

# Final Project

*Oilivia Wagner | Weijia Xiong | Wurongyan Zhang | Yiling Yang*

## Contents

Abstract . . . . .	2
Introduction . . . . .	2
Exploratory data analysis . . . . .	2
Method . . . . .	7
<b>try adding rank</b>	<b>41</b>
Results . . . . .	51
Conclusions . . . . .	51
Discussion . . . . .	51
Figures and tables . . . . .	52
References . . . . .	54
Appendix . . . . .	55

## Abstract

## Introduction

In our project, we examined the whether gender discrimination existed in setting salaries for people in the academia or higher education institutions.

## Exploratory data analysis

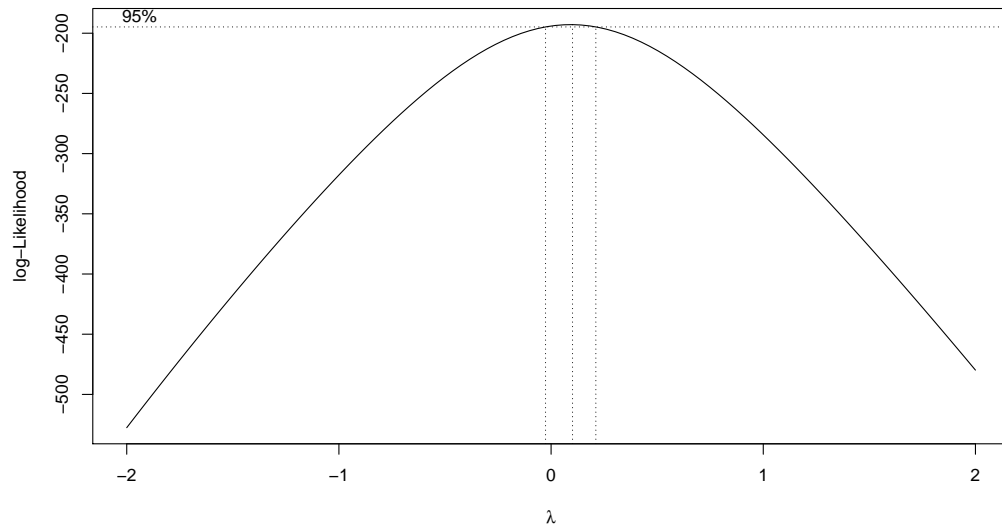
Data used in our studies were collected from 261 individuals who work in the academia or higher education institutions. The raw dataset contains following features.

- **Dept:** 1= Biochemistry/Molecular Biology 2= Physiology 3= Genetics 4= Pediatrics 5= Medicine 6= Surgery
- **Gender:** 1= Male, 0= Female
- **Clin:** 1= Primarily clinical emphasis, 0= Primarily research emphasis
- **Cert:** 1= Board certified, 0= not certified
- **Prate:** Publication rate (# publications on cv) / (# years between CV date and MD date)
- **Exper:** # years since obtaining MD
- **Rank:** 1= Assistant, 2= Associate, 3= Full professor (a proxy for productivity)
- **Sal94:** Salary in academic year 1994
- **Sal95:** Salary after increment to Sal94

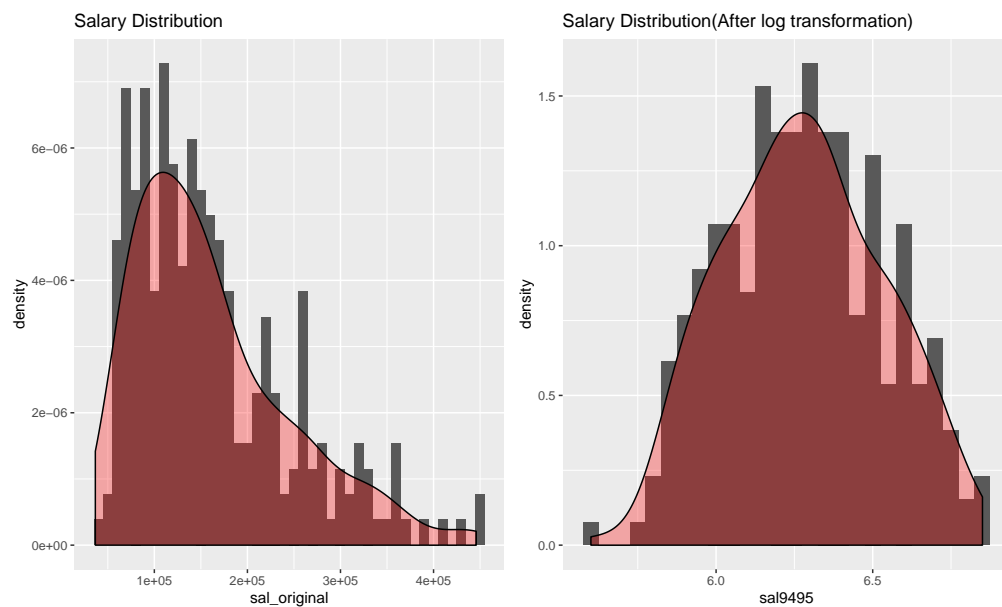
Table 1 contains the summary of variables in the dataset. Fortunately, there are no missing values in our dataset.

We then need to examine each interested variable against the main effect and main interest.

## Salary

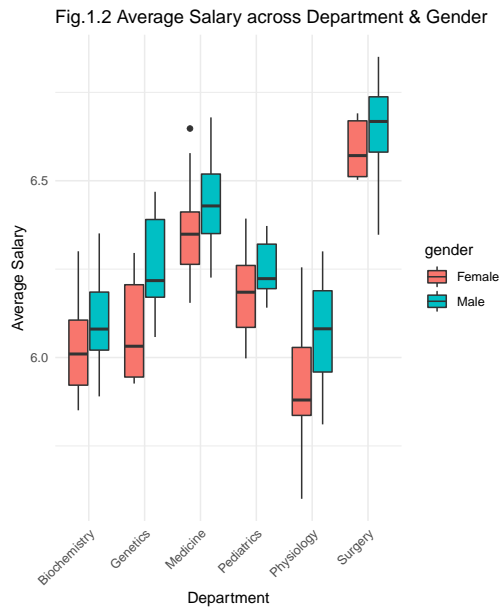
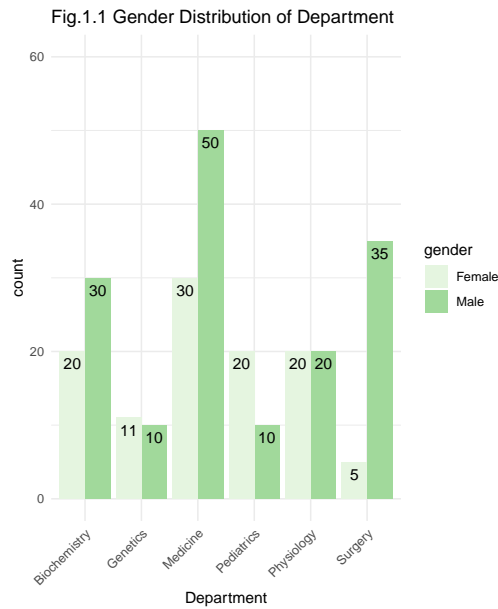


Since  $\lambda = 0$ , we use log transformation.



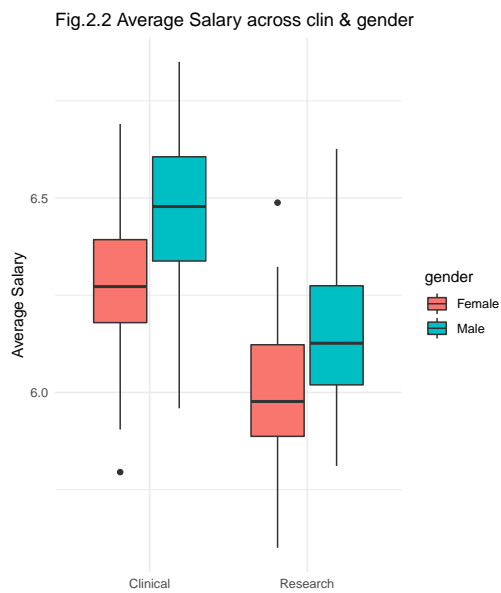
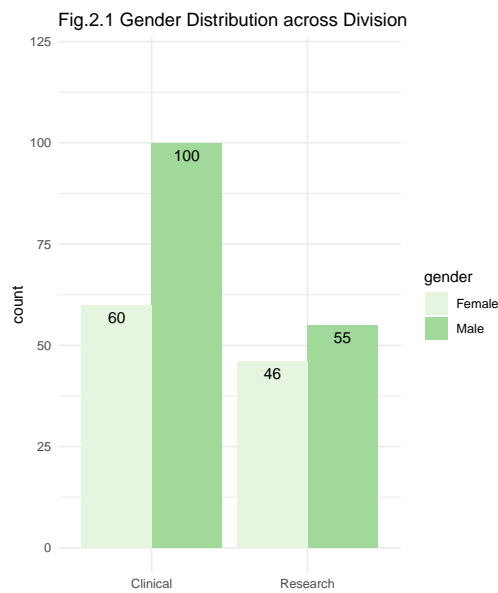
## Department

Fig 1.1. shows that the gender ratios in the department of Genetics and Physiology are very balanced, compared with department of Medicine and Surgery, which are very imbalanced. The differences in department of Biochemistry and Pediatrics are moderate. Fig 1.2 includes our main interest the salary into account. It seems that across all department, male earn more than female do. However, before further analysis, we cannot tell whether those difference are significant.



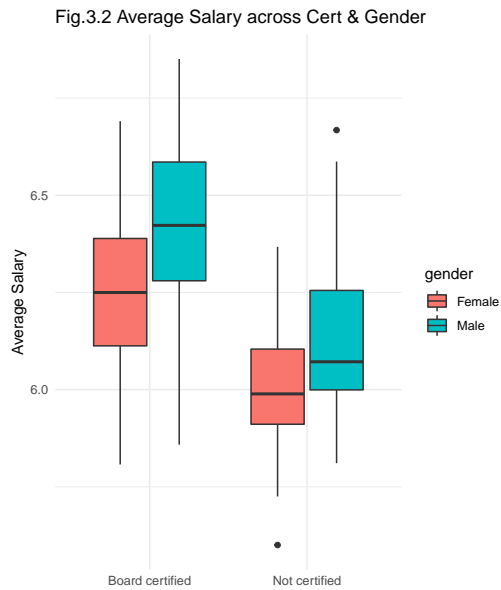
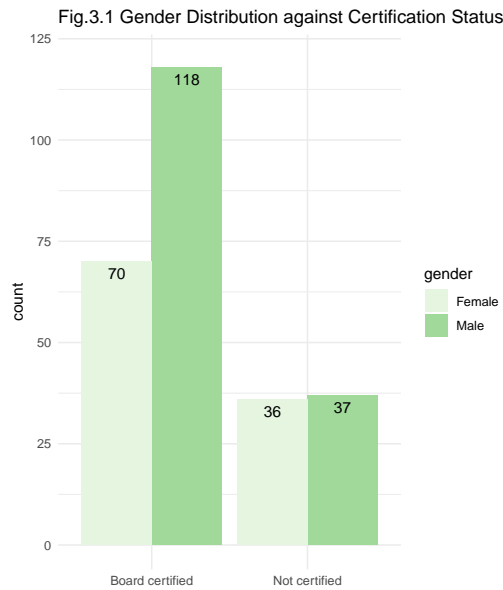
## Clinical or Research Division

Fig 2.1 shows that in either division, there are more male than female, especially the clinical division. Regarding salary, again, male in either division earned more than female do.



