

Lab1

Melwin Poovakottu

Friday, January 20, 2017

This is an R Markdown document. Markdown is a simple formatting syntax for authoring HTML, PDF, and MS Word documents. For more details on using R Markdown see <http://rmarkdown.rstudio.com>.

When you click the **Knit** button a document will be generated that includes both content as well as the output of any embedded R code chunks within the document. You can embed an R code chunk like this:

```
library(ggplot2)
library(dplyr)
```

```
##
## 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

load("ceo_w203.RData")
```

The market value has continuous variable. While running a summary command you see there is big difference in the mean(=3450) and Median(=1200). This means that the data is positively skewed to the right with a lot of data points close to the median and a relatively high number of outliers with higher market value.

We see that there are some negative market value data points. These data points also have the profits as -1. This could potentially be instances of missing data points and default value of -1 is recorded. So we will be considering them as NA.

```
summary(CEO$mktval)

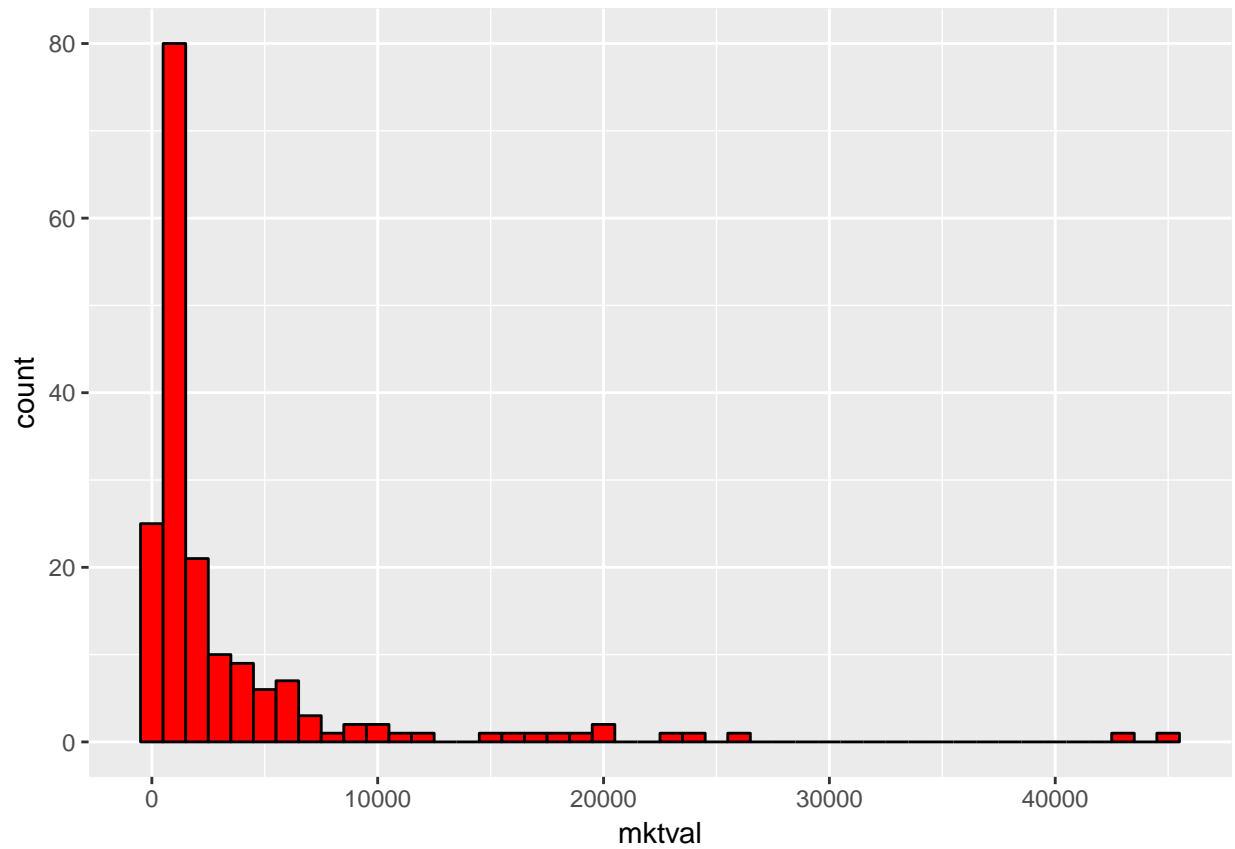
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      -1     567    1200    3450    3200   45400

CEO$mktval = ifelse(CEO$mktval == -1, NA, CEO$mktval)
```

