)
```

Load data

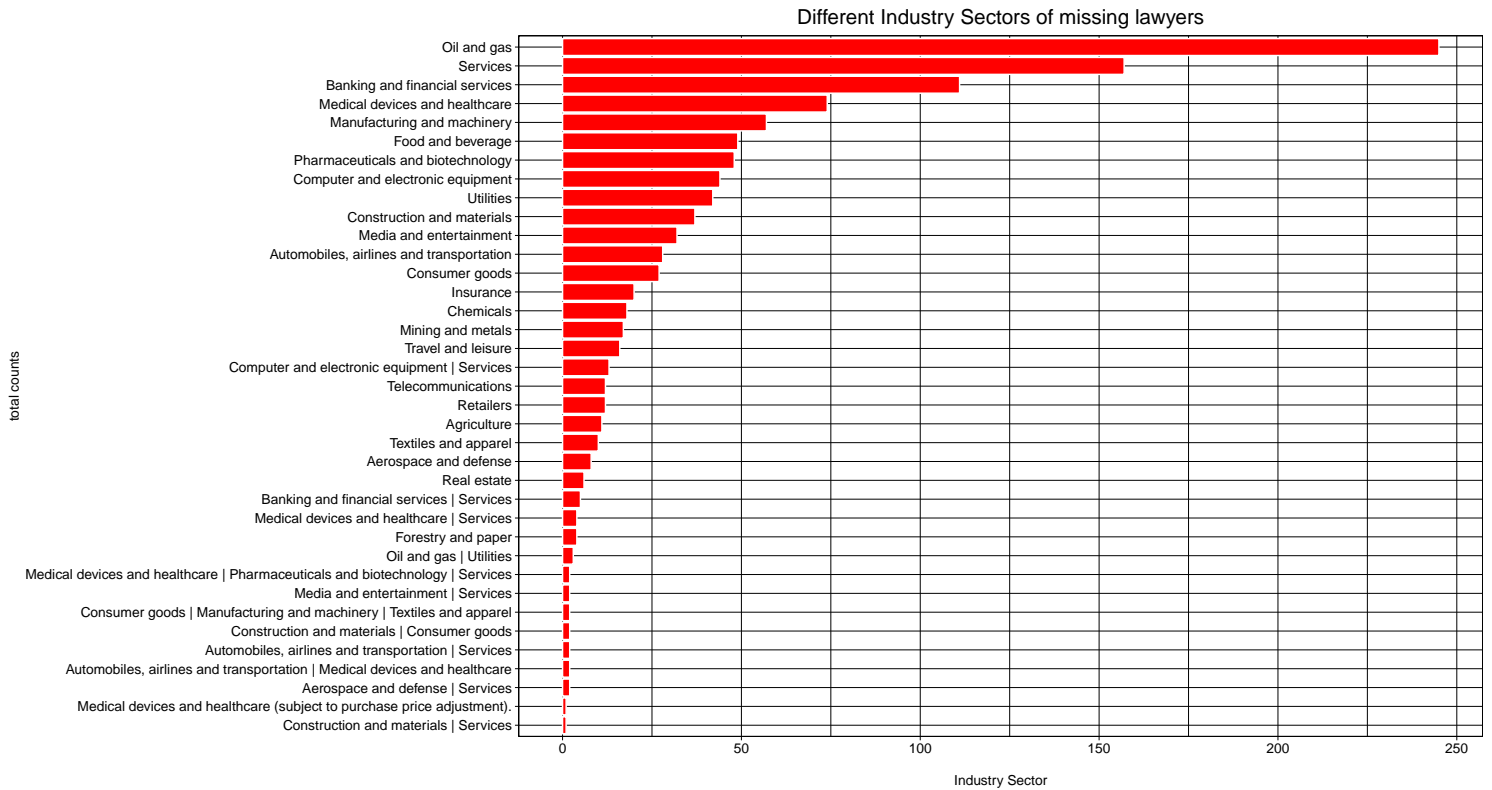
```
1 > deal <- read_csv("../data/deal/deal(1).csv", col_types = cols(.default = "c"))
2 > # View(deal)
3 >
4 > merge_deal_lawyer <- read_csv("../data/deal_lawyer/merge_deal_lawyer.csv")
5 > # View(merge_deal_lawyer)
6 >
7 > distin_com_lawyer <- read_csv("../data/lawyer/keep_MA_first.csv", col_types = cols(.default = "c"))
8 > # View(distin_com_lawyer)
9 >
10 > map_index <- read_csv("../data/deal/map_index.csv")
11 > # View(map_index)
12 >
13 > encoded_merge_dl <- merge_deal_lawyer %>%
14 +   left_join(map_index, by = c(deal = "Deal_name")) %>%
15 +   select(Deal_number, everything(), -deal)
16 > # View(encoded_merge_dl)
17 >
18 > encoded_deal <- deal %>%
19 +   left_join(map_index, by = c(`Deal name` = "Deal_name")) %>%
20 +   select(Deal_number, everything(), -`Deal name`)
21 > # View(encoded_deal)
```

How many lawyers are unknown (NA and “not disclosed”)?

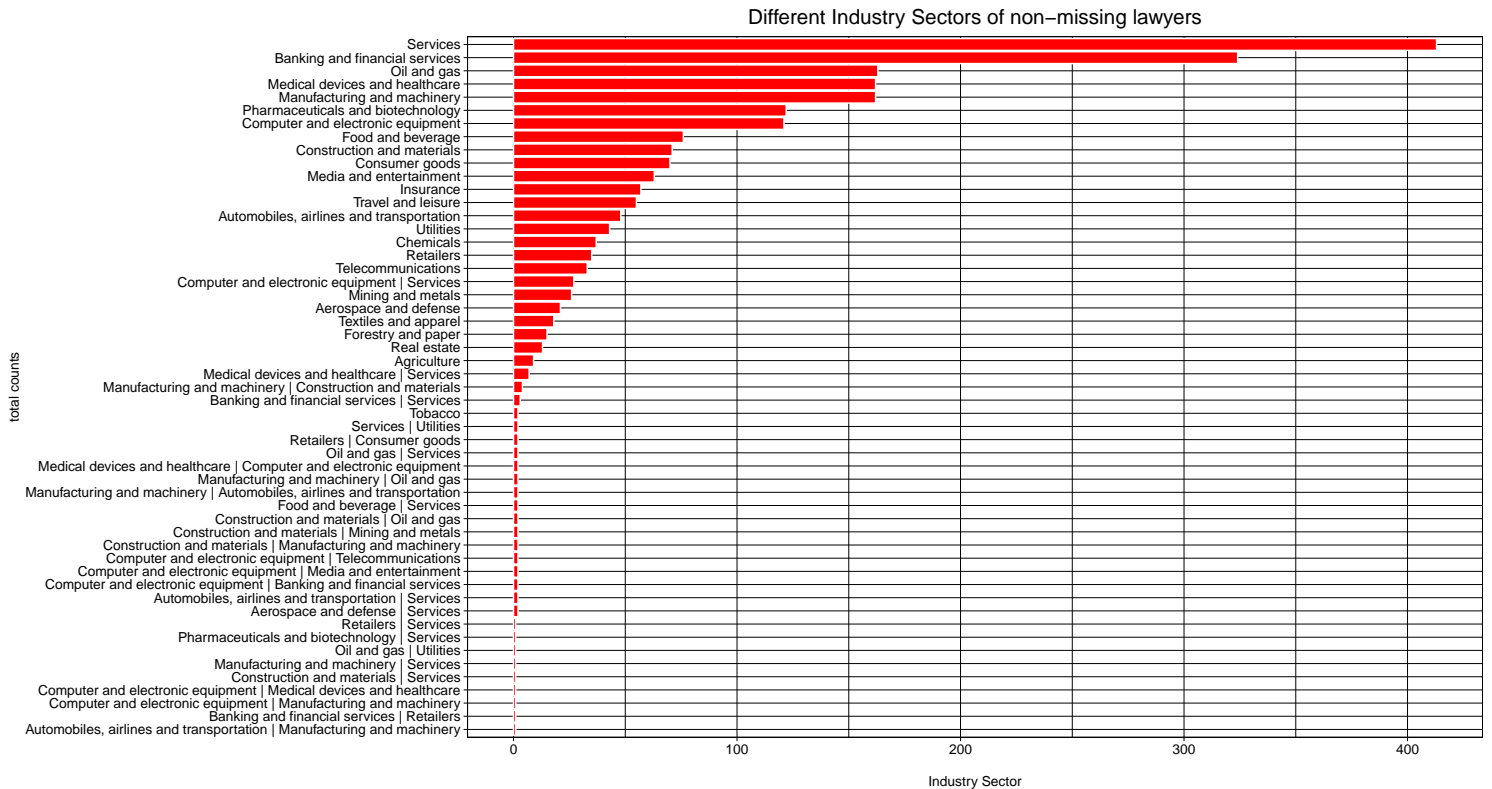
Distribution of the missing and non-missing lawyers

What is the distribution of the missing lawyers by deal value (category), year, and industry sector? Histogram for each plus mean, standard, median.

Counts of Missing Lawyers



Counts of non-missing lawyers



Chi-Square Test of Independence

The top 15 most industry sector among missing lawyer deals

Selecting by n