## Certification status

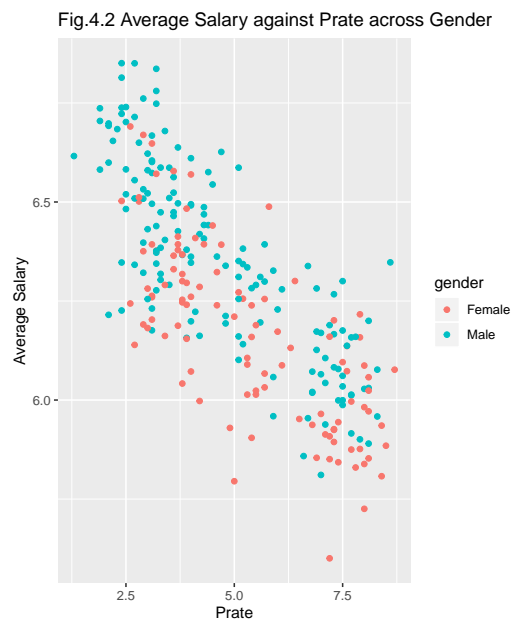
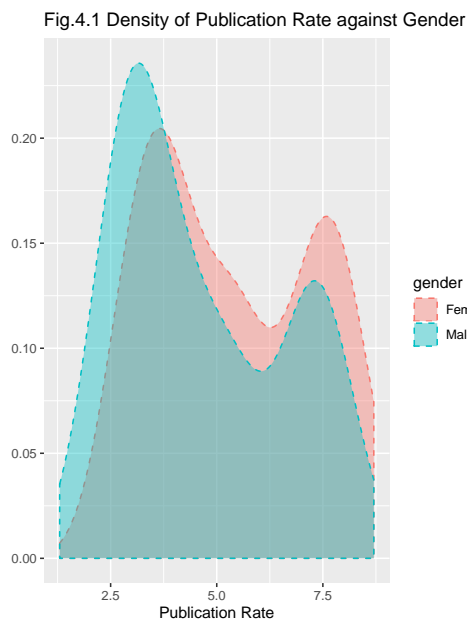
Fig.3. shows that the amount of certified male outnumbers surely the amount of certified female. However, for those without certification, those are just even. For the salary, male again earn more than female do regardless of his or her certification status.



## Prate

Fig 4.1 shows the density plot of publication rate for both male and female. No obvious difference was observed.

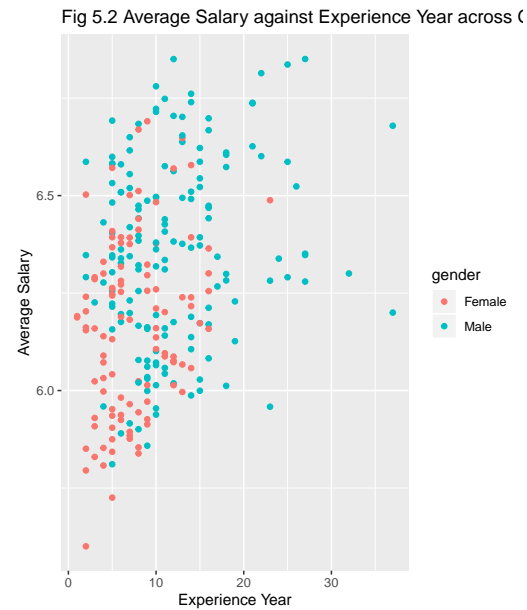
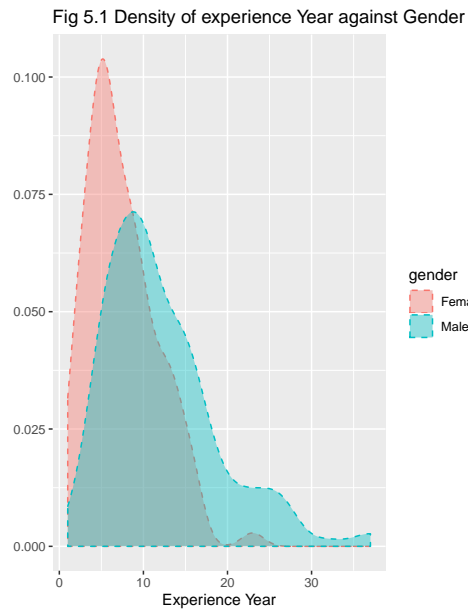
Fig 4.2 implies that there might be linear trend between average salary and publication rate. However, it's hard to tell whether there is any difference on the effect on the salary regarding the gender.



## Exper

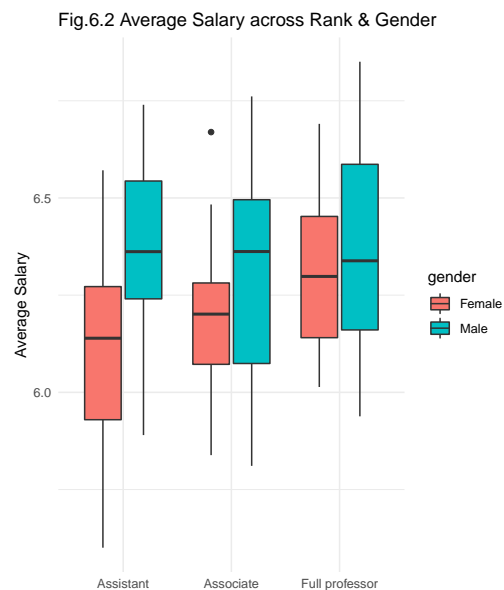
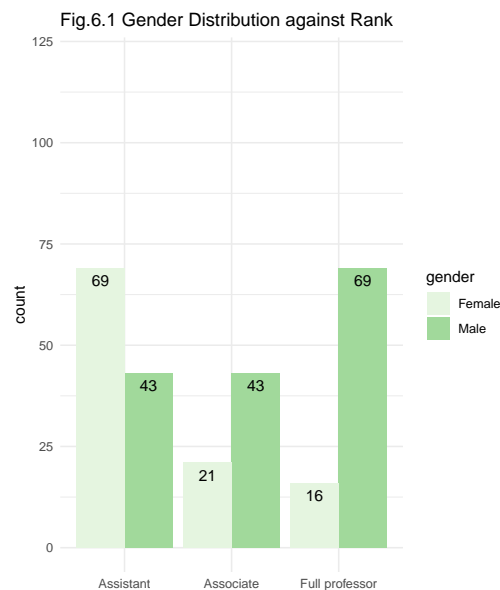
Fig.5.1 shows that the densities of experience year for both male and female are both very skewed, and seems that female have a more skewed trend. Although there are not obvious linear trend between salary and experience

year, the salary for male spread more widely than those for female and individual with high salary are dominately male.

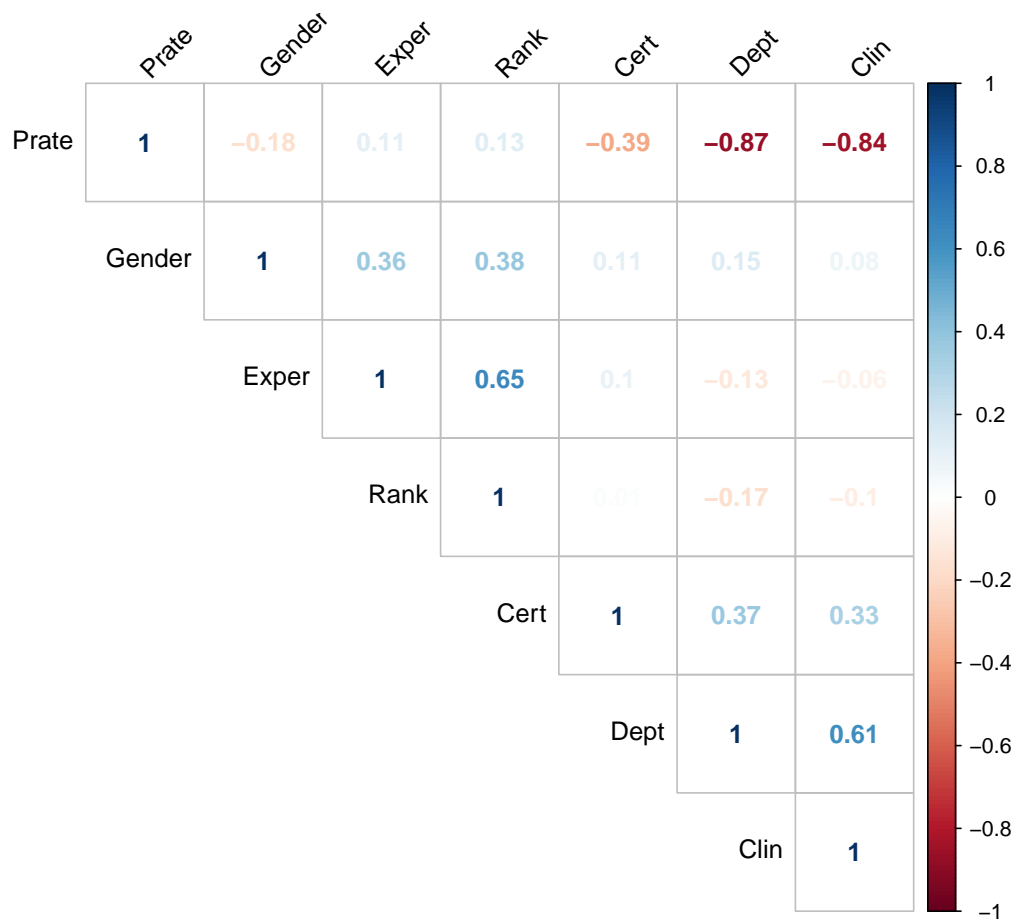


## Rank

Fig 6.1 shows that there are many female assistant professor and less associate and full professor than male. Without surprising, in either rank, male earn more than female.



Lastly, after examining the correlation matrix, we can see that there are some highly related variables ( $r > 80\%$ ). They are 1. Department (Dept) and Publish Rate (Prate) and 2. Clin and Prate. Some have morderate correlation such as 1. Experience Year (Exper) annd Cert (Certification status) and 2.Department (Dept) and Clin.



Those imply potential collinearities. As we go through those variables, there are some outliers. In later section, we will examine further about them.