The positive skew is clear in the histogram. Since there is a large variation in the values we can basically try to use the log of market value for our analysis.

```
ggplot(CEO,aes(mktval))+ geom_histogram(binwidth = 1000,fill ="red",col="black")

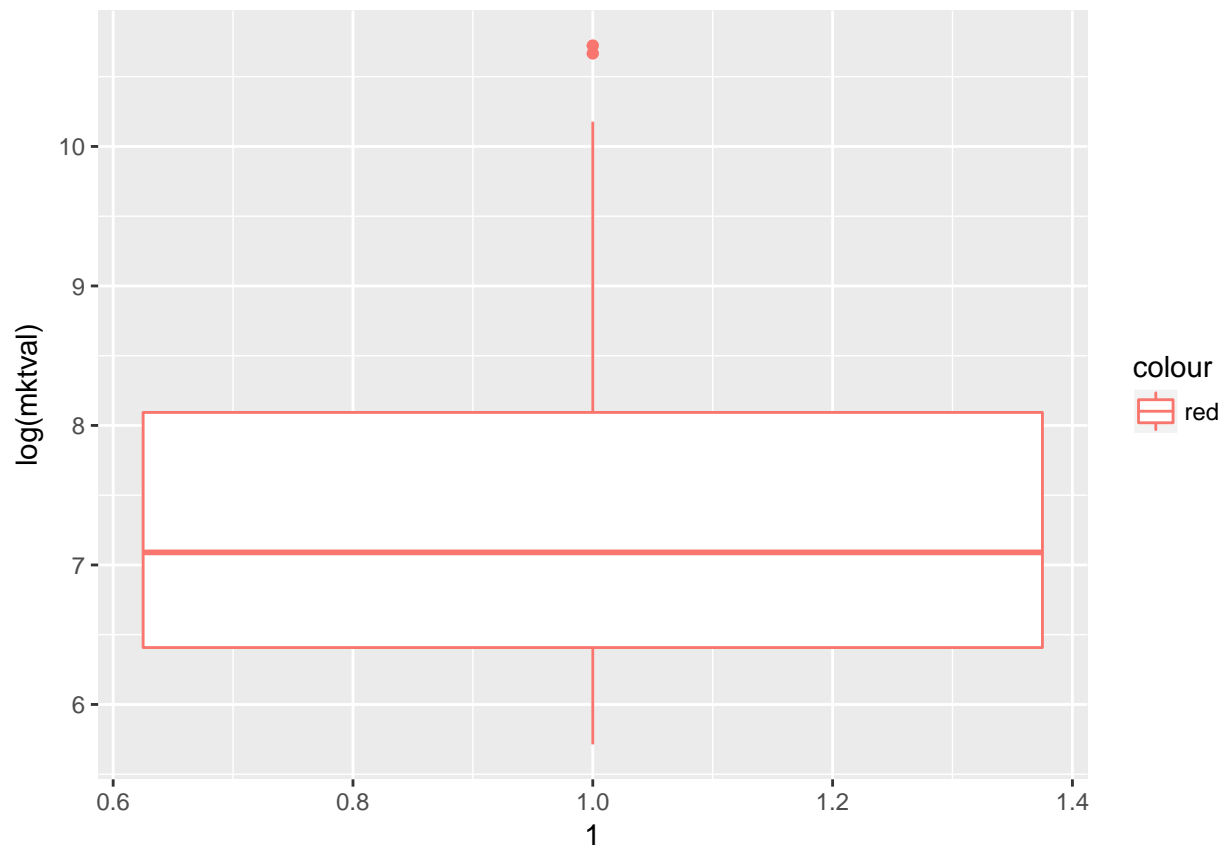
## Warning: Removed 5 rows containing non-finite values (stat_bin).
```