```

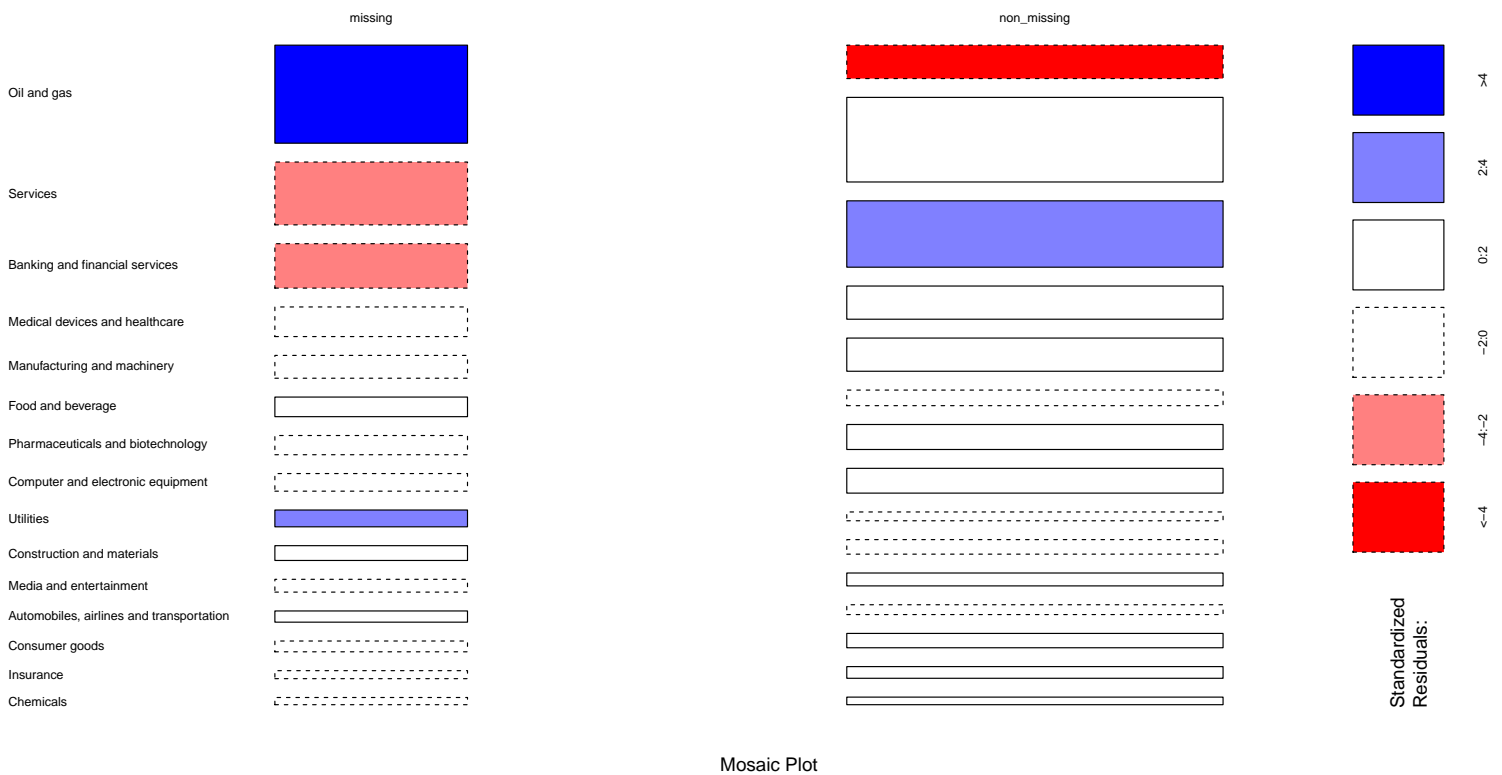
1 > contingency %>%
2 +   kbl(caption = "Contingency Table of missing and non-missing lawyer deals by 15 Industry",
3 +       booktabs = T) %>%
4 +   kable_styling(latex_options = c("striped", "hold_position"))

```

Table 1: Contingency Table of missing and non-missing lawyer deals by 15 Industry

	missing	non_missing
Oil and gas	245	163
Services	157	413
Banking and financial services	111	324
Medical devices and healthcare	74	162
Manufacturing and machinery	57	162
Food and beverage	49	76
Pharmaceuticals and biotechnology	48	122
Computer and electronic equipment	44	121
Utilities	42	43
Construction and materials	37	71
Media and entertainment	32	63
Automobiles, airlines and transportation	28	48
Consumer goods	27	70
Insurance	20	57
Chemicals	18	37

Pearson's chi-squared of missing and non-missing lawyer deals by Industry



- Blue color indicates that the observed value is higher than the expected value if the data were random
- Red color specifies that the observed value is lower than the expected value if the data were random

From this plot generated by **Pearson's chi-squared**, we can tell that **oil and gas** has more missing lawyers than expected while **oil and gas** has much lower deals of non-missing lawyers than expected.

```

1 > chisq <- chisq.test(contingency)
2 > chisq

```

```

##
## Pearson's Chi-squared test
##
## data:  contingency
## X-squared = 176.33, df = 14, p-value < 2.2e-16

```

From this $p - value \approx 0$, we can tell that industry sector are missing versus non-missing are statistically significantly associated.

If we want to know the most contributing cells to the total Chi-square score, you just have to calculate the Chi-square statistic for each cell:

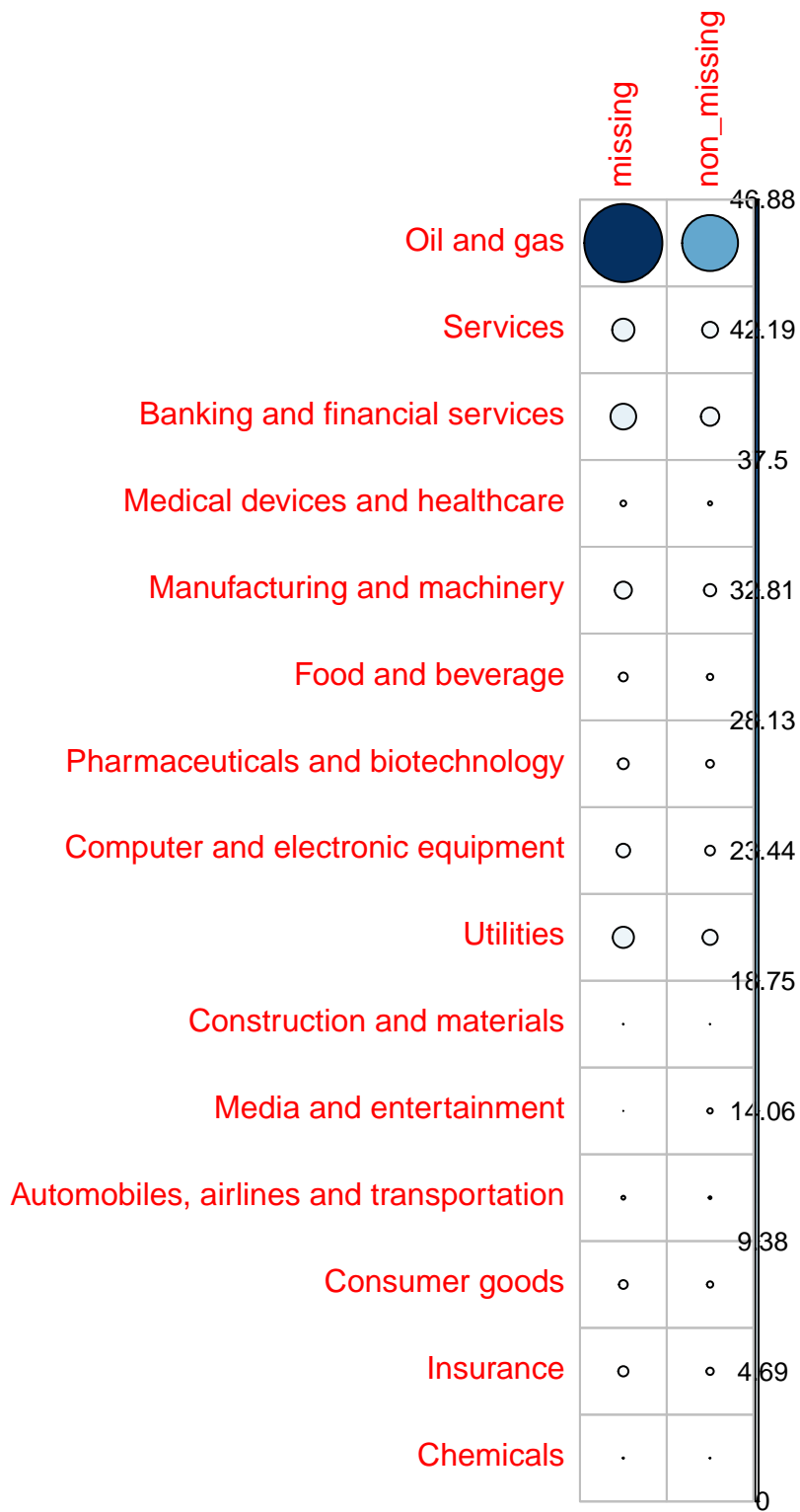
$$r = \frac{o - e}{\sqrt{e}}$$

$$contrib = \frac{r^2}{\chi^2}$$