## Method

```
##
## Call:
## lm(formula = sal9495 ~ rank, data = lawsuit)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.62108 -0.19816 -0.00634  0.17243  0.51821
##
## Coefficients:
```

```

##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)      6.22148    0.02348 265.019 < 2e-16 ***
## rankAssociate     0.04770    0.03893   1.225  0.22155
## rankFull professor 0.13634    0.03574   3.815  0.00017 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.2484 on 258 degrees of freedom
## Multiple R-squared:  0.05374,    Adjusted R-squared:  0.0464
## F-statistic: 7.326 on 2 and 258 DF,  p-value: 0.0008047
##
## Call:
## lm(formula = sal9495 ~ gender, data = lawsuit)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.5628 -0.1829 -0.0039  0.1671  0.5275
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)  6.16317    0.02298 268.246 < 2e-16 ***
## genderMale   0.19265    0.02981   6.462  5.1e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.2366 on 259 degrees of freedom
## Multiple R-squared:  0.1388, Adjusted R-squared:  0.1355
## F-statistic: 41.75 on 1 and 259 DF,  p-value: 5.103e-10
##

```



```
## Call:
## lm(formula = sal9495 ~ gender + dept, data = lawsuit)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.35272 -0.09716 -0.01250  0.08171  0.30390
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    6.01909    0.02055 292.901 < 2e-16 ***
## genderMale      0.10260    0.01678   6.116 3.60e-09 ***
## deptGenetics    0.10046    0.03301   3.044 0.00258 **
## deptMedicine    0.32471    0.02284  14.215 < 2e-16 ***
## deptPediatrics  0.15267    0.02960   5.158 5.02e-07 ***
## deptPhysiology -0.06597    0.02693  -2.450 0.01496 *
## deptSurgery     0.53834    0.02727  19.743 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1267 on 254 degrees of freedom
## Multiple R-squared:  0.7578, Adjusted R-squared:  0.752
## F-statistic: 132.4 on 6 and 254 DF,  p-value: < 2.2e-16
## Analysis of Variance Table
##
## Model 1: sal9495 ~ gender
## Model 2: sal9495 ~ gender + dept
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1      259 14.4927
## 2      254  4.0765  5     10.416 129.8 < 2.2e-16 ***
```

```

## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Call:
## lm(formula = sal9495 ~ gender + clin, data = lawsuit)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.50456 -0.12512 -0.00948  0.12282  0.49713
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   6.29514    0.02065  304.872 < 2e-16 ***
## genderMale    0.16859    0.02336   7.218 5.89e-12 ***
## clinResearch -0.30410    0.02355 -12.912 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1847 on 258 degrees of freedom
## Multiple R-squared:  0.4769, Adjusted R-squared:  0.4728
## F-statistic: 117.6 on 2 and 258 DF,  p-value: < 2.2e-16
## Analysis of Variance Table
##
## Model 1: sal9495 ~ gender
## Model 2: sal9495 ~ gender + clin
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1     259 14.4927
## 2     258  8.8035  1     5.6891 166.73 < 2.2e-16 ***
## ---

```

```

## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

##
## Call:
## lm(formula = sal9495 ~ gender + cert, data = lawsuit)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.55933 -0.13835 -0.00961  0.15540  0.50971
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      6.25144    0.02228 280.648 < 2e-16 ***
## genderMale        0.16642    0.02617   6.360 9.11e-10 ***
## certNot certified -0.25991    0.02863  -9.078 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.2063 on 258 degrees of freedom
## Multiple R-squared:  0.3473, Adjusted R-squared:  0.3423
## F-statistic: 68.65 on 2 and 258 DF,  p-value: < 2.2e-16

## Analysis of Variance Table
##
## Model 1: sal9495 ~ gender
## Model 2: sal9495 ~ gender + cert
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1      259 14.493
## 2      258 10.984   1    3.5087 82.416 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

```
##
## Call:
## lm(formula = sal9495 ~ gender + exper, data = lawsuit)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.51649 -0.18686  0.02018  0.16638  0.51473
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  6.100037   0.029194 208.950  < 2e-16 ***
## genderMale   0.153775   0.031384   4.900  1.7e-06 ***
## exper        0.008428   0.002480   3.399  0.000784 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.2319 on 258 degrees of freedom
## Multiple R-squared:  0.1757, Adjusted R-squared:  0.1693
## F-statistic: 27.5 on 2 and 258 DF,  p-value: 1.488e-11
## Analysis of Variance Table
##
## Model 1: sal9495 ~ gender
## Model 2: sal9495 ~ gender + exper
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1      259 14.493
## 2      258 13.872  1    0.62108 11.552 0.000784 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
```

```
## Call:
## lm(formula = sal9495 ~ gender + rank, data = lawsuit)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.55399 -0.18936  0.00746  0.17699  0.51766
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      6.154393    0.025482 241.520 < 2e-16 ***
## genderMale        0.174740    0.032210   5.425 1.34e-07 ***
## rankAssociate    -0.002612    0.038094  -0.069  0.9454
## rankFull professor  0.061583    0.036611   1.682  0.0938 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.2358 on 257 degrees of freedom
## Multiple R-squared:  0.151, Adjusted R-squared:  0.1411
## F-statistic: 15.23 on 3 and 257 DF,  p-value: 3.747e-09
## Analysis of Variance Table
##
## Model 1: sal9495 ~ gender
## Model 2: sal9495 ~ gender + rank
##   Res.Df    RSS Df Sum of Sq    F Pr(>F)
## 1     259 14.493
## 2     257 14.288  2    0.20422 1.8366 0.1614
```

There is no need to consider prate since the correlation of prate and dept is -0.87. Otherwise there might be collinearity.

From linear regression results, rank is not significant. From partial F test, p is more than 0.05, which shows that

model with rank is not superior.

Since the changes of estimate of genderMale are more than 10% and all p-values from partial F test are less than 0.05, dept, clin, cert and exper are all confounders. So we need to add these covariates into the model.

First we need to add department....

```
##
## Call:
## lm(formula = sal9495 ~ gender + dept, data = lawsuit)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.35272 -0.09716 -0.01250  0.08171  0.30390
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    6.01909    0.02055 292.901 < 2e-16 ***
## genderMale     0.10260    0.01678   6.116 3.60e-09 ***
## deptGenetics   0.10046    0.03301   3.044 0.00258 **
## deptMedicine   0.32471    0.02284  14.215 < 2e-16 ***
## deptPediatrics 0.15267    0.02960   5.158 5.02e-07 ***
## deptPhysiology -0.06597    0.02693  -2.450 0.01496 *
## deptSurgery    0.53834    0.02727  19.743 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1267 on 254 degrees of freedom
## Multiple R-squared:  0.7578, Adjusted R-squared:  0.752
## F-statistic: 132.4 on 6 and 254 DF,  p-value: < 2.2e-16
##
```

```
## Call:
## lm(formula = sal9495 ~ gender + dept + clin, data = lawsuit)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.315740 -0.083385 -0.003147  0.075969  0.309553
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    6.12952    0.02539 241.453 < 2e-16 ***
## genderMale      0.10156    0.01554   6.537 3.44e-10 ***
## deptGenetics    0.04401    0.03175   1.386 0.166977
## deptMedicine    0.23677    0.02503   9.460 < 2e-16 ***
## deptPediatrics  0.07587    0.02980   2.546 0.011481 *
## deptPhysiology -0.09789    0.02540  -3.854 0.000148 ***
## deptSurgery     0.43195    0.02999  14.401 < 2e-16 ***
## clinResearch   -0.12477    0.01898  -6.572 2.81e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.1173 on 253 degrees of freedom
## Multiple R-squared:  0.7931, Adjusted R-squared:  0.7874
## F-statistic: 138.5 on 7 and 253 DF,  p-value: < 2.2e-16
## Analysis of Variance Table
##
## Model 1: sal9495 ~ gender + dept
## Model 2: sal9495 ~ gender + dept + clin
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1      254 4.0765
```

```
## 2      253 3.4821  1    0.59448 43.194 2.811e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

##
## Call:
## lm(formula = sal9495 ~ gender + dept + cert, data = lawsuit)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.278961 -0.084532 -0.001179  0.077452  0.274284
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      6.09277    0.02079  293.077 < 2e-16 ***
## genderMale        0.09339    0.01515   6.164 2.78e-09 ***
## deptGenetics      0.11555    0.02978   3.880 0.000133 ***
## deptMedicine      0.28065    0.02133  13.156 < 2e-16 ***
## deptPediatrics    0.09115    0.02780   3.279 0.001187 **
## deptPhysiology    -0.07712    0.02428  -3.176 0.001681 **
## deptSurgery       0.47954    0.02569  18.668 < 2e-16 ***
## certNot certified -0.13629    0.01754  -7.771 1.96e-13 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

##
## Residual standard error: 0.1141 on 253 degrees of freedom
## Multiple R-squared:  0.8044, Adjusted R-squared:  0.799
## F-statistic: 148.7 on 7 and 253 DF,  p-value: < 2.2e-16

## Analysis of Variance Table
##
```



```

## Model 1: sal9495 ~ gender + dept
## Model 2: sal9495 ~ gender + dept + cert
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1     254 4.0765
## 2     253 3.2911  1    0.78546 60.382 1.959e-13 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

##
## Call:
## lm(formula = sal9495 ~ gender + dept + exper, data = lawsuit)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.261287 -0.066217  0.000115  0.061314  0.277094
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    5.877078   0.018514 317.446 < 2e-16 ***
## genderMale      0.034869   0.013529   2.577   0.0105 *
## deptGenetics    0.132781   0.024959   5.320 2.29e-07 ***
## deptMedicine    0.368417   0.017480  21.076 < 2e-16 ***
## deptPediatrics  0.209135   0.022647   9.235 < 2e-16 ***
## deptPhysiology -0.044616   0.020330  -2.195   0.0291 *
## deptSurgery     0.584001   0.020788  28.093 < 2e-16 ***
## exper           0.014613   0.001046  13.968 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

##
## Residual standard error: 0.09538 on 253 degrees of freedom

```

```
## Multiple R-squared:  0.8632, Adjusted R-squared:  0.8594
## F-statistic: 228.1 on 7 and 253 DF,  p-value: < 2.2e-16

## Analysis of Variance Table

##

## Model 1: sal9495 ~ gender + dept
## Model 2: sal9495 ~ gender + dept + exper

##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1     254 4.0765
## 2     253 2.3017  1     1.7749 195.09 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Result: Compared with other covariates, F value of adding exper is the largest. So we add exper into our model.

```
##

## Call:
## lm(formula = sal9495 ~ gender + dept + exper, data = lawsuit)
##

## Residuals:

##      Min       1Q   Median       3Q      Max
## -0.261287 -0.066217  0.000115  0.061314  0.277094
##

## Coefficients:

##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    5.877078   0.018514 317.446 < 2e-16 ***
## genderMale      0.034869   0.013529   2.577  0.0105 *
## deptGenetics    0.132781   0.024959   5.320 2.29e-07 ***
## deptMedicine    0.368417   0.017480  21.076 < 2e-16 ***
## deptPediatrics  0.209135   0.022647   9.235 < 2e-16 ***
## deptPhysiology -0.044616   0.020330  -2.195  0.0291 *
## deptSurgery     0.584001   0.020788  28.093 < 2e-16 ***
```

```
## exper          0.014613    0.001046  13.968  < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.09538 on 253 degrees of freedom
## Multiple R-squared:  0.8632, Adjusted R-squared:  0.8594
## F-statistic: 228.1 on 7 and 253 DF,  p-value: < 2.2e-16
```