The boxplot for the log of market value shows a more uniform distribution of data points

```
ggplot(CEO,aes(y=log(mktval),x=1,color='red'))+ geom_boxplot()
```

```
## Warning: Removed 5 rows containing non-finite values (stat_boxplot).
```



We can see a high correlation between market value and profits, which is expected for most companies. Similarly we see a good correlation between market value and salary. Also we see that if the CEO has a graduate degree he is more likely to get a better salary than without a graduate degree. On plotting the salary and profits against mktval, we see there are few data points that are outliers which are affecting the regression line.

```
cor(CEO$mktval, CEO, use="complete.obs")
```

```
##          salary      age   college    grad   comten   ceoten
## [1,] 0.4082068 0.1239606 0.004211891 0.1139215 0.1462181 0.01317969
##          profits mktval
## [1,] 0.9184209      1
```

```
ggplot(CEO, aes(x = log(mktval), y = log(profits), color = as.factor(1))) + geom_point() + stat_smooth(m
```

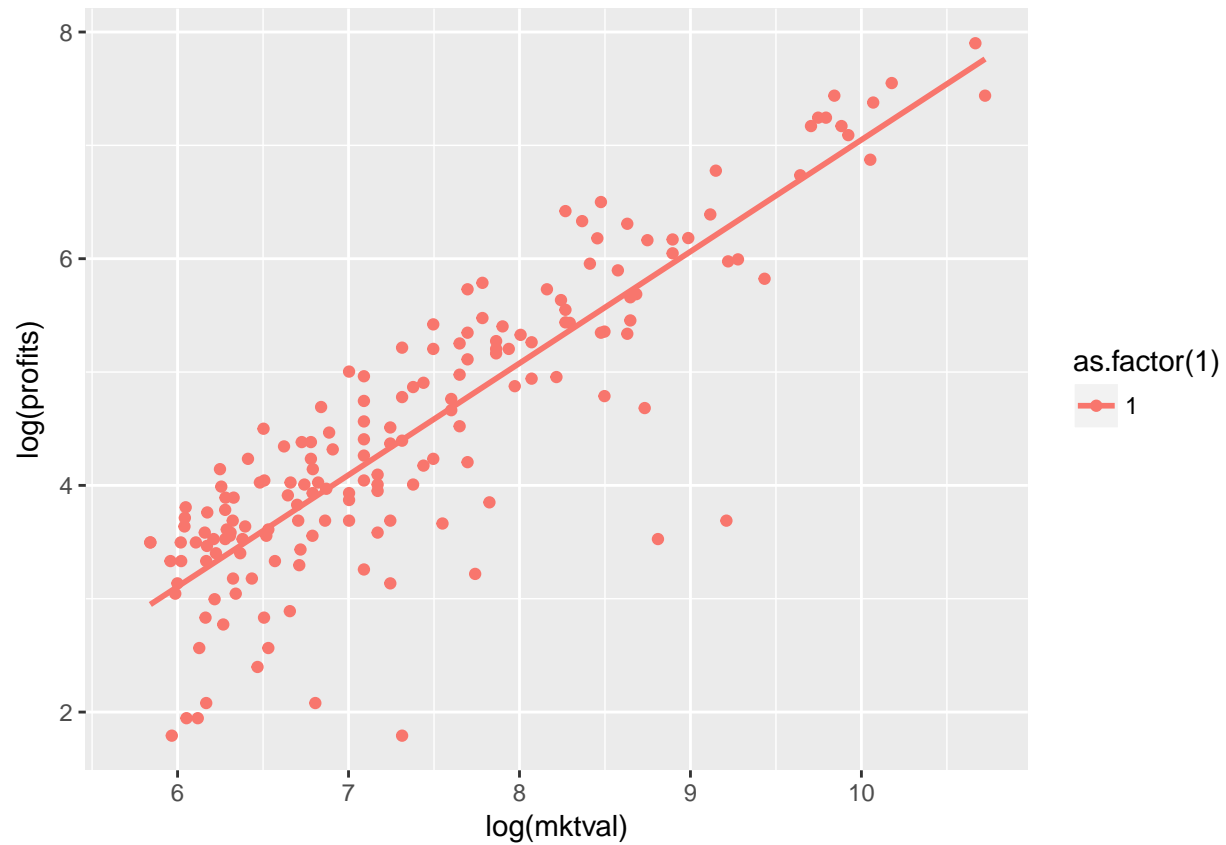
```
## Warning in log(profits): NaNs produced
```

```
## Warning in log(profits): NaNs produced
```

```
## Warning in log(profits): NaNs produced
```

```
## Warning: Removed 15 rows containing non-finite values (stat_smooth).
```