```

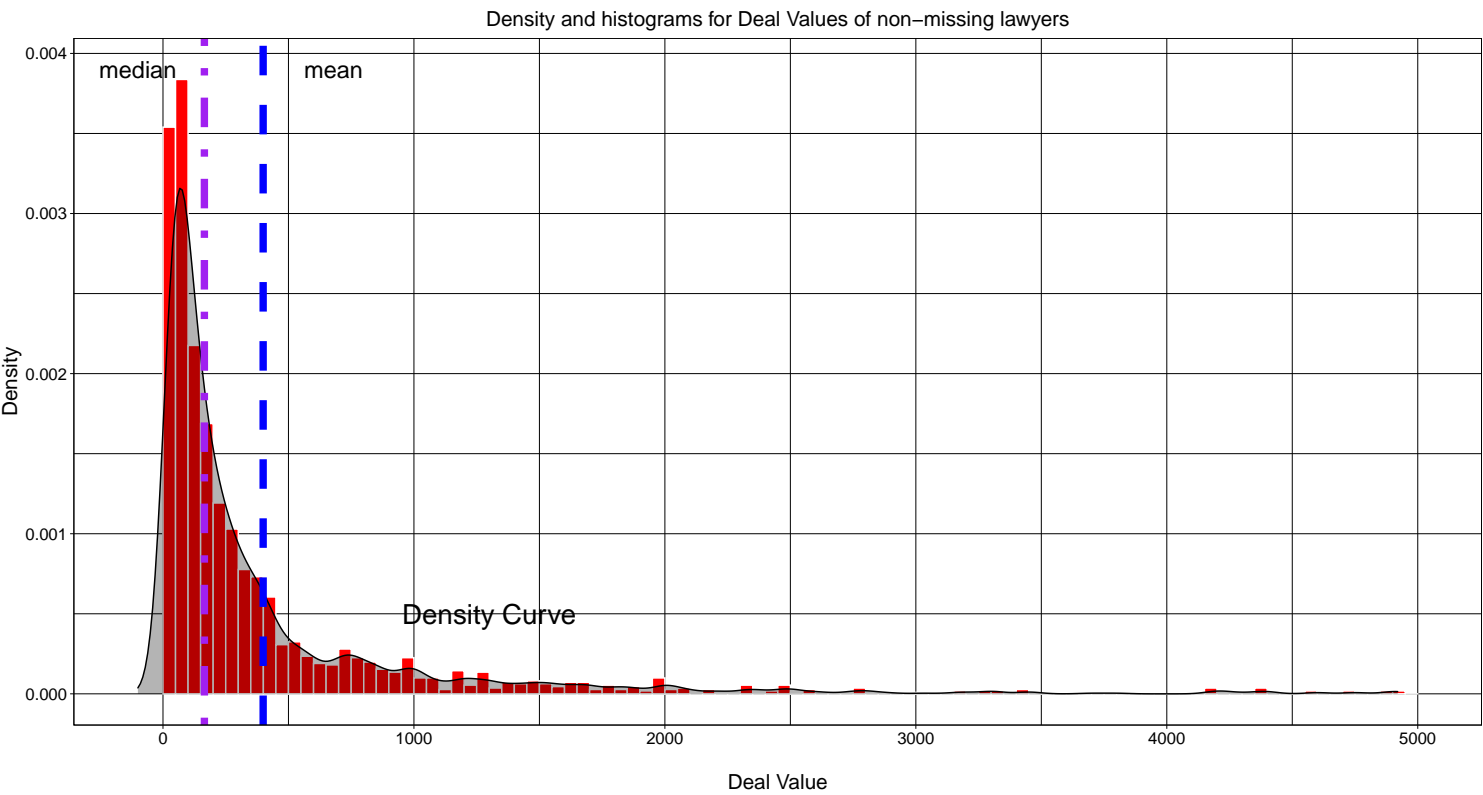
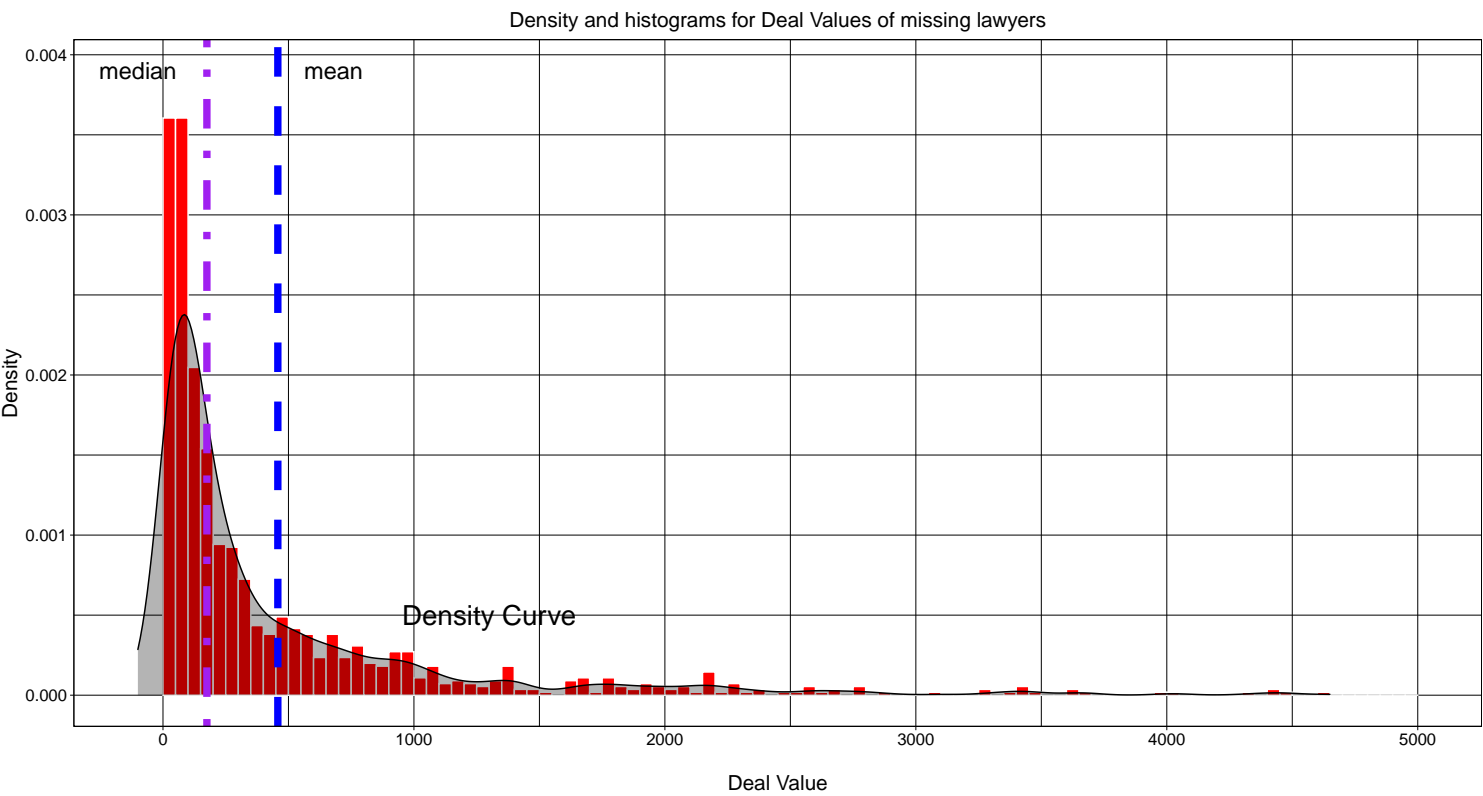
1 > # Visualize the contribution
2 > corrplot(contrib, is.corr = FALSE, outline = T)

```

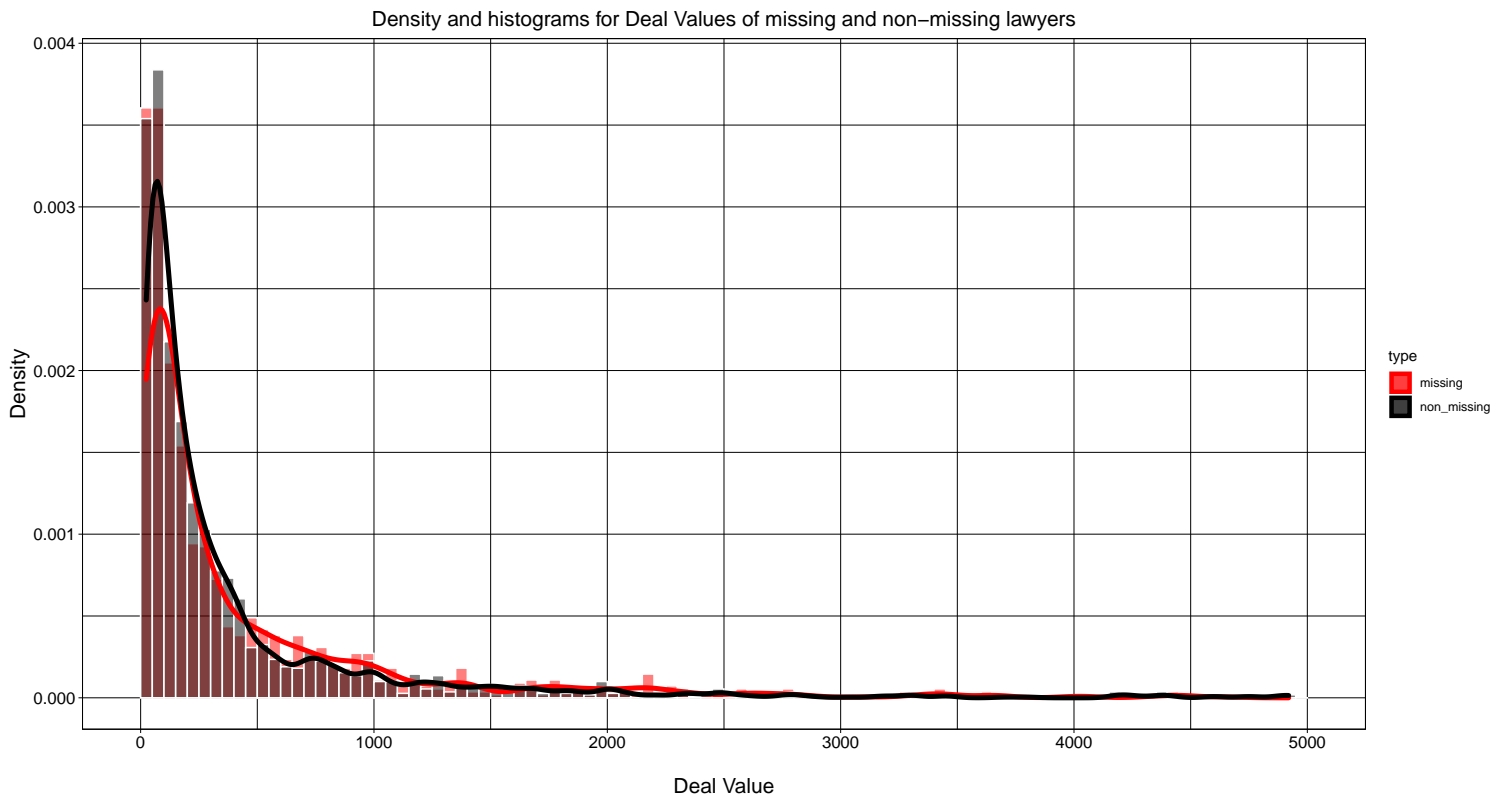


From this percentage contribution correlation plot, we can tell that oil and gas makes huge statistically significant to chi-squared values, which means that oil and gas industry sector is a main factor to determine missing versus non-missing lawyer in deals.

Distribution and Density curve of deal values of missing lawyers



Put them together to compare



Two-sample Kolmogorov–Smirnov test

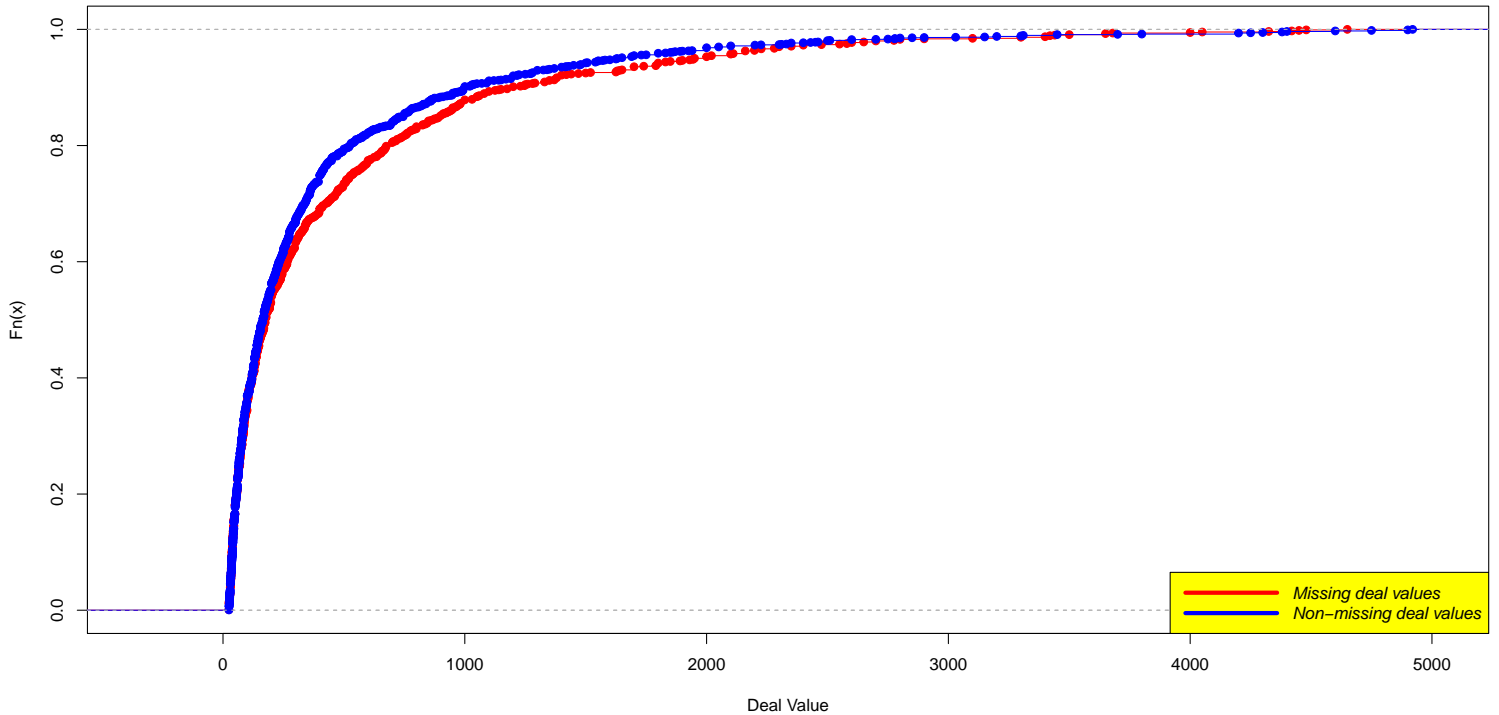
Empirical distributions of deal values for missing and non-missing lawyer deals