Then we try to look for next covariate.

```
##
## Call:
## lm(formula = sal9495 ~ gender + dept + exper + clin, data = lawsuit)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.23172 -0.05501  0.00399  0.05391  0.37148
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)   5.9835666  0.0204581 292.479  < 2e-16 ***
## genderMale     0.0352073  0.0119141   2.955 0.003422 **
## deptGenetics   0.0791242  0.0228446   3.464 0.000626 ***
## deptMedicine   0.2849665  0.0181866  15.669  < 2e-16 ***
## deptPediatrics 0.1358926  0.0216792   6.268 1.57e-09 ***
## deptPhysiology -0.0750153  0.0182475  -4.111 5.33e-05 ***
## deptSurgery    0.4831773  0.0217262  22.239  < 2e-16 ***
## exper          0.0143278  0.0009219  15.542  < 2e-16 ***
## clinResearch  -0.1171971  0.0136012  -8.617 7.70e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
## Residual standard error: 0.08399 on 252 degrees of freedom
## Multiple R-squared:  0.8944, Adjusted R-squared:  0.891
## F-statistic: 266.7 on 8 and 252 DF,  p-value: < 2.2e-16

## Analysis of Variance Table

##

## Model 1: sal9495 ~ gender + dept + exper
## Model 2: sal9495 ~ gender + dept + exper + clin
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1      253 2.3017
## 2      252 1.7779   1    0.52381 74.247 7.696e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

##

## Call:
## lm(formula = sal9495 ~ gender + dept + exper + cert, data = lawsuit)
##

## Residuals:

##      Min       1Q   Median       3Q      Max
## -0.22497 -0.05412  0.00002  0.05689  0.35517

##

## Coefficients:

##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    5.9461021   0.0187608  316.943 < 2e-16 ***
## genderMale      0.0334723   0.0121359   2.758 0.00624 **
## deptGenetics    0.1417147   0.0224149   6.322 1.16e-09 ***
## deptMedicine    0.3306264   0.0163909  20.171 < 2e-16 ***
## deptPediatrics  0.1567611   0.0213655   7.337 2.99e-12 ***
## deptPhysiology -0.0550454   0.0182822  -3.011 0.00287 **
```

```

## deptSurgery          0.5346466  0.0196625  27.191  < 2e-16 ***
## exper                0.0133753  0.0009513  14.060  < 2e-16 ***
## certNot certified -0.1054425  0.0133379  -7.905  8.35e-14 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.08555 on 252 degrees of freedom
## Multiple R-squared:  0.8904, Adjusted R-squared:  0.8869
## F-statistic: 255.9 on 8 and 252 DF,  p-value: < 2.2e-16

## Analysis of Variance Table
##
## Model 1: sal9495 ~ gender + dept + exper
## Model 2: sal9495 ~ gender + dept + exper + cert
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1      253 2.3017
## 2      252 1.8443  1    0.45739 62.497 8.354e-14 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Result: Compared with other covariates, F value of adding clin is the largest. So we add exper into our model.

```

##
## Call:
## lm(formula = sal9495 ~ gender + dept + exper + clin, data = lawsuit)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.23172 -0.05501  0.00399  0.05391  0.37148
##
## Coefficients:
##
##              Estimate Std. Error t value Pr(>|t|)

```

```
## (Intercept)      5.9835666  0.0204581 292.479 < 2e-16 ***
## genderMale       0.0352073  0.0119141   2.955 0.003422 **
## deptGenetics     0.0791242  0.0228446   3.464 0.000626 ***
## deptMedicine     0.2849665  0.0181866  15.669 < 2e-16 ***
## deptPediatrics   0.1358926  0.0216792   6.268 1.57e-09 ***
## deptPhysiology  -0.0750153  0.0182475  -4.111 5.33e-05 ***
## deptSurgery      0.4831773  0.0217262  22.239 < 2e-16 ***
## exper            0.0143278  0.0009219  15.542 < 2e-16 ***
## clinResearch     -0.1171971  0.0136012  -8.617 7.70e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.08399 on 252 degrees of freedom
## Multiple R-squared:  0.8944, Adjusted R-squared:  0.891
## F-statistic: 266.7 on 8 and 252 DF,  p-value: < 2.2e-16
```

Then we try to look for next covariate :

```
##
## Call:
## lm(formula = sal9495 ~ gender + dept + exper + clin + cert, data = lawsuit)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.18571 -0.05097  0.00044  0.04429  0.42863
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    6.0314905   0.0194750  309.705 < 2e-16 ***
## genderMale     0.0339569   0.0107443   3.160 0.00177 **
## deptGenetics   0.0929245   0.0206775   4.494 1.07e-05 ***
```

```

## deptMedicine      0.2615613  0.0166801  15.681  < 2e-16 ***
## deptPediatrics    0.0986776  0.0201405   4.899  1.72e-06 ***
## deptPhysiology    -0.0806445  0.0164703  -4.896  1.75e-06 ***
## deptSurgery       0.4516913  0.0200155  22.567  < 2e-16 ***
## exper             0.0132859  0.0008423  15.774  < 2e-16 ***
## clinResearch      -0.1039933  0.0123844  -8.397  3.40e-15 ***
## certNot certified -0.0915366  0.0119239  -7.677  3.63e-13 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.07574 on 251 degrees of freedom
## Multiple R-squared:  0.9144, Adjusted R-squared:  0.9114
## F-statistic: 298.1 on 9 and 251 DF,  p-value: < 2.2e-16

## Analysis of Variance Table
##
## Model 1: sal9495 ~ gender + dept + exper + clin
## Model 2: sal9495 ~ gender + dept + exper + clin + cert
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1      252 1.7779
## 2      251 1.4398   1    0.33805 58.932 3.632e-13 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Result:

p value of partial F test is more than 0.05, which shows that the model adding cert is superior. Therefore, we need to add cert in the model.

**Final main effect model:**

```

##
## Call:

```

```
## lm(formula = sal9495 ~ gender + dept + exper + clin + cert, data = lawsuit)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.18571 -0.05097  0.00044  0.04429  0.42863
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    6.0314905   0.0194750  309.705 < 2e-16 ***
## genderMale      0.0339569   0.0107443   3.160  0.00177 **
## deptGenetics    0.0929245   0.0206775   4.494 1.07e-05 ***
## deptMedicine    0.2615613   0.0166801  15.681 < 2e-16 ***
## deptPediatrics  0.0986776   0.0201405   4.899 1.72e-06 ***
## deptPhysiology -0.0806445   0.0164703  -4.896 1.75e-06 ***
## deptSurgery     0.4516913   0.0200155  22.567 < 2e-16 ***
## exper           0.0132859   0.0008423  15.774 < 2e-16 ***
## clinResearch    -0.1039933   0.0123844  -8.397 3.40e-15 ***
## certNot certified -0.0915366   0.0119239  -7.677 3.63e-13 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.07574 on 251 degrees of freedom
## Multiple R-squared:  0.9144, Adjusted R-squared:  0.9114
## F-statistic: 298.1 on 9 and 251 DF,  p-value: < 2.2e-16

add interaction

##
## Call:
## lm(formula = sal9495 ~ gender * dept + exper + clin + cert, data = lawsuit)
##
```



```

## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.18564 -0.04880  0.00124  0.04409  0.42850
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      6.0325227   0.0227388 265.297 < 2e-16 ***
## genderMale        0.0302541   0.0223541   1.353 0.177168
## deptGenetics      0.0833594   0.0290913   2.865 0.004524 **
## deptMedicine      0.2637483   0.0240299  10.976 < 2e-16 ***
## deptPediatrics    0.0976654   0.0259728   3.760 0.000212 ***
## deptPhysiology    -0.0920384   0.0244470  -3.765 0.000209 ***
## deptSurgery       0.4716680   0.0400878  11.766 < 2e-16 ***
## exper            0.0132652   0.0008508  15.591 < 2e-16 ***
## clinResearch      -0.1034518   0.0126267  -8.193 1.40e-14 ***
## certNot certified -0.0895941   0.0121654  -7.365 2.67e-12 ***
## genderMale:deptGenetics  0.0190307   0.0406550   0.468 0.640125
## genderMale:deptMedicine -0.0018252   0.0282104  -0.065 0.948465
## genderMale:deptPediatrics 0.0032796   0.0368694   0.089 0.929191
## genderMale:deptPhysiology 0.0225352   0.0328368   0.686 0.493183
## genderMale:deptSurgery  -0.0201824   0.0426198  -0.474 0.636244
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.0763 on 246 degrees of freedom
## Multiple R-squared:  0.9149, Adjusted R-squared:  0.9101
## F-statistic: 188.9 on 14 and 246 DF, p-value: < 2.2e-16
## Analysis of Variance Table
##

```

```
## Model 1: sal9495 ~ gender + dept + exper + clin + cert
```

```
## Model 2: sal9495 ~ gender * dept + exper + clin + cert
```

```
##   Res.Df    RSS Df Sum of Sq    F Pr(>F)
```

```
## 1      251 1.4398
```

```
## 2      246 1.4320  5 0.0078142 0.2685  0.93
```

Interaction terms are not significant since p values are more than 0.05.

And partial F test p value > 0.05

no interaction between gender and department.

```
##
```

```
## Call:
```

```
## lm(formula = sal9495 ~ gender * exper + dept + clin + cert, data = lawsuit)
```

```
##
```

```
## Residuals:
```

```
##      Min       1Q   Median       3Q      Max
```

```
## -0.17887 -0.04469  0.00158  0.04131  0.40799
```

```
##
```

```
## Coefficients:
```

```
##              Estimate Std. Error t value Pr(>|t|)
```

```
## (Intercept)      5.974101   0.023092 258.710 < 2e-16 ***
```

```
## genderMale        0.105216   0.019578   5.374 1.76e-07 ***
```

```
## exper             0.020034   0.001770  11.319 < 2e-16 ***
```

```
## deptGenetics      0.095713   0.020005   4.784 2.93e-06 ***
```

```
## deptMedicine      0.267482   0.016188  16.524 < 2e-16 ***
```

```
## deptPediatrics    0.108945   0.019621   5.552 7.19e-08 ***
```

```
## deptPhysiology    -0.073740   0.016007  -4.607 6.52e-06 ***
```

```
## deptSurgery       0.457157   0.019396  23.569 < 2e-16 ***
```

```
## clinResearch      -0.104183   0.011975  -8.700 4.53e-16 ***
```

```
## certNot certified -0.087238   0.011573  -7.538 8.78e-13 ***
```

```
## genderMale:exper  -0.008377   0.001951  -4.294 2.51e-05 ***
```

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.07324 on 250 degrees of freedom
## Multiple R-squared:  0.9203, Adjusted R-squared:  0.9171
## F-statistic: 288.8 on 10 and 250 DF,  p-value: < 2.2e-16

## Analysis of Variance Table
##
## Model 1: sal9495 ~ gender + dept + exper + clin + cert
## Model 2: sal9495 ~ gender * exper + dept + clin + cert
##   Res.Df    RSS Df Sum of Sq    F  Pr(>F)
## 1      251 1.4398
## 2      250 1.3409  1  0.098909 18.441 2.51e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

genderMale:exper is significant since p value is less than 0.05.

partial F test p value < 0.05 interaction between gender and exper is significant.