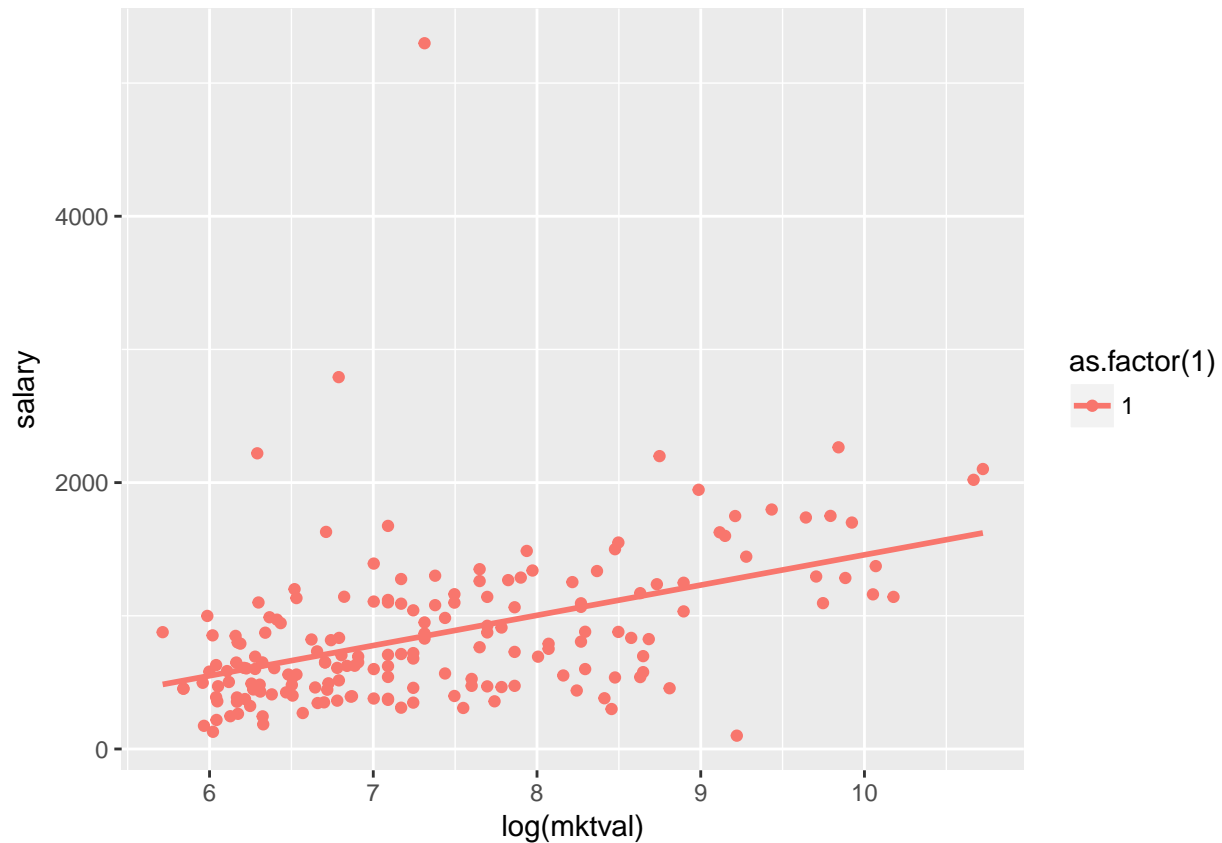
```
## Warning: Removed 15 rows containing missing values (geom_point).
```



```
ggplot(CEO, aes(x = log(mktval), y = salary, color = as.factor(1))) +geom_point() + stat_smooth(method = "lm", se = FALSE)
```

```
## Warning: Removed 5 rows containing non-finite values (stat_smooth).
```

```
## Warning: Removed 5 rows containing missing values (geom_point).
```