```

1 > plot(ecdf(dat %>%
2 +   filter(type == "missing") %>%
3 +   pull(value1)), col = "red", main = "Empirical Distribution Functions of missing and non-missing deal values",
4 +   xlab = "Deal Value")
5 > lines(ecdf(dat %>%
6 +   filter(type == "non_missing") %>%
7 +   pull(value1)), col = "blue")
8 > legend(x = "bottomright", legend = c("Missing deal values", "Non-missing deal values"),
9 +   col = c("red", "blue"), lty = c(1, 1), lwd = 4, bg = "yellow", seg.len = 6, text.font = 3)

```


Empirical Distribution Functions of missing and non-missing deal values



Two-sample Kolmogorov–Smirnov test

The empirical distribution function F_n for n independent and identically distributed (i.i.d.) ordered observations X_i is defined as

$$F_n(x) = \frac{1}{n} \sum_{i=1}^n I_{[-\infty, x]}(X_i)$$

where $I_{[-\infty, x]}(X_i)$ is the indicator function, equal to 1 if $X_i \leq x$ and equal to 0 otherwise.

The Kolmogorov–Smirnov test may also be used to test whether two underlying one-dimensional probability distributions differ. In this case, the Kolmogorov–Smirnov statistic is

$$D_{n,m} = \sup_x |F_{1,n}(x) - F_{2,m}(x)|$$

where $F_{1,n}$ and $F_{2,m}$ are the empirical distribution functions of the first and the second sample respectively, and \sup is the supremum function.

```
1 > miss <- dat %>%
2 +   filter(type == "missing") %>%
3 +   pull(value1)
4 > non_miss <- dat %>%
5 +   filter(type == "non_missing") %>%
6 +   pull(value1)
7 > ks.test(miss, non_miss)
```

```
## Warning in ks.test(miss, non_miss): p-value will be approximate in the presence
## of ties
```

```
##
## Two-sample Kolmogorov-Smirnov test
##
## data: miss and non_miss
## D = 0.069539, p-value = 0.001608
## alternative hypothesis: two-sided
```

From this p-value, 0.001608, we can say that missing and non-missing lawyer deal values are from different distribution since the null hypothesis is two samples are drawn from the same distribution. However, this p-value is not too small and from the ECDF plot, they are quite close, in light of p-value isn't always reliable, so I think the distribution of missing lawyer deal values are almost very close that of non-missing counterpart.

Mean, Median and Standard Deviation of deal values

```

1 > miss_non_miss %>%
2 +   group_by(type) %>%
3 +   summarise(`:=`(min, min(value1, na.rm = T)), `:=`(mean, mean(value1, na.rm = T)),
4 +             `:=`(median, median(value1, na.rm = T)), `:=`(max, max(value1, na.rm = T)),
5 +             `:=`(sd, sd(value1, na.rm = T))) %>%
6 +   kbl(caption = "Summary Stats Table comparing missing and non-missing", booktabs = T) %>%
7 +   kable_styling(latex_options = c("striped", "hold_position"))

```

Table 2: Summary Stats Table comparing missing and non-missing

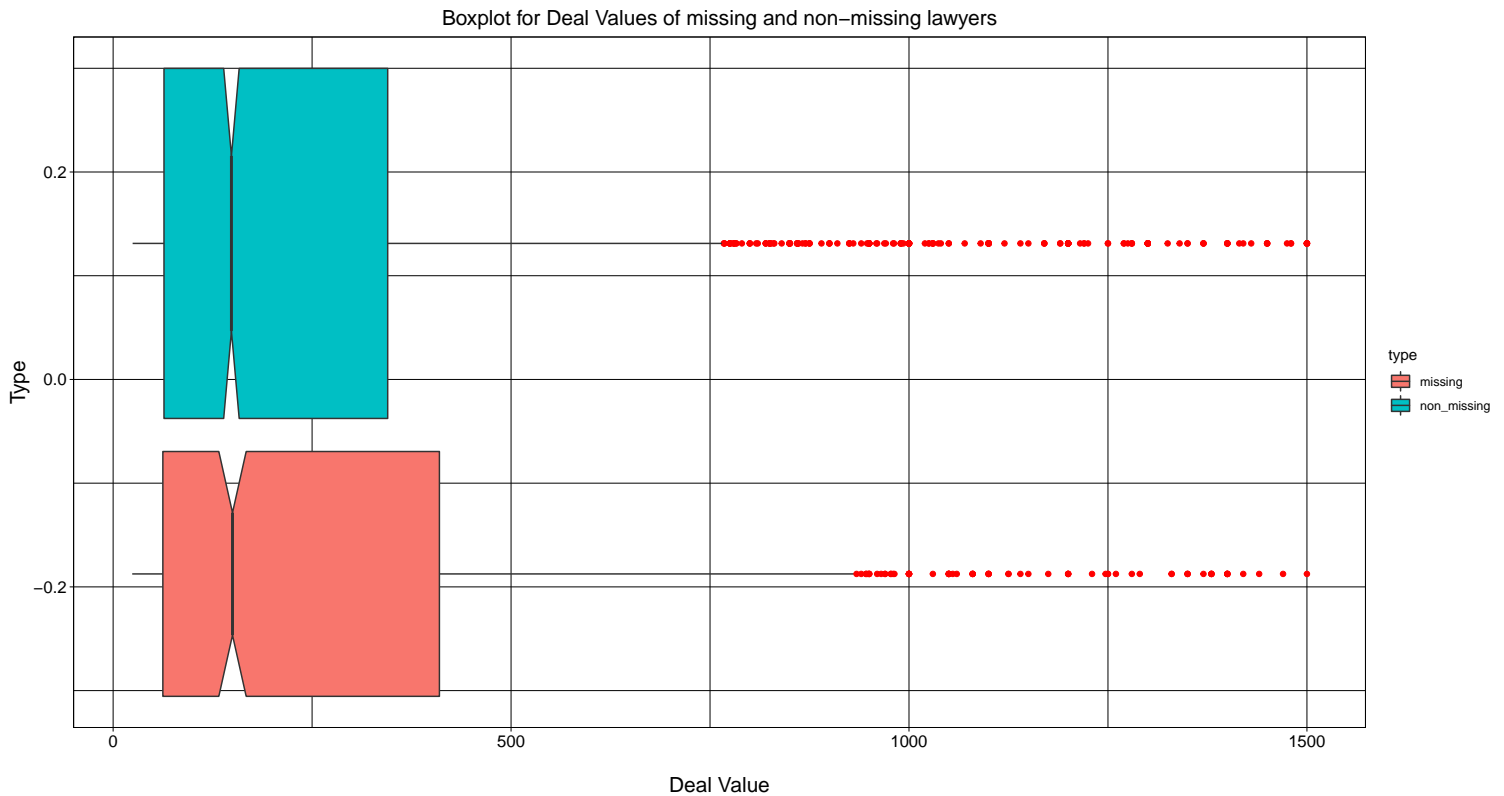
type	min	mean	median	max	sd
missing	24.0	687.6272	180	39000	2099.451
non_missing	24.5	494.4269	169	32700	1308.376

Boxplot for missing and non-missing