##
## Call:
## lm(formula = sal9495 ~ gender * clin + exper + cert + dept, data = lawsuit)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.19595 -0.04838 -0.00007  0.04304  0.41441
##
## Coefficients:
##
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      6.0430697   0.0196756 307.135  < 2e-16 ***
```

```

## genderMale          0.0106440  0.0135578   0.785  0.43315
## clinResearch        -0.1324656  0.0159951  -8.282  7.40e-15 ***
## exper               0.0132387  0.0008316  15.921  < 2e-16 ***
## certNot certified   -0.0949870  0.0118359  -8.025  3.95e-14 ***
## deptGenetics         0.1019910  0.0206728   4.934  1.47e-06 ***
## deptMedicine         0.2643889  0.0164962  16.027  < 2e-16 ***
## deptPediatrics       0.0975262  0.0198844   4.905  1.69e-06 ***
## deptPhysiology      -0.0787295  0.0162721  -4.838  2.29e-06 ***
## deptSurgery          0.4605185  0.0200139  23.010  < 2e-16 ***
## genderMale:clinResearch 0.0551454  0.0199793   2.760  0.00621 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.07476 on 250 degrees of freedom
## Multiple R-squared:  0.917, Adjusted R-squared:  0.9137
## F-statistic: 276.1 on 10 and 250 DF, p-value: < 2.2e-16

## Analysis of Variance Table

##
## Model 1: sal9495 ~ gender + dept + exper + clin + cert
## Model 2: sal9495 ~ gender * clin + exper + cert + dept
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1      251  1.4398
## 2      250  1.3972   1  0.042578 7.6183 0.006206 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

??? wird no main effect

##
## Call:
## lm(formula = sal9495 ~ gender * cert + exper + clin + dept, data = lawsuit)

```

```
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.18588 -0.05139  0.00078  0.04455  0.42827
##
## Coefficients:
##
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      6.031796   0.020117  299.842 < 2e-16 ***
## genderMale        0.033539   0.012668   2.647  0.00863 **
## certNot certified -0.092297   0.017038  -5.417  1.42e-07 ***
## exper            0.013285   0.000844  15.740 < 2e-16 ***
## clinResearch     -0.104079   0.012485  -8.336  5.16e-15 ***
## deptGenetics      0.093072   0.020853   4.463  1.22e-05 ***
## deptMedicine      0.261518   0.016727  15.634 < 2e-16 ***
## deptPediatrics    0.098545   0.020291   4.856  2.11e-06 ***
## deptPhysiology    -0.080561   0.016556  -4.866  2.02e-06 ***
## deptSurgery       0.451731   0.020066  22.513 < 2e-16 ***
## genderMale:certNot certified 0.001370   0.021885   0.063  0.95013
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.07589 on 250 degrees of freedom
## Multiple R-squared:  0.9144, Adjusted R-squared:  0.911
## F-statistic: 267.2 on 10 and 250 DF,  p-value: < 2.2e-16
## Analysis of Variance Table
##
## Model 1: sal9495 ~ gender + dept + exper + clin + cert
## Model 2: sal9495 ~ gender * cert + exper + clin + dept
##   Res.Df    RSS Df Sum of Sq    F Pr(>F)
```

```
## 1    251 1.4398
## 2    250 1.4398  1 2.2572e-05 0.0039 0.9501
```

Interaction term is not significant since p value is more than 0.05.

And partial F test p value > 0.05

no interaction between gender and cert.

---

### Add Rank

```
##
```

```
## Call:
```

```
## lm(formula = sal9495 ~ gender + rank + exper + cert + clin +
##      dept, data = lawsuit)
```

```
##
```

```
## Residuals:
```

```
##      Min      1Q   Median      3Q      Max
## -0.17303 -0.03848 -0.00936  0.03798  0.45196
```

```
##
```

```
## Coefficients:
```

```
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      6.033505   0.017199 350.801 < 2e-16 ***
## genderMale        0.012882   0.009812   1.313    0.19
## rankAssociate      0.067332   0.011779   5.716 3.10e-08 ***
## rankFull professor 0.111107   0.013125   8.466 2.22e-15 ***
## exper             0.008863   0.000906   9.783 < 2e-16 ***
## certNot certified -0.094875   0.010622  -8.932 < 2e-16 ***
## clinResearch      -0.104170   0.010943  -9.520 < 2e-16 ***
## deptGenetics       0.092985   0.018250   5.095 6.90e-07 ***
## deptMedicine       0.269652   0.014758  18.272 < 2e-16 ***
## deptPediatrics     0.101673   0.017856   5.694 3.48e-08 ***
```

```

## deptPhysiology      -0.087875    0.014561   -6.035 5.73e-09 ***
## deptSurgery         0.466910    0.017767   26.280 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.06685 on 249 degrees of freedom
## Multiple R-squared:  0.9339, Adjusted R-squared:  0.931
## F-statistic: 319.7 on 11 and 249 DF,  p-value: < 2.2e-16

grank interaction

##
## Call:
## lm(formula = sal9495 ~ gender * rank + exper + cert + clin +
##      dept, data = lawsuit)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.16333 -0.04040 -0.00537  0.03823  0.43343
##
## Coefficients:
##
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      6.0203633   0.0179001  336.331 < 2e-16 ***
## genderMale        0.0372396   0.0137840    2.702  0.00738 **
## rankAssociate     0.0865708   0.0169521    5.107 6.55e-07 ***
## rankFull professor 0.1411406   0.0197970    7.129 1.11e-11 ***
## exper            0.0090856   0.0009028   10.064 < 2e-16 ***
## certNot certified -0.0956067   0.0106813   -8.951 < 2e-16 ***
## clinResearch     -0.0985156   0.0110875   -8.885 < 2e-16 ***
## deptGenetics      0.0922859   0.0181030    5.098 6.84e-07 ***
## deptMedicine      0.2716021   0.0146819   18.499 < 2e-16 ***

```

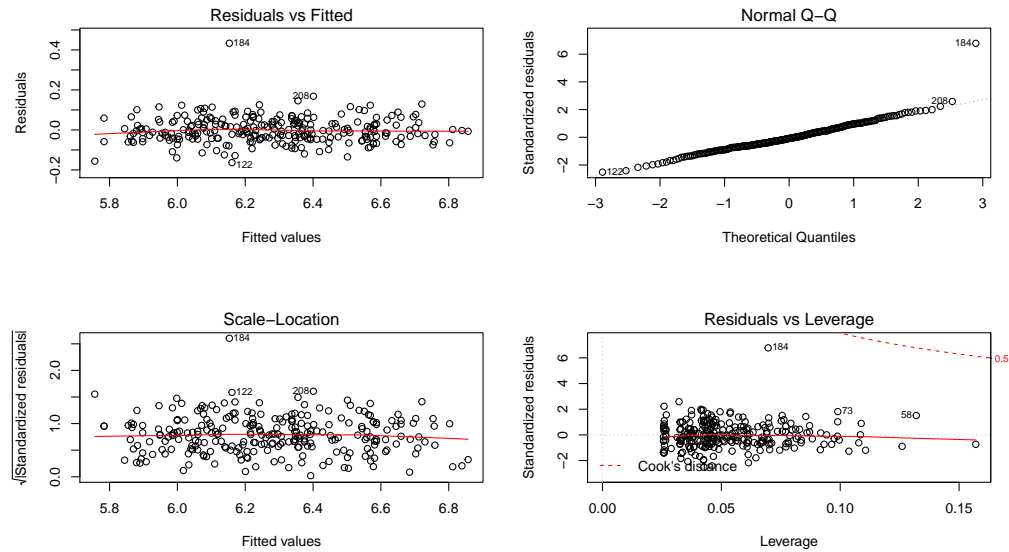
```

## deptPediatrics          0.1042339  0.0177642   5.868 1.41e-08 ***
## deptPhysiology         -0.0877718  0.0144357  -6.080 4.53e-09 ***
## deptSurgery            0.4656938  0.0176337  26.409 < 2e-16 ***
## genderMale:rankAssociate -0.0414716  0.0223749  -1.853 0.06501 .
## genderMale:rankFull professor -0.0526354  0.0233270  -2.256 0.02492 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.06627 on 247 degrees of freedom
## Multiple R-squared:  0.9355, Adjusted R-squared:  0.9322
## F-statistic: 275.8 on 13 and 247 DF,  p-value: < 2.2e-16

## Analysis of Variance Table
##
## Model 1: sal9495 ~ gender + rank + exper + cert + clin + dept
## Model 2: sal9495 ~ gender * rank + exper + cert + clin + dept
##   Res.Df    RSS Df Sum of Sq    F Pr(>F)
## 1      249 1.1126
## 2      247 1.0847  2  0.027927 3.1796 0.04331 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```