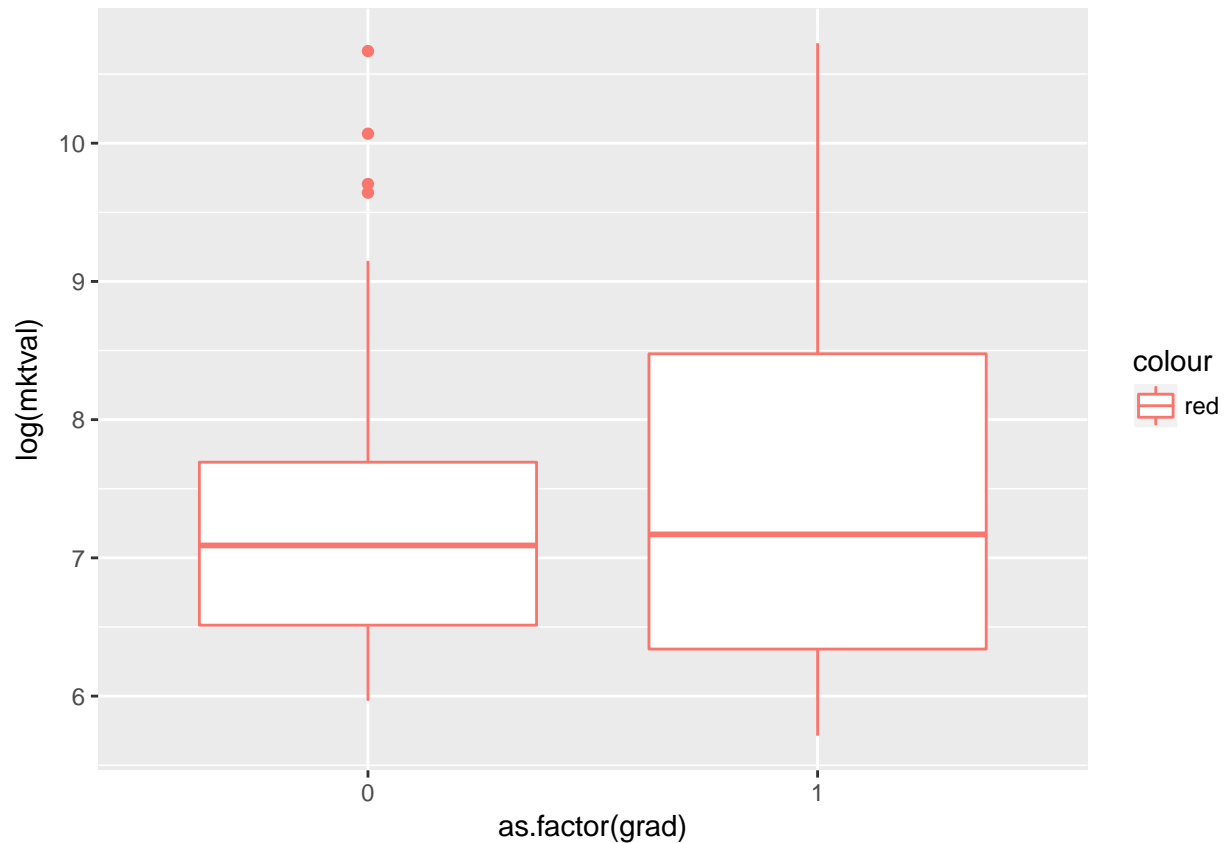
```
cor(log(CEO$mktval),CEO$salary,use = "complete.obs")
```

```
## [1] 0.4438906
```

From the boxplot for log of mktval vs grad, we can see that there is greater variation in the market price if you have a grad degree

```
ggplot(CEO,aes(y=log(mktval),x=as.factor(grad),color='red'))+ geom_boxplot()
```

```
## Warning: Removed 5 rows containing non-finite values (stat_boxplot).
```