```

1 > miss_non_miss %>%
2 +   filter(value1 <= 1500) %>%
3 +   ggplot(aes(x = value1)) + geom_boxplot(aes(fill = type), na.rm = F, notch = T,
4 +   varwidth = T, orientation = "y", outlier.colour = "red") + theme_linedraw() +
5 +   labs(x = "\nDeal Value", y = "Type", title = "Boxplot for Deal Values of missing and non-missing lawyers")
6 +   theme(axis.text.x = element_text(size = 12), axis.text.y = element_text(size = 12),
7 +   axis.title = element_text(size = 15), plot.title = element_text(hjust = 0.5,
8 +   size = 15))

```

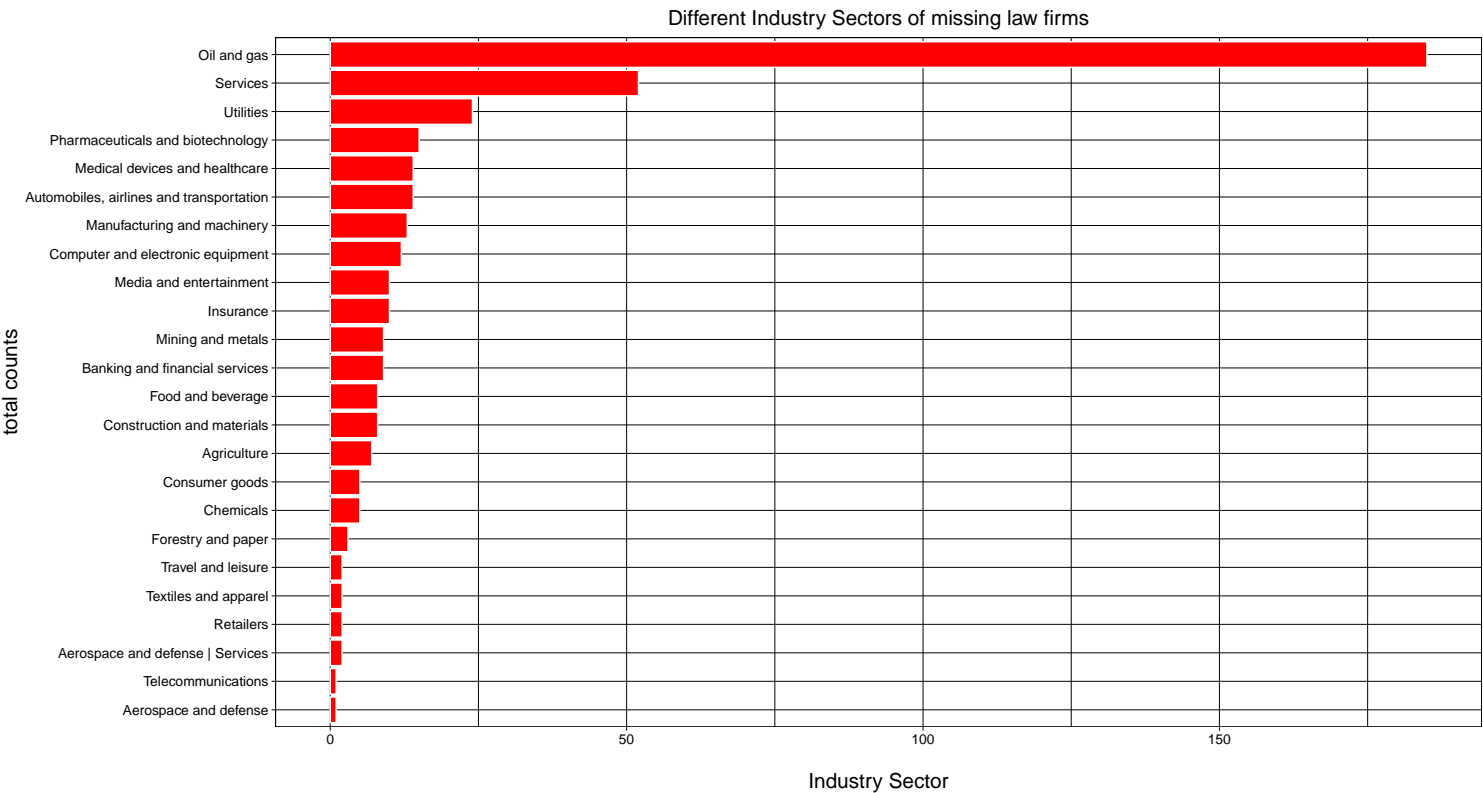


How many law firms are unknown (NA and “not disclosed”)?

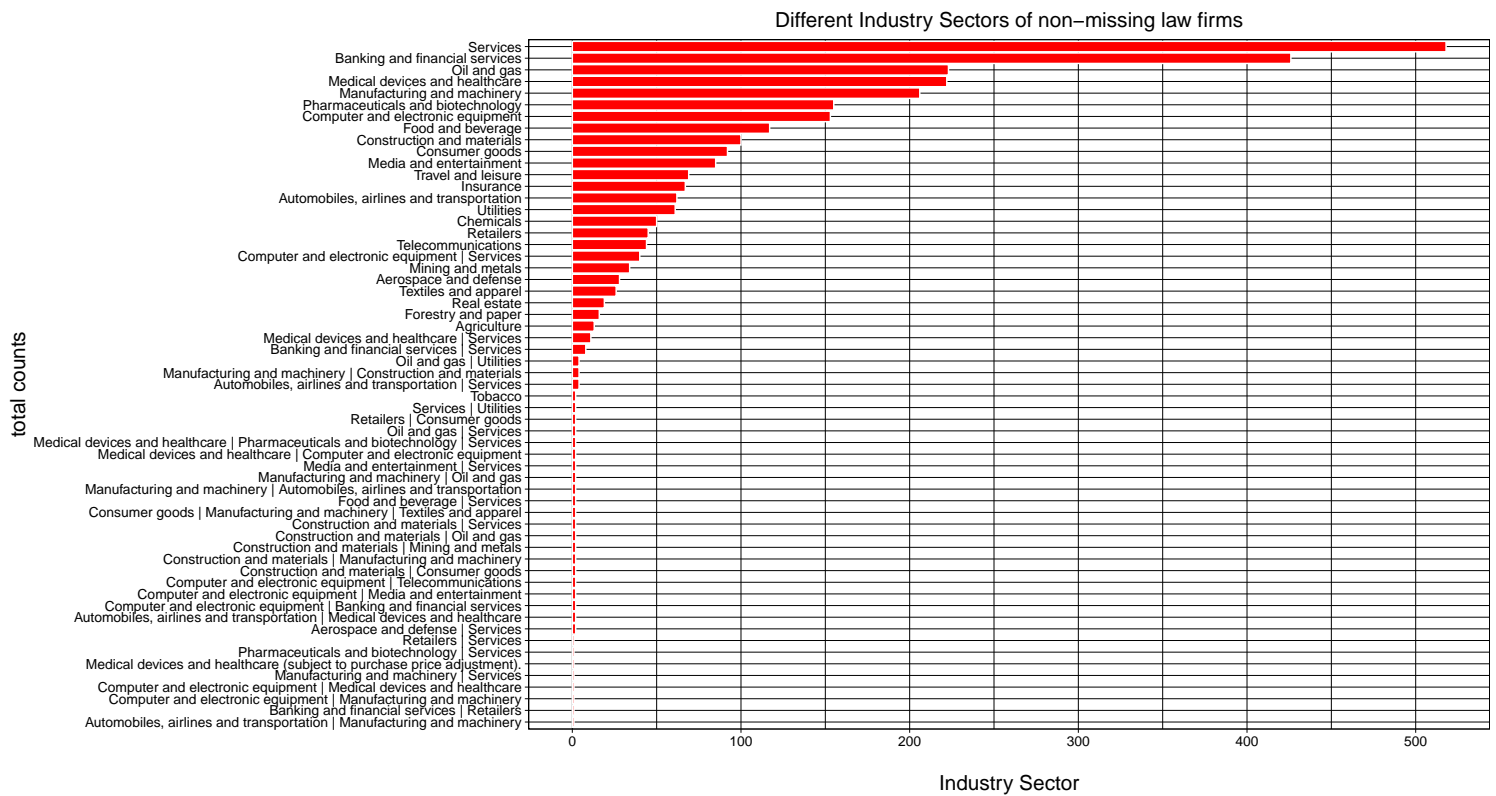
What is the distribution of the missing law firms by the size (in dollars) of the deal, year, and industry? Histogram each.

Extract corresponding dataset

Bar plots for missing law firms



Bar plots for non_missing law firms



Histogram and density curves to compare them

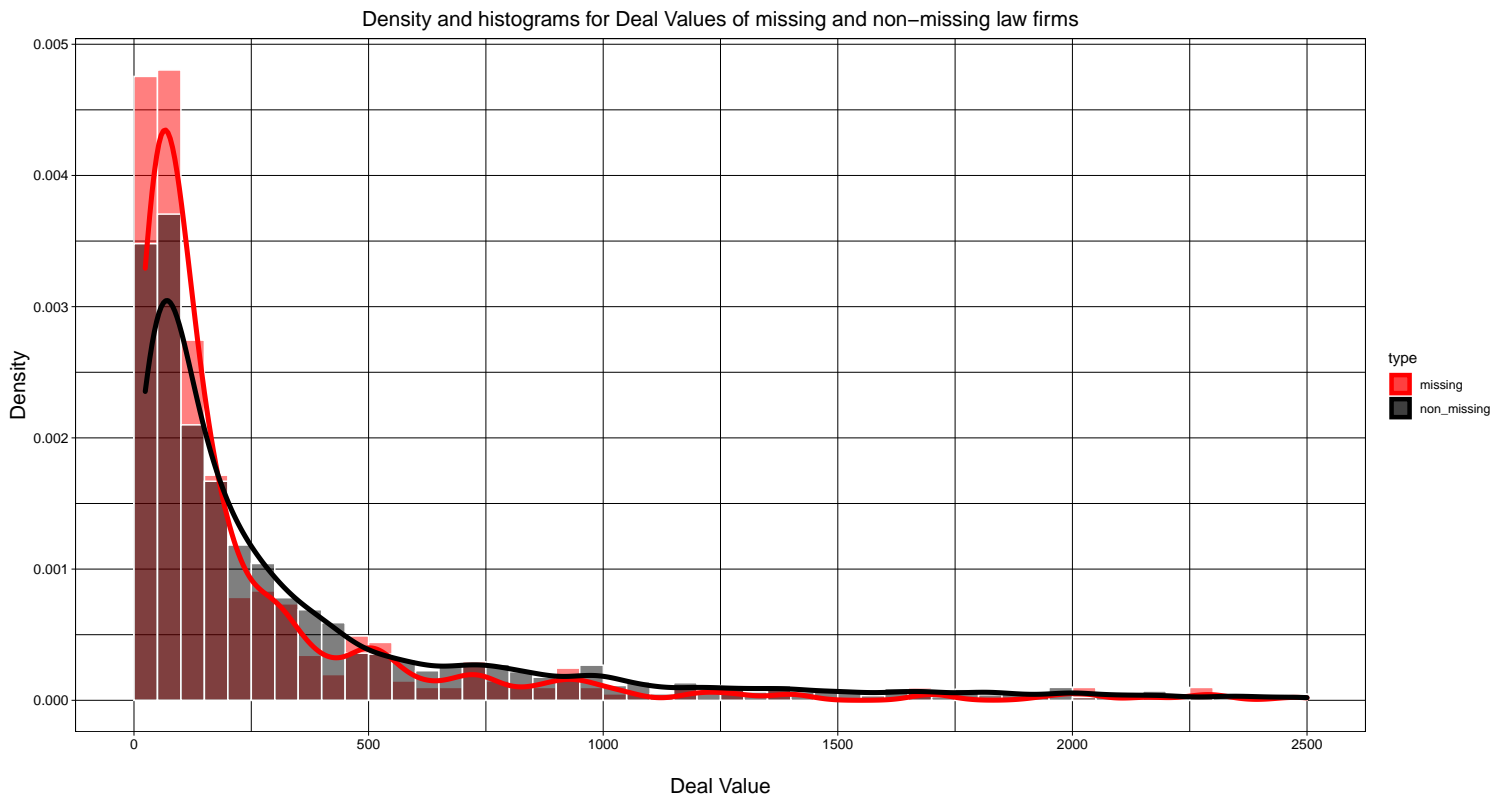
```

1 > missing_law_firm <- missing_law_firm %>%
2 +   mutate(`:=`(type, "missing"))
3 > non_missing_law_firm <- non_missing_law_firm %>%
4 +   mutate(`:=`(type, "non_missing"))
5 > miss_non_miss_law_firm <- rbind(missing_law_firm, non_missing_law_firm)
6 > View(miss_non_miss_law_firm)

```

Combine two dataset marked by types

Plot it based on types



Mean, Median and Standard Deviation

```
1 > miss_non_miss_law_firm %>%
2 +   group_by(type) %>%
3 +   summarise(`:=`(min, min(value1, na.rm = T)), `:=`(mean, mean(value1, na.rm = T)),
4 +             `:=`(median, median(value1, na.rm = T)), `:=`(max, max(value1, na.rm = T)),
5 +             `:=`(sd, sd(value1, na.rm = T))) %>%
6 +   kbl(caption = "Summary Stats Table comparing missing and non-missing", booktabs = T) %>%
7 +   kable_styling(latex_options = c("striped", "hold_position"))
```

Table 3: Summary Stats Table comparing missing and non-missing

type	min	mean	median	max	sd
missing	25	332.3452	107	9730	882.8552
non_missing	24	590.8620	185	39000	1694.6415

How many deal attorneys are missing biographical information?

What is the distribution of the missing data – are these mostly from earlier years, for example?

Load the no-matched data set

```
1 > com_deal_lawyer_firm_last_no_match <- read_csv("../data/confidence_match/com_deal_lawyer_firm_last_no_match.csv")
2 > # View(com_deal_lawyer_firm_last_no_match)
```

attorneys without biographical info

```
1 > att_miss_bio <- com_deal_lawyer_firm_last_no_match %>%
2 +   drop_na(First) %>%
3 +   mutate(`:=`(year, str_match(string = `Signing date`, pattern = "\\d\\d\\d\\d\\d\\d")))
4 > # View(att_miss_bio)
```