## Potentially influential observations of

## `lm(formula = sal9495 ~ gender * rank + exper + cert + clin + dept, data = lawsuit) :`

##

##      dfb.1\_ dfb.gndM dfb.rnkA dfb.rnFp dfb.expr dfb.crNc dfb.clnR dfb.dptG

## 19    0.08   -0.01     0.05     0.09    -0.27    -0.08    -0.01     0.04

## 39   -0.02    0.00     0.00    -0.02     0.00     0.01     0.01     0.01

## 82   -0.02    0.13     0.14     0.11     0.08    -0.18    -0.10     0.00

## 109   0.00    0.00     0.00     0.00     0.00     0.00     0.00     0.00

## 122   -0.17    0.15     0.15     0.11     0.04     0.01     0.13     0.04

## 184   -0.85    1.15\_\*    0.19     0.24    -0.26     0.83    1.06\_\*    0.14

## 204    0.15   -0.20    -0.12    -0.12     0.03    -0.06    -0.08    -0.02

## 208    0.10   -0.24    -0.17    -0.20     0.21    -0.04    -0.09    -0.01

##      dfb.dptM dfb.dptPd dfb.dptPh dfb.dptS dfb.gM:A dfb.gM:p dffit

## 19    0.03     0.00     0.04     0.03    -0.02     0.03    -0.32

## 39    0.02     0.01     0.01     0.02     0.00     0.02    -0.03

## 82   -0.07    -0.04    -0.30    -0.11    -0.09    -0.10    -0.54

## 109   0.00     0.00     0.00     0.00     0.00     0.00     0.00

## 122   0.08    -0.25     0.03     0.06    -0.12    -0.13    -0.54

## 184   1.03\_\*    0.58     0.22     0.50    -0.80    -0.66    2.05\_\*

```

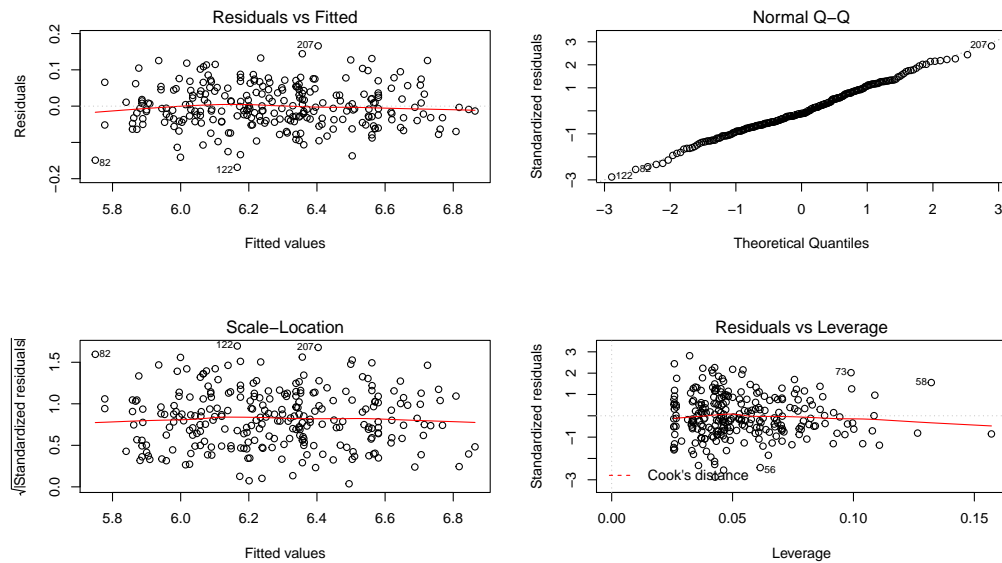
## 204  0.07    -0.07    -0.02    -0.03    0.11    0.11    0.37
## 208  0.11    -0.06    -0.01    -0.02    0.11    0.12    0.48
##      cov.r    cook.d hat
## 19   1.22_*   0.01    0.16
## 39   1.17_*   0.00    0.10
## 82   0.79_*   0.02    0.05
## 109  1.19_*   0.00    0.11
## 122  0.77_*   0.02    0.04
## 184  0.06_*   0.25    0.07
## 204  0.82_*   0.01    0.03
## 208  0.74_*   0.02    0.03

## # A tibble: 8 x 19
##   .rownames    dfb.1_ dfb.gndM dfb.rnkA dfb.rnFp dfb.expr dfb.crNc dfb.clnR
##   <chr>         <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
## 1 19           7.77e-2 -8.32e-3 0.0533  9.07e-2 -2.71e-1 -0.0824 -0.00962
## 2 39          -1.53e-2 1.07e-3 0.00141 -2.07e-2 -9.84e-4 0.00854 0.0103
## 3 82          -2.36e-2 1.26e-1 0.136   1.08e-1 8.33e-2 -0.176  -0.0982
## 4 109         -5.27e-4 6.43e-5 0.00185 1.08e-4 1.44e-4 -0.00112 0.00115
## 5 122         -1.74e-1 1.54e-1 0.145   1.13e-1 3.51e-2 0.0147  0.134
## 6 184         -8.50e-1 1.15e+0 0.187   2.35e-1 -2.56e-1 0.826   1.06
## 7 204          1.49e-1 -2.00e-1 -0.119  -1.23e-1 2.77e-2 -0.0588 -0.0795
## 8 208          9.89e-2 -2.38e-1 -0.173  -2.03e-1 2.13e-1 -0.0422 -0.0877
## # ... with 11 more variables: dfb.dptG <dbl>, dfb.dptM <dbl>,
## #   dfb.dptPd <dbl>, dfb.dptPh <dbl>, dfb.dptS <dbl>, dfb.gM.A <dbl>,
## #   dfb.gM.p <dbl>, dffit <dbl>, cov.r <dbl>, cook.d <dbl>, hat <dbl>

##
## Call:
## lm(formula = sal9495 ~ gender * rank + exper + dept + clin +
##     cert, data = newlawsuit)

```

```
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.168540 -0.037171 -0.006902  0.042172  0.166046
##
## Coefficients:
##                                Estimate Std. Error t value Pr(>|t|)
## (Intercept)                   6.0341149   0.0162845  370.542 < 2e-16 ***
## genderMale                     0.0229280   0.0126055    1.819  0.0701 .
## rankAssociate                  0.0837058   0.0153288    5.461 1.16e-07 ***
## rankFull professor             0.1369305   0.0179046    7.648 4.60e-13 ***
## exper                          0.0092948   0.0008166   11.383 < 2e-16 ***
## deptGenetics                   0.0900340   0.0163672    5.501 9.46e-08 ***
## deptMedicine                   0.2578824   0.0133973   19.249 < 2e-16 ***
## deptPediatrics                 0.0948508   0.0161068    5.889 1.27e-08 ***
## deptPhysiology                 -0.0905967   0.0130548   -6.940 3.46e-11 ***
## deptSurgery                    0.4577307   0.0159755   28.652 < 2e-16 ***
## clinResearch                   -0.1091391   0.0101222  -10.782 < 2e-16 ***
## certNot certified              -0.1035838   0.0097139  -10.663 < 2e-16 ***
## genderMale:rankAssociate       -0.0251920   0.0203422   -1.238  0.2167
## genderMale:rankFull professor -0.0387336   0.0211680   -1.830  0.0685 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.0599 on 246 degrees of freedom
## Multiple R-squared:  0.9472, Adjusted R-squared:  0.9445
## F-statistic: 339.8 on 13 and 246 DF,  p-value: < 2.2e-16
```



ge interaction

##

## Call:

```
## lm(formula = sal9495 ~ gender * exper + rank + cert + clin +
##     dept, data = lawsuit)
```

##

## Residuals:

```
##      Min      1Q   Median      3Q      Max
## -0.16065 -0.03930 -0.00493  0.03550  0.43455
```

##

## Coefficients:

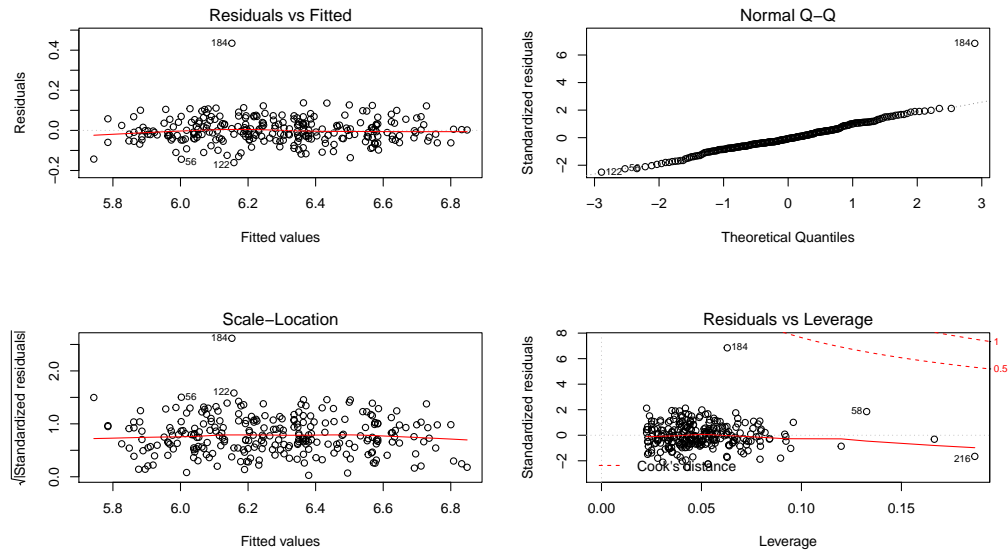
```
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    5.993406   0.020845  287.519 < 2e-16 ***
## genderMale      0.064466   0.018456   3.493 0.000566 ***
## exper           0.013887   0.001773   7.834 1.38e-13 ***
## rankAssociate    0.059115   0.011824   5.000 1.09e-06 ***
## rankFull professor 0.104018   0.013056   7.967 5.90e-14 ***
## certNot certified -0.091083   0.010484  -8.688 5.09e-16 ***
## clinResearch    -0.104087   0.010735  -9.696 < 2e-16 ***
```

```

## deptGenetics      0.094885    0.017913    5.297 2.60e-07 ***
## deptMedicine      0.273386    0.014522   18.825 < 2e-16 ***
## deptPediatrics    0.109302    0.017671    6.185 2.54e-09 ***
## deptPhysiology    -0.082534    0.014377   -5.741 2.75e-08 ***
## deptSurgery       0.469915    0.017453   26.924 < 2e-16 ***
## genderMale:exper  -0.005864    0.001790   -3.276 0.001204 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.06558 on 248 degrees of freedom
## Multiple R-squared:  0.9366, Adjusted R-squared:  0.9336
## F-statistic: 305.4 on 12 and 248 DF,  p-value: < 2.2e-16

## Analysis of Variance Table
##
## Model 1: sal9495 ~ gender + rank + exper + cert + clin + dept
## Model 2: sal9495 ~ gender * exper + rank + cert + clin + dept
##   Res.Df    RSS Df Sum of Sq    F  Pr(>F)
## 1      249 1.1126
## 2      248 1.0665  1  0.046146 10.731 0.001204 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```



## Potentially influential observations of

## `lm(formula = sal9495 ~ gender * exper + rank + cert + clin + dept, data = lawsuit) :`

##

## `dfb.1_ dfb.gndM dfb.expr dfb.rnkA dfb.rnFp dfb.crNc dfb.clnR dfb.dptG`

## 19 0.01 0.03 -0.03 0.01 0.07 -0.03 0.00 0.02

## 58 -0.13 -0.12 0.13 -0.02 -0.35 -0.11 0.13 0.08

## 91 0.01 0.00 0.00 0.00 0.00 0.01 -0.01 -0.03

## 109 0.00 0.00 0.00 0.01 0.00 -0.01 0.01 0.02

## 122 -0.20 0.16 0.14 0.06 0.02 0.03 0.11 0.05

## 184 -0.87 0.93 0.37 -0.63 -0.36 0.89 0.92 0.18

## 216 0.54 -0.39 -0.65 0.14 0.03 -0.03 -0.26 -0.11

## `dfb.dptM dfb.dptPd dfb.dptPh dfb.dptS dfb.gnM: dffit cov.r cook.d`

## 19 0.01 0.01 0.02 0.01 -0.04 -0.14 1.26\_\* 0.00

## 58 0.08 0.07 0.29 0.04 0.17 0.73\_\* 1.01 0.04

## 91 -0.01 0.00 0.00 0.00 -0.01 -0.03 1.16\_\* 0.00

## 109 0.00 0.00 0.00 0.00 0.00 0.03 1.16\_\* 0.00

## 122 0.08 -0.23 0.05 0.07 -0.15 -0.53 0.79\_\* 0.02

## 184 1.03\_\* 0.60 0.28 0.55 -0.64 1.97\_\* 0.07\_\* 0.24

## 216 -0.34 -0.21 -0.13 -0.23 0.62 -0.80\_\* 1.12 0.05

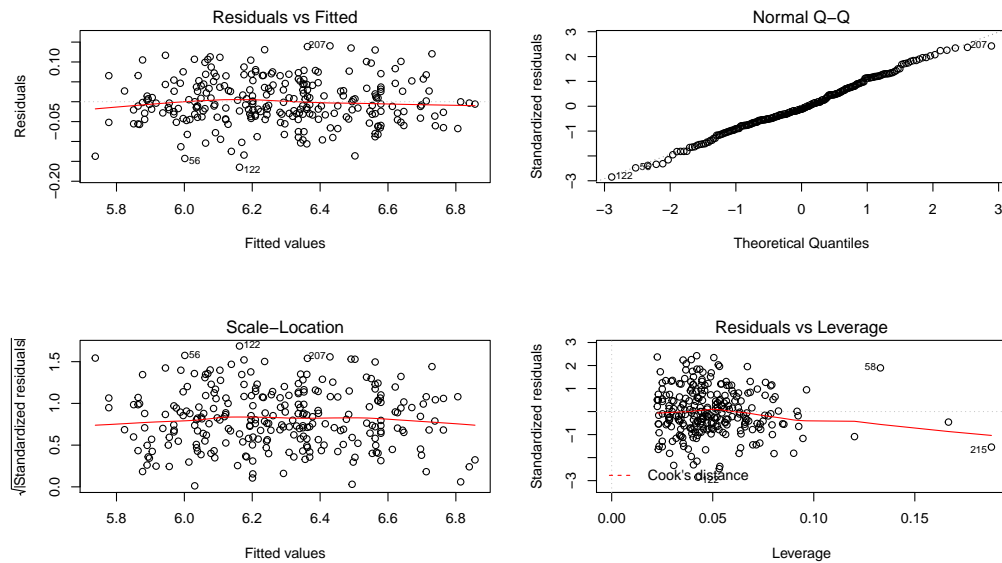
```
##      hat
## 19  0.17_*
## 58  0.13
## 91  0.09
## 109 0.09
## 122 0.04
## 184 0.06
## 216 0.19_*