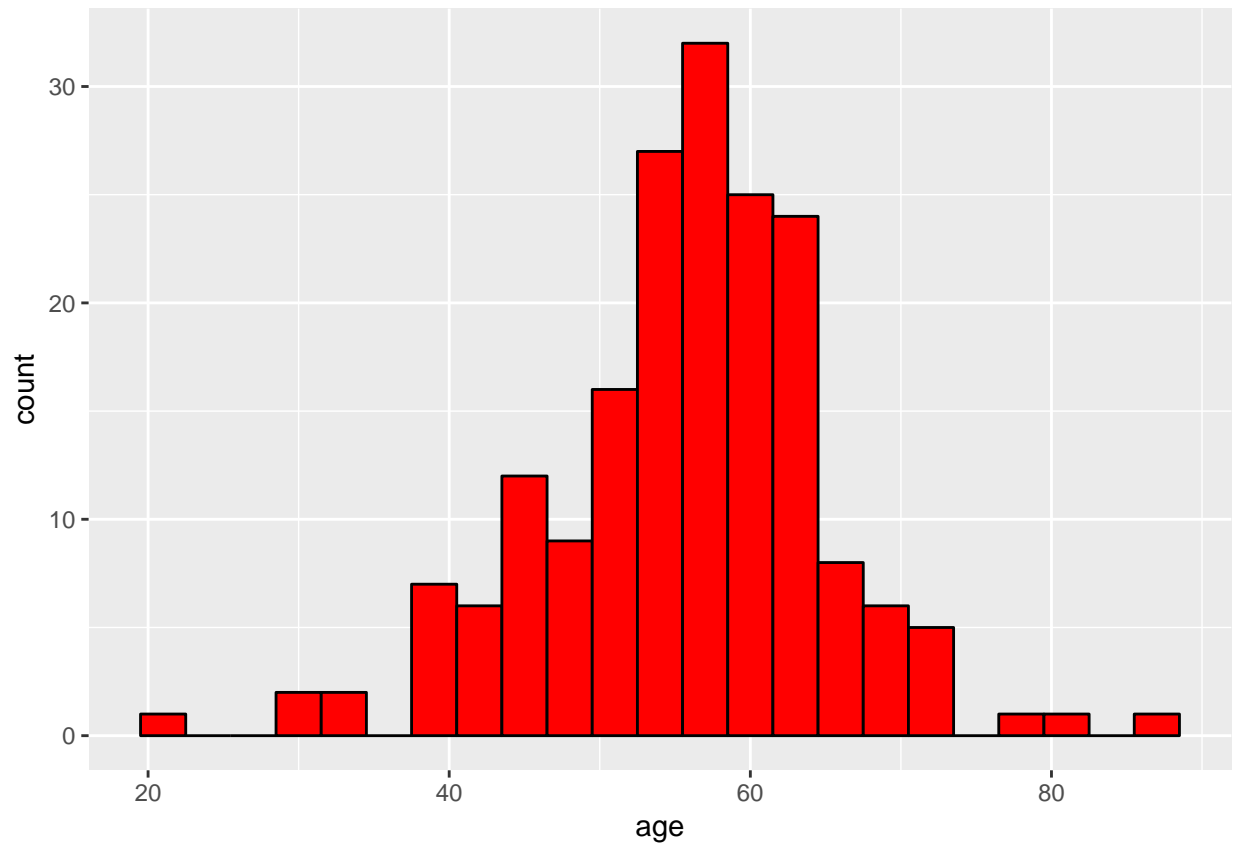
Age

Age is a discrete variable. Looking at the histogram plot, we can see that the variable is close to a normally distributed curve with the mean(55) and median(57). But you can see a drastic drop after the age of 65. This likely corresponds to the age of retirement in the US.

```
summary(CEO$age)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      21.00  51.00   57.00   55.78  61.00   86.00
```

```
ggplot(CEO,aes(age,color = as.factor(1)))+ geom_histogram(binwidth = 3,fill ="red",col="black")
```



If we simply plot a graph of age vs salary do not get any particular co-relation. The corelatoin is also very low(0.13) We do see that if you have a grad degree you are likely to get paid more as you in advance your career.

```
ggplot(CEO,aes(x=age,y=salary,color = as.factor(1)))+ geom_point()+ stat_smooth(method = "lm")
```