```

5 >
6 > att_miss_bio %>%
7 +   distinct(First, Last) %>%
8 +   nrow

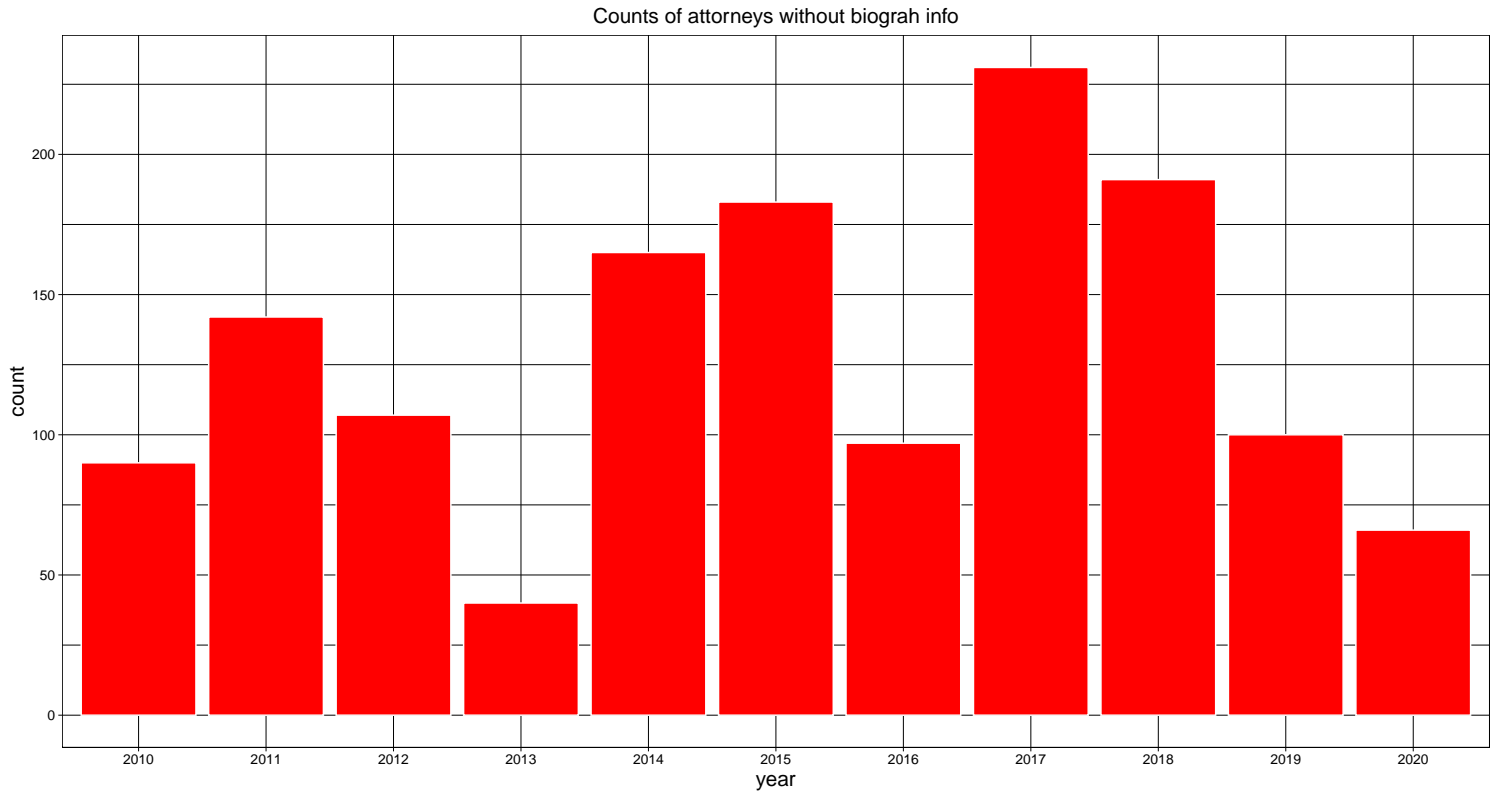
```

```
## [1] 904
```

There are 904 attorneys without biographical info.

Distribution of Years

```
## Warning: Ignoring unknown parameters: binwidth, bins, pad
```



Distribution of deal type

```
## Warning: Ignoring unknown parameters: binwidth, bins, pad
```

Distribution of attorneys with biograph info

```

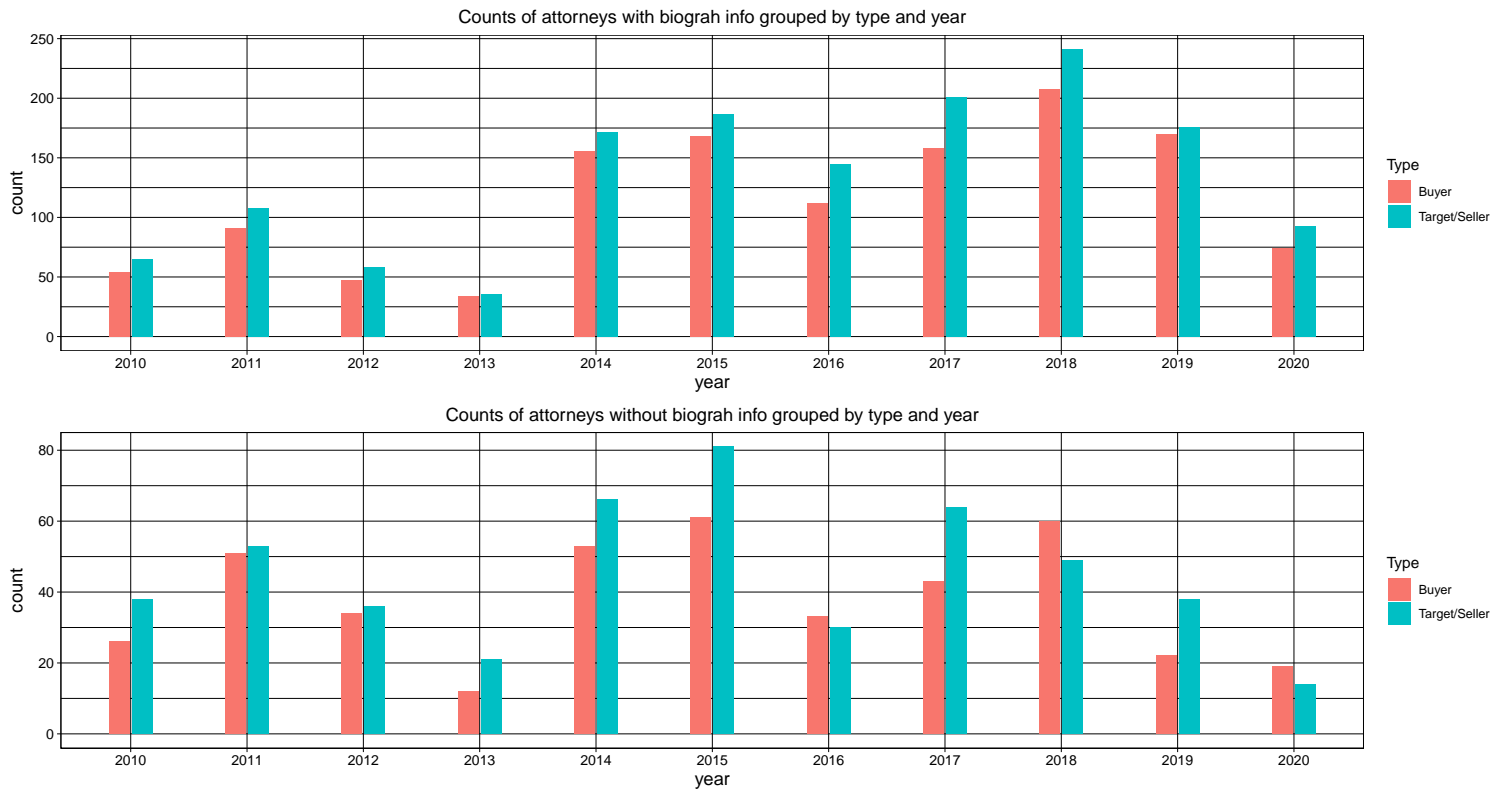
1 > three_matched_stacked <- read_csv("../data/confidence_match/three_matched_stacked.csv",
2 +   col_types = cols(.default = "c"))
3 > # View(three_matched_stacked)
4 >
5 > # If duplicates, keep attorneys from MA only
6 > dis_three_matched_stacked <- three_matched_stacked %>%
7 +   distinct(Deal_number, `Signing date`, First, Last, Law_Firm, .keep_all = T)
8 > # View(dis_three_matched_stacked)

```

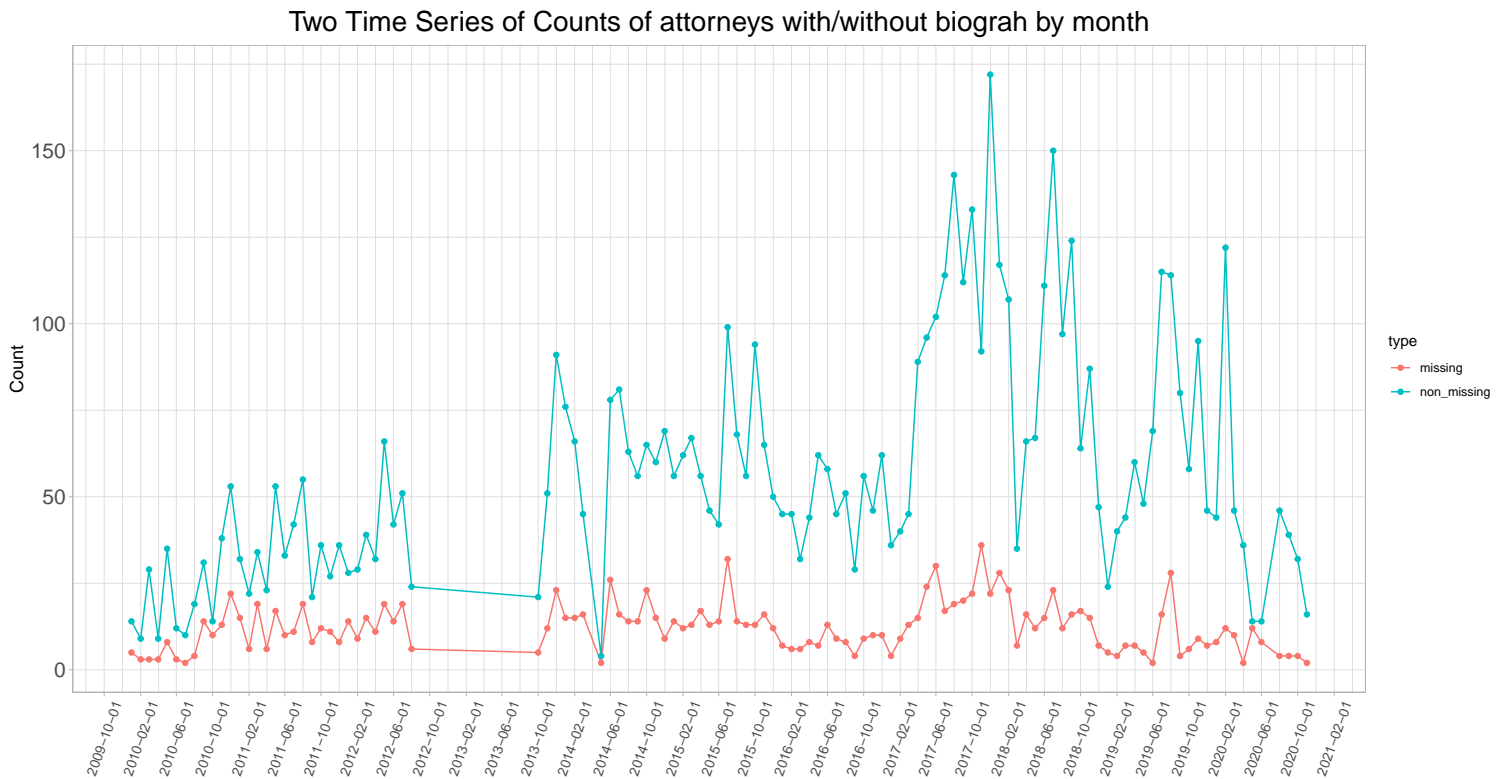
```
## Warning: Ignoring unknown parameters: binwidth, bins, pad
```

Compare them together

```
1 > p_with/p_without
```



Two time series by months of counts of attorneys with and without biograph info



From these two time series by month, we can tell numbers of attorneys with/without biograph info have very close trend but absolute number. The missing number is less than non-missing number, but they have the same trend.

How many deals are missing attorneys' biographical information?