## # A tibble: 7 x 18
##   .rownames    dfb.1_ dfb.gndM dfb.expr dfb.rnkA dfb.rnFp dfb.crNc dfb.clnR
##   <chr>        <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>   <dbl>
## 1 19          0.00609 3.07e-2 -2.53e-2 0.0123  6.80e-2 -0.0343 -0.00240
## 2 58         -0.130  -1.20e-1 1.26e-1 -0.0151 -3.50e-1 -0.108  0.135
## 3 91          0.00531 4.84e-3 5.65e-5 0.00357 9.21e-4 0.00962 -0.0127
## 4 109        -0.00340 -8.95e-4 3.26e-3 0.00900 9.27e-4 -0.0104 0.00884
## 5 122        -0.196   1.62e-1 1.39e-1 0.0643  1.97e-2 0.0326 0.106
## 6 184        -0.866   9.28e-1 3.73e-1 -0.627  -3.61e-1 0.891  0.925
## 7 216         0.541  -3.90e-1 -6.51e-1 0.144   3.15e-2 -0.0326 -0.264
## # ... with 10 more variables: dfb.dptG <dbl>, dfb.dptM <dbl>,
## #   dfb.dptPd <dbl>, dfb.dptPh <dbl>, dfb.dptS <dbl>, dfb.gnM. <dbl>,
## #   dffit <dbl>, cov.r <dbl>, cook.d <dbl>, hat <dbl>

##
## Call:
## lm(formula = sal9495 ~ gender * exper + rank + dept + clin +
##     cert, data = newlawsuit)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.164921 -0.036222 -0.006406  0.040199  0.140506
```

```
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      6.009705   0.018933 317.415 < 2e-16 ***
## genderMale        0.049018   0.016779   2.921  0.00381 **
## exper             0.013290   0.001601   8.298 6.91e-15 ***
## rankAssociate      0.065807   0.010707   6.146 3.16e-09 ***
## rankFull professor 0.108273   0.011795   9.179 < 2e-16 ***
## deptGenetics       0.091972   0.016170   5.688 3.62e-08 ***
## deptMedicine       0.259894   0.013225  19.652 < 2e-16 ***
## deptPediatrics     0.099663   0.015997   6.230 2.00e-09 ***
## deptPhysiology     -0.086228   0.012983  -6.641 1.96e-10 ***
## deptSurgery        0.461249   0.015791  29.209 < 2e-16 ***
## clinResearch       -0.113043   0.009759 -11.584 < 2e-16 ***
## certNot certified  -0.099511   0.009526 -10.446 < 2e-16 ***
## genderMale:exper   -0.004838   0.001621  -2.984  0.00313 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.05918 on 247 degrees of freedom
## Multiple R-squared:  0.9483, Adjusted R-squared:  0.9458
## F-statistic: 377.6 on 12 and 247 DF, p-value: < 2.2e-16
```





## try adding rank

```
## Analysis of Variance Table
##
## Model 1: sal9495 ~ gender + dept + exper + clin + cert + gender * exper
## Model 2: sal9495 ~ gender + dept + exper + clin + cert + gender * exper +
##      rank
##   Res.Df    RSS Df Sum of Sq    F    Pr(>F)
## 1     250 1.3409
## 2     248 1.0665  2    0.2744 31.904 4.676e-13 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Call:
## lm(formula = sal9495 ~ gender + dept + exper + clin + cert +
##      gender * exper + rank, data = lawsuit)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
```

```
## -0.16065 -0.03930 -0.00493 0.03550 0.43455
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    5.993406   0.020845 287.519 < 2e-16 ***
## genderMale      0.064466   0.018456   3.493 0.000566 ***
## deptGenetics    0.094885   0.017913   5.297 2.60e-07 ***
## deptMedicine    0.273386   0.014522  18.825 < 2e-16 ***
## deptPediatrics  0.109302   0.017671   6.185 2.54e-09 ***
## deptPhysiology -0.082534   0.014377  -5.741 2.75e-08 ***
## deptSurgery     0.469915   0.017453  26.924 < 2e-16 ***
## exper          0.013887   0.001773   7.834 1.38e-13 ***
## clinResearch    -0.104087   0.010735  -9.696 < 2e-16 ***
## certNot certified -0.091083   0.010484  -8.688 5.09e-16 ***
## rankAssociate   0.059115   0.011824   5.000 1.09e-06 ***
## rankFull professor 0.104018   0.013056   7.967 5.90e-14 ***
## genderMale:exper -0.005864   0.001790  -3.276 0.001204 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.06558 on 248 degrees of freedom
## Multiple R-squared:  0.9366, Adjusted R-squared:  0.9336
## F-statistic: 305.4 on 12 and 248 DF,  p-value: < 2.2e-16
##
## Call:
## lm(formula = sal9495 ~ gender + dept + exper + clin + cert +
##     gender * exper, data = lawsuit)
##
## Residuals:
```

```

##      Min      1Q   Median      3Q      Max
## -0.17887 -0.04469  0.00158  0.04131  0.40799
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      5.974101   0.023092 258.710 < 2e-16 ***
## genderMale        0.105216   0.019578   5.374 1.76e-07 ***
## deptGenetics      0.095713   0.020005   4.784 2.93e-06 ***
## deptMedicine      0.267482   0.016188  16.524 < 2e-16 ***
## deptPediatrics    0.108945   0.019621   5.552 7.19e-08 ***
## deptPhysiology    -0.073740   0.016007  -4.607 6.52e-06 ***
## deptSurgery       0.457157   0.019396  23.569 < 2e-16 ***
## exper            0.020034   0.001770  11.319 < 2e-16 ***
## clinResearch      -0.104183   0.011975  -8.700 4.53e-16 ***
## certNot certified -0.087238   0.011573  -7.538 8.78e-13 ***
## genderMale:exper  -0.008377   0.001951  -4.294 2.51e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.07324 on 250 degrees of freedom
## Multiple R-squared:  0.9203, Adjusted R-squared:  0.9171
## F-statistic: 288.8 on 10 and 250 DF,  p-value: < 2.2e-16

```

---

### Final Model with Interaction:

```

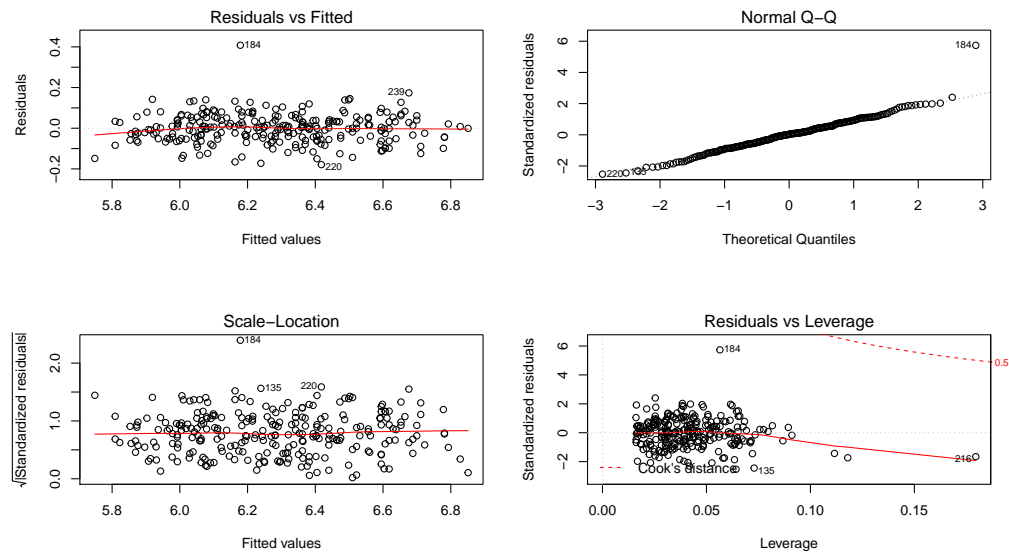
##
## Call:
## lm(formula = sal9495 ~ gender * exper + dept + clin + cert, data = lawsuit)
##
## Residuals:

```

```

##      Min      1Q   Median      3Q      Max
## -0.17887 -0.04469  0.00158  0.04131  0.40799
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      5.974101   0.023092 258.710 < 2e-16 ***
## genderMale        0.105216   0.019578   5.374 1.76e-07 ***
## exper            0.020034   0.001770  11.319 < 2e-16 ***
## deptGenetics      0.095713   0.020005   4.784 2.93e-06 ***
## deptMedicine      0.267482   0.016188  16.524 < 2e-16 ***
## deptPediatrics    0.108945   0.019621   5.552 7.19e-08 ***
## deptPhysiology    -0.073740   0.016007  -4.607 6.52e-06 ***
## deptSurgery       0.457157   0.019396  23.569 < 2e-16 ***
## clinResearch      -0.104183   0.011975  -8.700 4.53e-16 ***
## certNot certified -0.087238   0.011573  -7.538 8.78e-13 ***
## genderMale:exper -0.008377   0.001951  -4.294 2.51e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.07324 on 250 degrees of freedom
## Multiple R-squared:  0.9203, Adjusted R-squared:  0.9171
## F-statistic: 288.8 on 10 and 250 DF,  p-value: < 2.2e-16

```



## Potentially influential observations of

## `lm(formula = sal9495 ~ gender * exper + dept + clin + cert, data = lawsuit)` :

##

##      dfb.1\_ dfb.gndM dfb.expr dfb.dptG dfb.dptM dfb.dptPd dfb.dptPh

## 19    0.00    0.25    -0.01    0.10    0.05    0.01    0.13

## 39    0.00    0.01    0.01    -0.01    -0.01    -0.01    -0.01

## 58   -0.02   -0.04    0.00    0.02    0.02    0.02    0.05

## 91    0.01    0.01    0.00    -0.04    -0.01    -0.01    0.00

## 122 -0.19   0.17    0.16    0.04    0.08    -0.21    0.05

## 135   0.36   -0.18   -0.39   -0.11   -0.21   -0.48   -0.10

## 184 -0.65   0.63    0.13    0.14    0.84    0.45    0.20

## 216   0.53   -0.37   -0.69   -0.11   -0.33   -0.19   -0.12

## 220   0.36   -0.15   -0.41   -0.12   -0.38   -0.16   -0.10

## 239 -0.01   0.01    0.01    0.00    0.00    0.00    0.00

##      dfb.dptS dfb.clnR dfb.crNc dfb.gnM: dffit    cov.r    cook.d hat

## 19    0.04    -0.02    -0.20    -0.25    -0.64\_\*   1.04    0.04    0.12

## 39   -0.01    -0.01    0.00    -0.01    0.02    1.14\_\*   0.00    0.08

## 58    0.02    0.03    -0.02    0.04    0.12    1.14\_\*   0.00    0.09

## 91   -0.01    -0.02    0.02    -0.01    -0.06    1.15\_\*   0.00    0.09