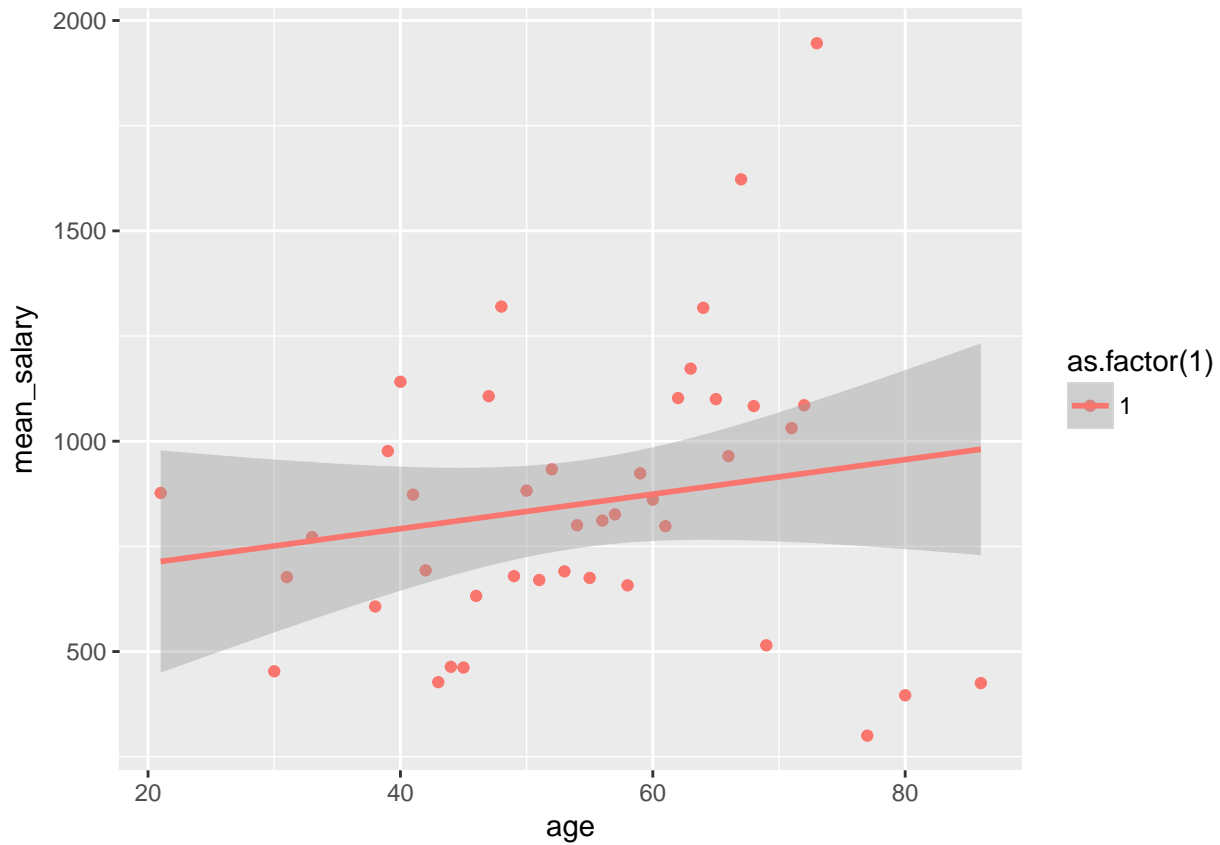


```
cor(CEO$age,CEO$salary)
```

```
## [1] 0.130081
```

Let us examine age variable closely. Since age is a discrete variable let us group the salary's of all the CEOs of the same age. We can take mean of all the salary's in each group. On plotting the mean salary against the age we see an interesting trend. We see that within the age range 29 to 75 there might be a correlation between average salary and age. This is also the range which has the maximum data points.

```
grouped = group_by(CEO, age)
mean_salary_age = summarise(grouped, mean_salary = mean(salary, na.rm = T))
ggplot(mean_salary_age, aes(x=age, y=mean_salary, color = as.factor(1))) + geom_point() + stat_smooth(method="lm")
```