```
1 > att_miss_bio %>%
2 +   distinct(Deal_number) %>%
3 +   nrow()
```

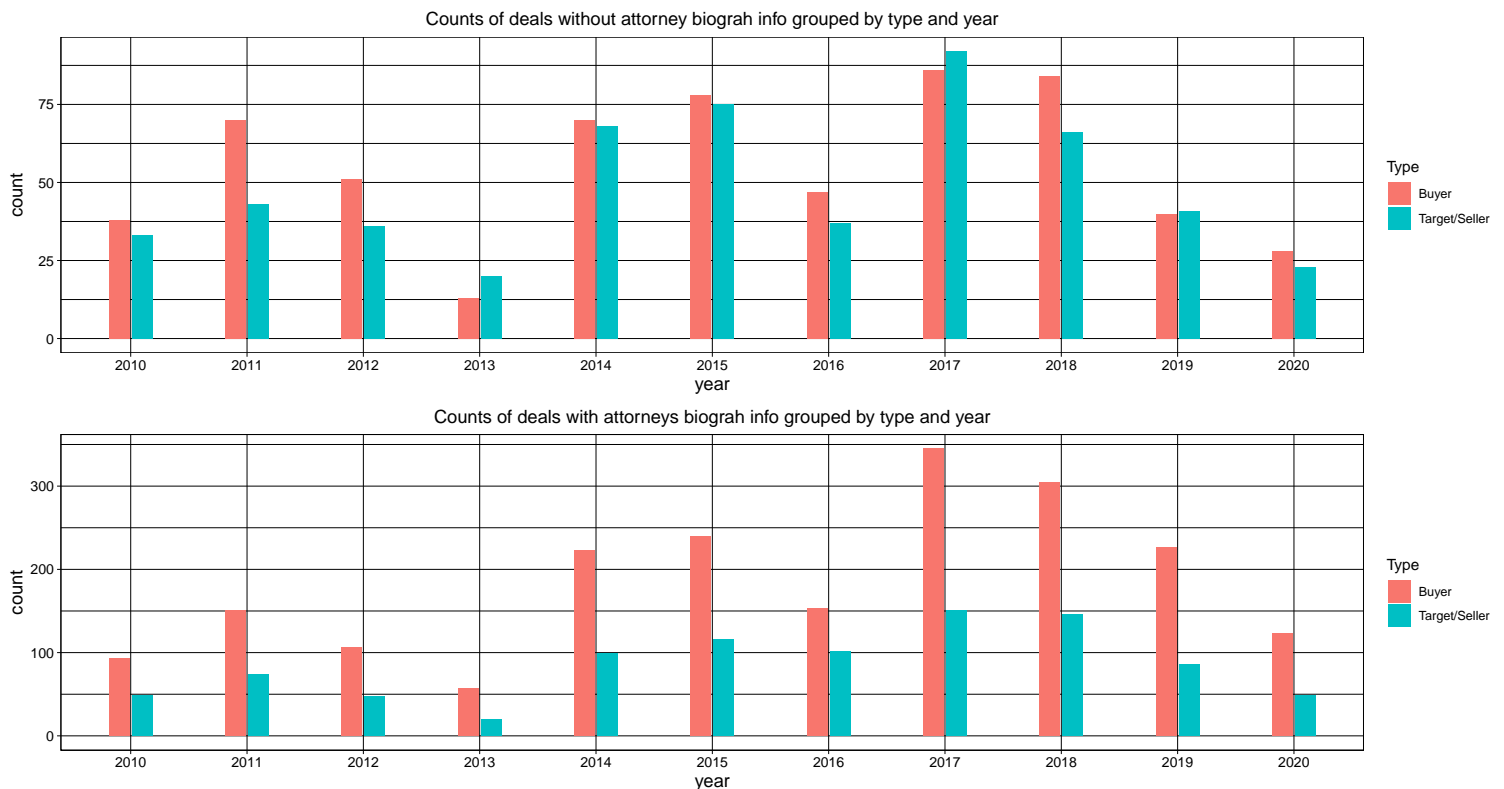
```
## [1] 1139
```

There are 1139 deals without attorneys' biograph info.

Counts of deals without attorneys' biograph info grouped by Year and Deal Type.

```
## Warning: Ignoring unknown parameters: binwidth, bins, pad
```

```
## Warning: Ignoring unknown parameters: binwidth, bins, pad
```



Create some summary tables

Create a summary table of the number of deals, gender breakdown, age distribution (of lawyers appearing on the deal); school distribution (i.e., most common law school attended), law firm distribution (i.e., most common law firms appearing on the deal); types of deals (e.g., health care; financial services, etc.); size of deal

Number of deals, gender, age

a summary table of the number of deals, gender breakdown, age distribution (of lawyers appearing on the deal)

We consider 25 years old as average ages when students graduate from law school.

```
1 > dis_three_matched_stacked %>%
2 +   drop_na(Gender, age_breaks) %>%
3 +   group_by(Gender, age_breaks) %>%
4 +   count() %>%
5 +   arrange(desc(n)) %>%
6 +   kbl(caption = "Summary Stats Table comparing missing and non-missing", booktabs = T) %>%
7 +   kable_styling(latex_options = c("striped", "hold_position"))
```


Table 4: Summary Stats Table comparing missing and non-missing

Gender	age_breaks	n
Male	45 < age <= 60	2713
Male	30 < age <= 45	1027
Male	age > 60	1011
Female	45 < age <= 60	295
Female	30 < age <= 45	213
Female	age > 60	68
Male	age <= 30	15
Female	age <= 30	3

school and law distribution

school distribution (i.e., most common law school attended), law firm distribution (i.e., most common law firms appearing on the deal)

Please see attachment csv file for full list data here.

```

1 > dis_three_matched_stacked %>%
2 +   drop_na(Gender, `Law School`) %>%
3 +   group_by(`Law School`, Law_Firm) %>%
4 +   count() %>%
5 +   arrange(desc(n)) %>%
6 +   write_csv(file = "../data/deal_lawyer/distri_law_sch_firm.csv")

```

```

1 > dis_three_matched_stacked %>%
2 +   drop_na(Gender, `Law School`) %>%
3 +   group_by(Gender, `Law School`) %>%
4 +   count() %>%
5 +   arrange(desc(n)) %>%
6 +   ungroup(Gender, `Law School`) %>%
7 +   top_n(50) %>%
8 +   ggplot(aes(y = reorder(`Law School`, n), x = n)) + geom_col(position = position_dodge2(),
9 +   width = 0.4, orientation = "y", aes(fill = Gender)) + labs(title = "Counts of lawyers grouped by gender and law school",
10 +   y = "Counts") + theme_linedraw() + theme(axis.text.x = element_text(size = 10),
11 +   axis.text.y = element_text(size = 10), axis.title = element_text(size = 15),
12 +   plot.title = element_text(hjust = 0.5, size = 15))

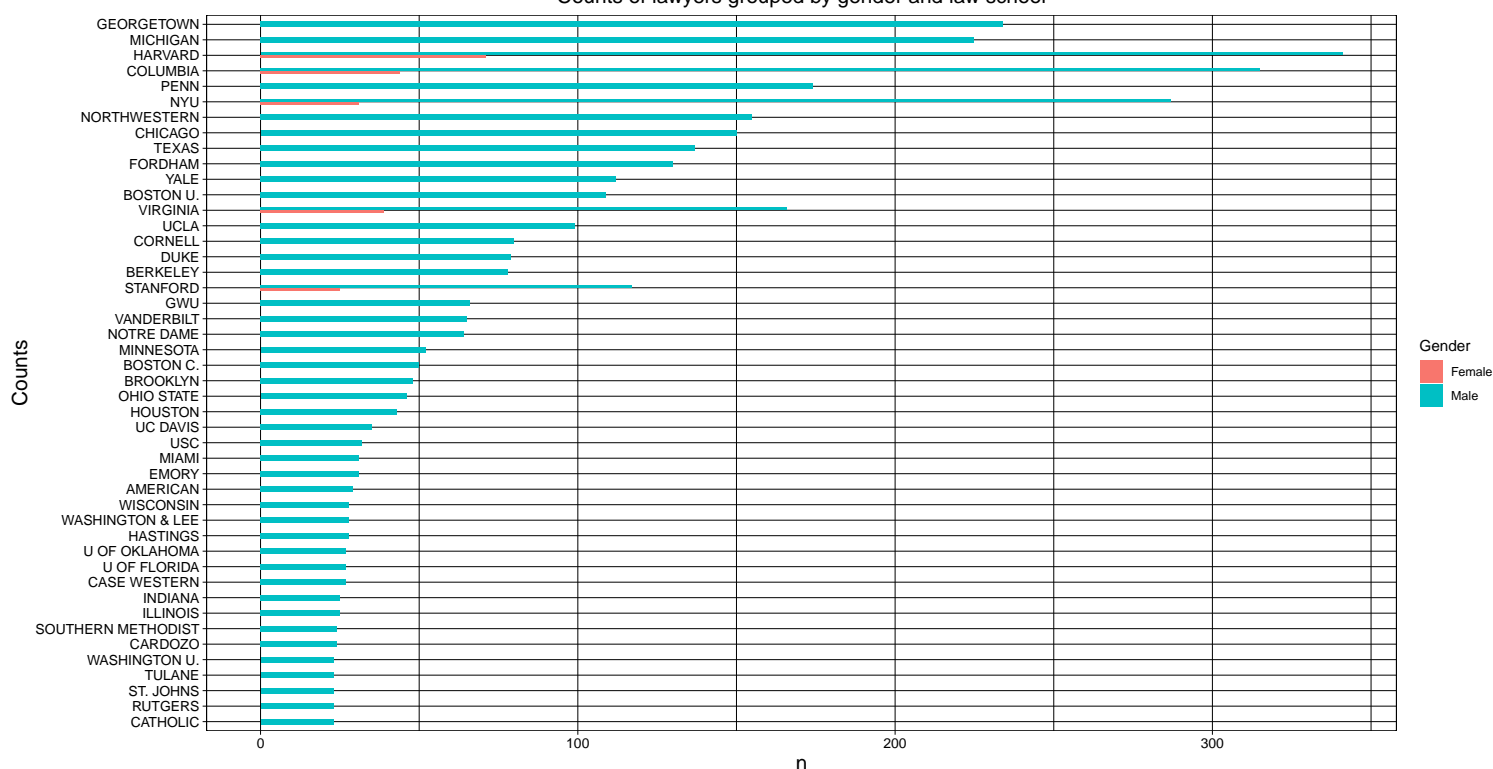
```

Selecting by n

Table 5: Top 30 most law school and law firm attended

Law School	Law_Firm	n
MICHIGAN	Kirkland & Ellis LLP	55
NORTHWESTERN	Kirkland & Ellis LLP	53
HARVARD	Latham & Watkins LLP	42
GEORGETOWN	Skadden, Arps, Slate, Meagher & Flom LLP	39
COLUMBIA	Wachtell, Lipton, Rosen & Katz	35
YALE	Wachtell, Lipton, Rosen & Katz	34
CHICAGO	Kirkland & Ellis LLP	32
NYU	Latham & Watkins LLP	32
FORDHAM	Skadden, Arps, Slate, Meagher & Flom LLP	30
COLUMBIA	Paul, Weiss, Rifkind, Wharton & Garrison LLP	29
HARVARD	Wachtell, Lipton, Rosen & Katz	28
COLUMBIA	Latham & Watkins LLP	27
NYU	Skadden, Arps, Slate, Meagher & Flom LLP	26
TEXAS	Latham & Watkins LLP	25
HARVARD	Sidley Austin LLP	24
NYU	Wachtell, Lipton, Rosen & Katz	24
PENN	Wachtell, Lipton, Rosen & Katz	24
COLUMBIA	Sullivan & Cromwell LLP	23
NYU	Willkie Farr & Gallagher LLP	23
UCLA	Jones Day	22
COLUMBIA	Kirkland & Ellis LLP	21
U OF OKLAHOMA	Alston & Bird LLP	21
HARVARD	Kirkland & Ellis LLP	20
MICHIGAN	Latham & Watkins LLP	20
BOSTON U.	Ropes & Gray LLP	19
HARVARD	Ropes & Gray LLP	19
NYU	Kirkland & Ellis LLP	19
TEXAS	Vinson & Elkins LLP	19
BOSTON U.	Goodwin Procter LLP	18
MINNESOTA	Faegre Baker Daniels LLP	18

Counts of lawyers grouped by gender and law school



types of deals

types of deals (e.g., health care; financial services, etc.); size of deal