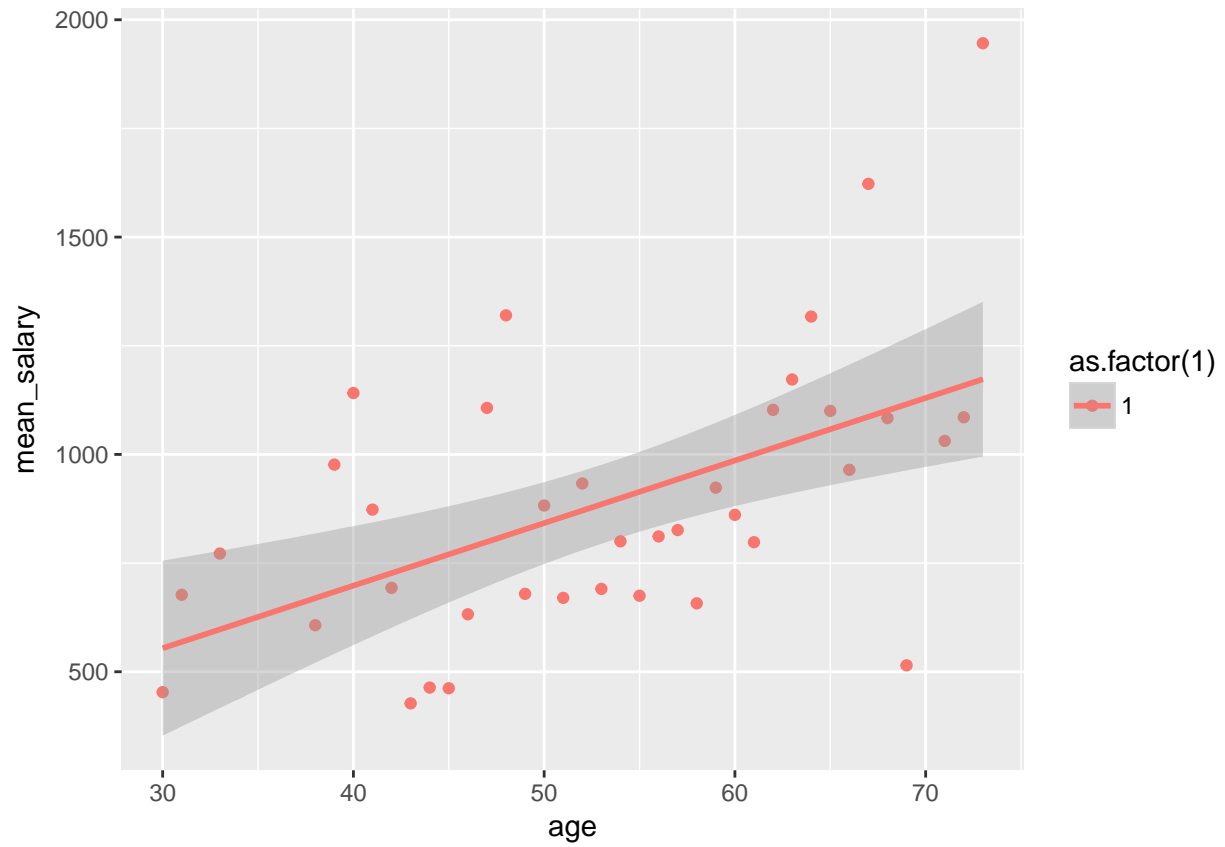



```
cor(mean_salary_age)
```

```
##           age mean_salary
## age      1.00000    0.17776
## mean_salary 0.17776    1.00000
```

Let Narrowing age from 29 to 75 We see a high co-relation between mean salary and the age within the age range 29 to 75

```
narrowed_mean_salary_age=mean_salary_age[mean_salary_age$age<=75 & mean_salary_age$age>=29,]
ggplot(narrowed_mean_salary_age,aes(x=age,y=mean_salary,color = as.factor(1)))+ geom_point() + stat_sm
```



```
cor(narrowed_mean_salary_age$age,narrowed_mean_salary_age$mean_salary)
```

```
## [1] 0.5321218
```