```
1 > type_value_deal %>%
2 +   group_by(`Industry sector`) %>%
3 +   count(sort = T) %>%
4 +   ungroup(`Industry sector`) %>%
5 +   filter(n >= 10) %>%
6 +   kbl(caption = "Industry Sector Count", booktabs = T) %>%
7 +   kable_styling(latex_options = c("striped", "hold_position"))
```

Table 6: Industry Sector Count

Industry sector	n
Services	1199
Banking and financial services	827
Medical devices and healthcare	531
Oil and gas	498
Manufacturing and machinery	455
Pharmaceuticals and biotechnology	389
Computer and electronic equipment	354
Food and beverage	282
Construction and materials	223
Consumer goods	203
Media and entertainment	181
Insurance	175
Travel and leisure	167
Utilities	135
Automobiles, airlines and transportation	121
Telecommunications	115
Retailers	108
Chemicals	99
Computer and electronic equipment Services	99
Mining and metals	78
Aerospace and defense	72
Textiles and apparel	60
Forestry and paper	34
Real estate	33
Agriculture	31
Banking and financial services Services	24
Medical devices and healthcare Services	21
Manufacturing and machinery Construction and materials	10
Services Utilities	10

```
1 > type_value_deal %>%
2 +   drop_na(value1) %>%
3 +   mutate(`:=`(value_breaks, case_when(value1 <= 250 ~ "value <= 250", value1 >
4 +     250 & value1 <= 500 ~ "250 < value <= 500", value1 > 500 & value1 <= 1500 ~
5 +     "500 < value <= 1500", value1 > 1500 ~ "1500 < value")) %>%
6 +   group_by(Type, value_breaks) %>%
7 +   count(sort = T) %>%
8 +   ungroup() %>%
9 +   kbl(caption = "Counts of deals grouped by types and value\\_breaks", booktabs = T) %>%
10 +   kable_styling(latex_options = c("striped", "hold_position"))
```

Table 7: Counts of deals grouped by types and value_breaks

Type	value_breaks	n
Target/Seller	value <= 250	1258
Buyer	value <= 250	1247
Target/Seller	250 < value <= 500	1140
Buyer	250 < value <= 500	1128
Target/Seller	1500 < value	975
Buyer	1500 < value	907

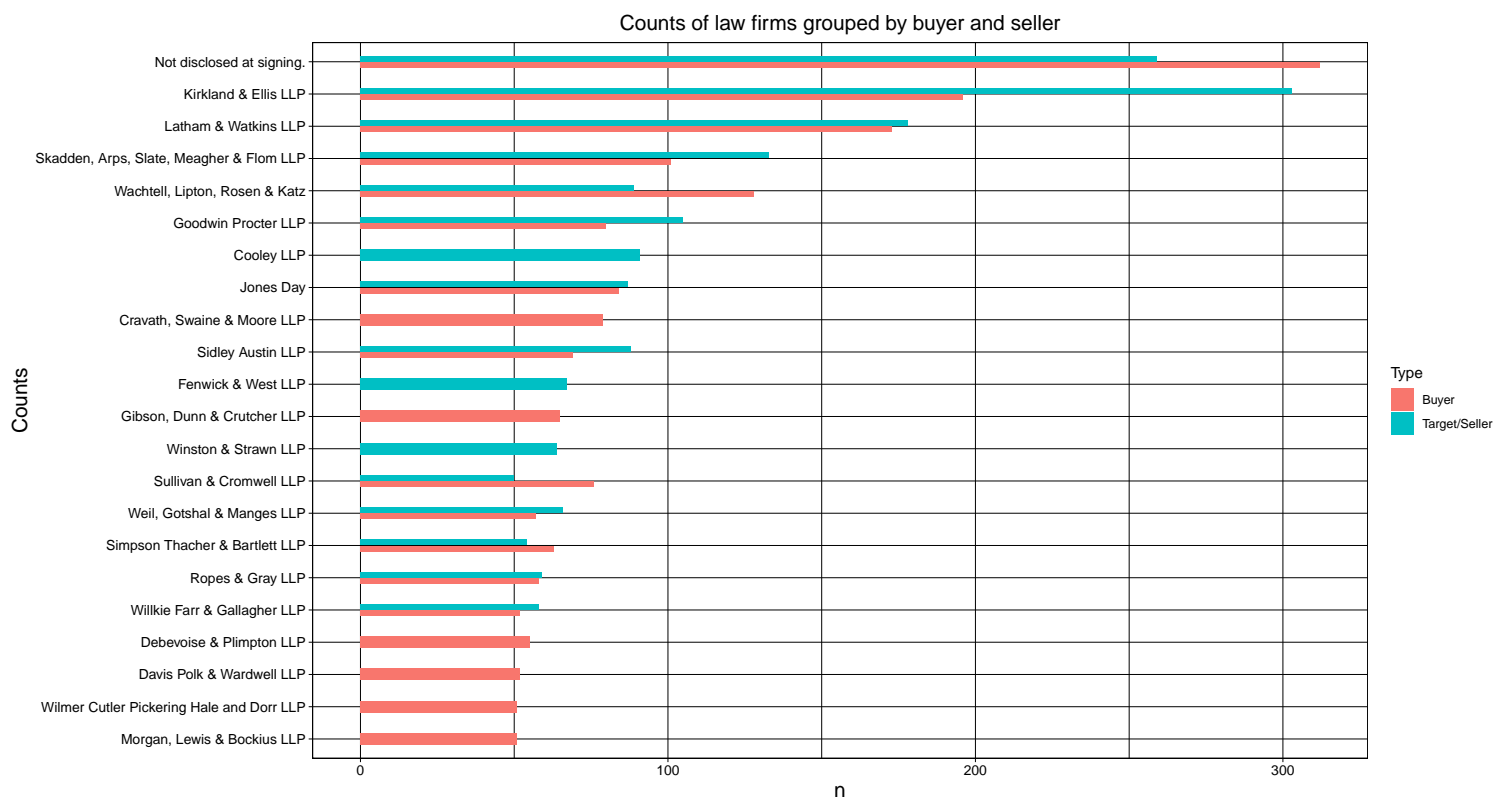
Buyer and seller of deals

Breakdown by buyer and seller in the deals. We think that law firms and M&A lawyers appear on both sides of the deal with roughly the same probability, but it would be helpful to know if this were true.

```

1 > encoded_merge_dl %>%
2 +   group_by(Law_Firm, Type) %>%
3 +   count(sort = T) %>%
4 +   ungroup(Law_Firm, Type) %>%
5 +   filter(n >= 50) %>%
6 +   ggplot(aes(y = reorder(Law_Firm, n), x = n)) + geom_col(position = position_dodge2(),
7 +   width = 0.4, orientation = "y", aes(fill = Type)) + labs(title = "Counts of law firms grouped by buyer and
8 +   y = "Counts") + theme_linedraw() + theme(axis.text.x = element_text(size = 10),
9 +   axis.text.y = element_text(size = 10), axis.title = element_text(size = 15),
10 +   plot.title = element_text(hjust = 0.5, size = 15))

```



```

1 > dis_three_matched_stacked %>%
2 +   filter(Source == "M&A") %>%
3 +   count(Type, sort = T) %>%
4 +   kbl(caption = "MA Lawyers grouped by Buyer/Seller", booktabs = T) %>%
5 +   kable_styling(latex_options = c("striped", "hold_position"))

```

Table 8: MA Lawyers grouped by Buyer/Seller

Type	n
Target/Seller	2685
Buyer	2660

Regression Analysis

Some basic regressions: e.g., regress gender on observable characteristics (e.g., firm, industry, size of deal). Explain what factors explain when women are more likely to appear on a deal.

First we used Multiple Logistic Regression with all variables

age coefficients

Law_school coefficient signs

Law_Firm coefficient signs

Industry Sector coefficient signs

Multiple Logistic Regression with all variables except for Law Schol and Law firm

Random Forest with the whole variables used

Variable Importance after fitting our model

Conclusion

From two method, Logistic Classification and Random Forest, we can tell that Age, Law_Firm, Law_school are obvious significant when predicting Gender. From huge positive coefficients bar plots, at least we can tell that Law_SchoolMeMPhis can make more likely Women on the deals. And age is another huge negative affects which can make huge Adverse conditions for women appearing on Deals. Deal Value and Type(buyer/Seller) aren't significant variables affecting women.

As for other law_school and law_firm coefficients, there are both positive and negative coefficients affecting women showing on Deals.

1. In total, there are 128 Law_school coefficients greater than 0, while 18 Law_school lower than 0, which means Law_school in general make positive affects on women showing on Deals.
2. However, there are 174 Law_Firm coefficients greater than 0, while 296 Law_Firm lower than 0, which means Law_Firm in general make negative affects on women showing on Deals.
3. There are 17 Industrial coefficients greater than 0, while 30 Industrial sector lower than 0, which means Industrial Sector in general make negative affects on women showing on Deals.