```
## 122 0.06      0.10      0.04      -0.15      -0.49      0.86_* 0.02 0.04
## 135 -0.22     -0.31      0.02      0.33      -0.69_* 0.86_* 0.04 0.07
## 184 0.46      0.73      0.66      -0.42      1.50_* 0.23_* 0.18 0.06
## 216 -0.23     -0.26     -0.01      0.61      -0.78_* 1.13    0.05 0.18_*
## 220 -0.25     -0.37      0.07      0.36      -0.67_* 0.84_* 0.04 0.06
## 239 0.22      -0.01     -0.01      0.00      0.39      0.83_* 0.01 0.03
```

	dfb.l	dfb.gndM	dfb.expr	dfb.dptG	dfb.dptM	dfb.dptPd	dfb.dptPh	dfb.dptS	dfb.clnR	dfb.crNc	dfb.gnM:	dffit	cov.r	cook.d	hat
19	0.0022	0.2483	-0.0053	0.1019	0.0534	0.0093	0.1347	0.0363	-0.0228	-0.2015	-0.2491	-0.6367	1.0372	0.0366	0.1180
39	0.0041	0.0073	0.0145	-0.0087	-0.0129	-0.0105	-0.0091	-0.0113	-0.0115	-0.0034	-0.0139	0.0245	1.1352	0.0001	0.0797
58	-0.0202	-0.0357	0.0018	0.0159	0.0197	0.0182	0.0523	0.0156	0.0284	-0.0172	0.0387	0.1155	1.1408	0.0012	0.0894
91	0.0080	0.0096	0.0017	-0.0408	-0.0088	-0.0070	-0.0040	-0.0068	-0.0200	0.0165	-0.0132	-0.0555	1.1483	0.0003	0.0911
122	-0.1873	0.1720	0.1610	0.0419	0.0786	-0.2098	0.0467	0.0628	0.1003	0.0399	-0.1524	-0.4872	0.8604	0.0212	0.0419
135	0.3633	-0.1842	-0.3857	-0.1101	-0.2134	-0.4827	-0.0987	-0.2173	-0.3129	0.0171	0.3297	-0.6929	0.8637	0.0428	0.0729
184	-0.6471	0.6344	0.1333	0.1445	0.8352	0.4492	0.1991	0.4609	0.7288	0.6551	-0.4150	1.5029	0.2347	0.1790	0.0565
216	0.5306	-0.3743	-0.6869	-0.1128	-0.3318	-0.1917	-0.1228	-0.2286	-0.2570	-0.0090	0.6050	-0.7786	1.1278	0.0547	0.1795
220	0.3636	-0.1548	-0.4125	-0.1235	-0.3818	-0.1584	-0.1024	-0.2476	-0.3676	0.0735	0.3592	-0.6687	0.8406	0.0398	0.0642
239	-0.0093	0.0148	0.0070	0.0028	-0.0043	0.0024	0.0016	0.2221	-0.0071	-0.0114	0.0016	0.3917	0.8291	0.0137	0.0254

## Remove outlier

```
##

## Call:
## lm(formula = sal9495 ~ gender * exper + dept + clin + cert, data = newlawsuit)
##

## Residuals:

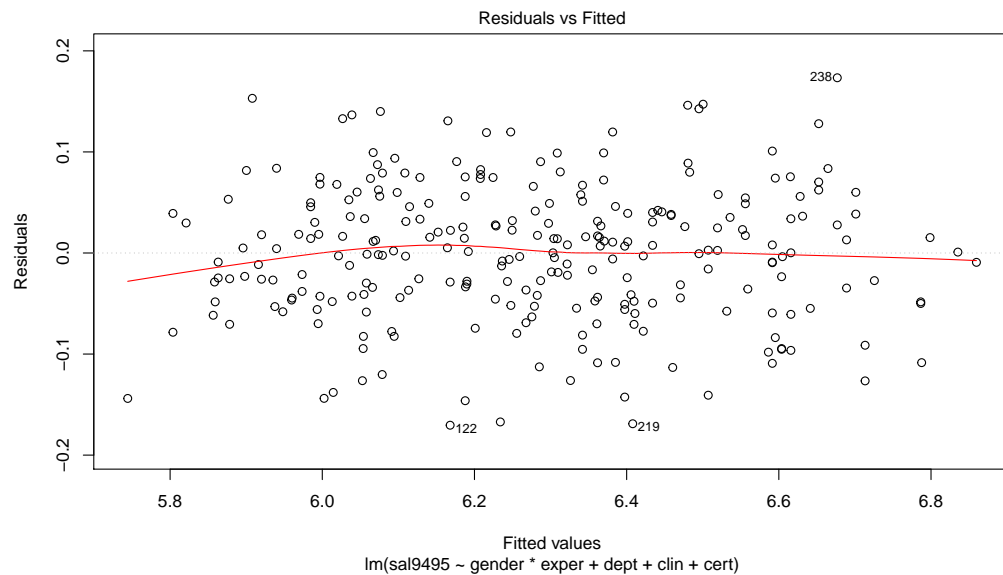
##      Min      1Q   Median      3Q      Max
## -0.170419 -0.044897  0.002285  0.045782  0.173369

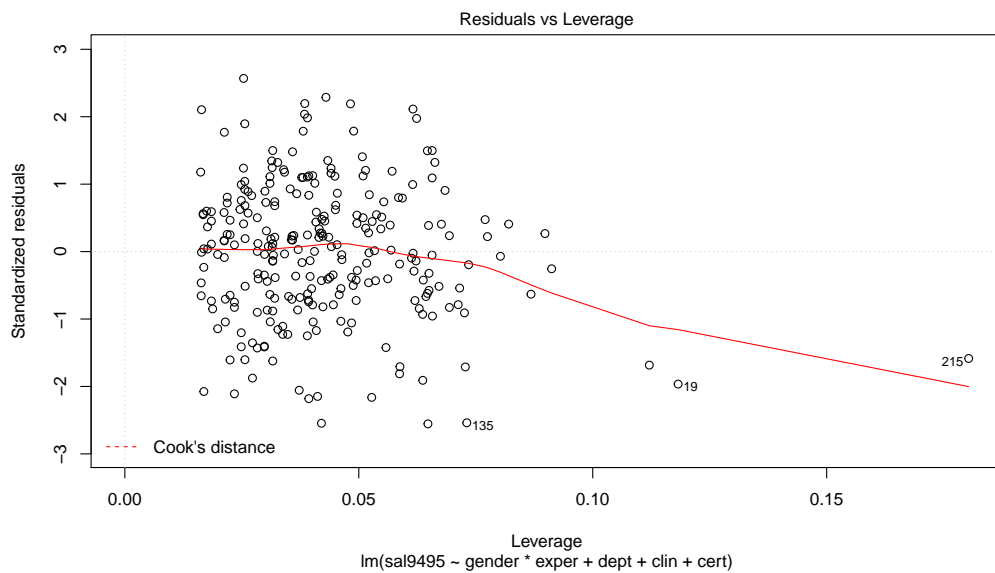
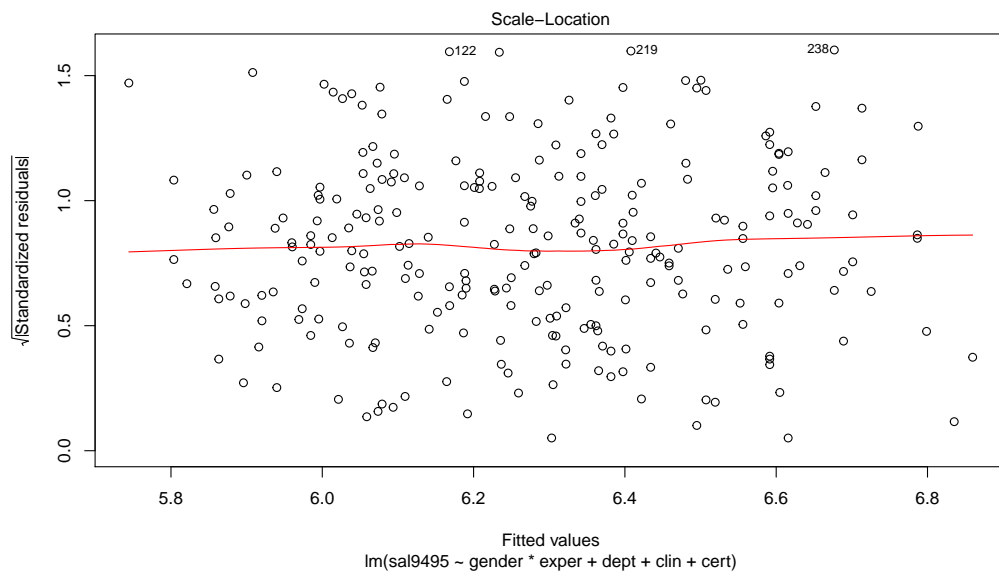
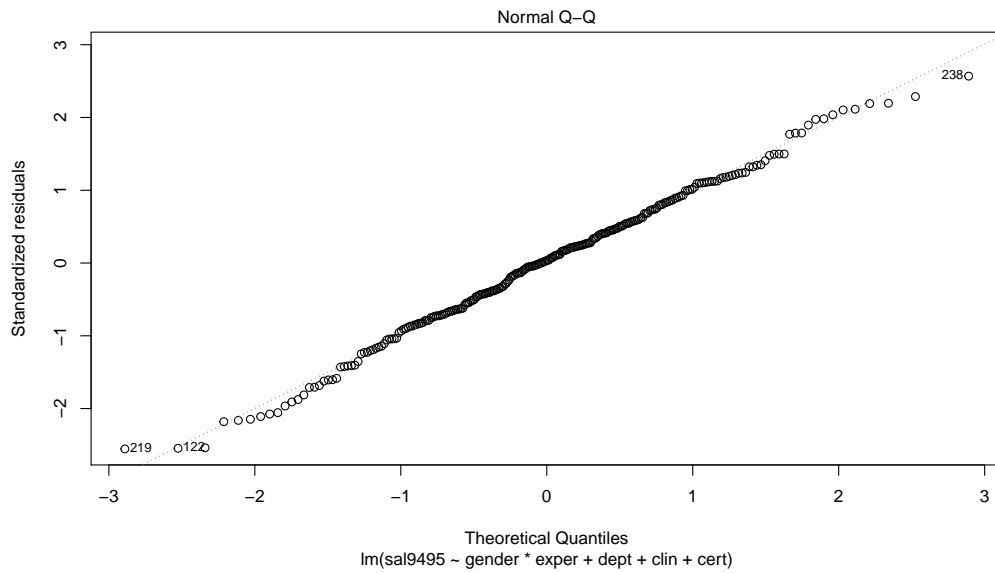
##

## Coefficients:

##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    5.988053   0.021682  276.180 < 2e-16 ***
## genderMale      0.093619   0.018378   5.094 6.93e-07 ***
## exper           0.019814   0.001653  11.986 < 2e-16 ***
## deptGenetics    0.093014   0.018685   4.978 1.20e-06 ***
## deptMedicine    0.254857   0.015255  16.707 < 2e-16 ***
## deptPediatrics  0.100715   0.018371   5.482 1.03e-07 ***
## deptPhysiology -0.076716   0.014955  -5.130 5.84e-07 ***
```

```
## deptSurgery      0.448810    0.018162   24.711   < 2e-16 ***
## clinResearch     -0.112333    0.011261   -9.976   < 2e-16 ***
## certNot certified -0.094318    0.010868   -8.678   5.33e-16 ***
## genderMale:exper -0.007621    0.001826   -4.174   4.13e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.06839 on 249 degrees of freedom
## Multiple R-squared:  0.9304, Adjusted R-squared:  0.9276
## F-statistic: 332.9 on 10 and 249 DF,  p-value: < 2.2e-16
```

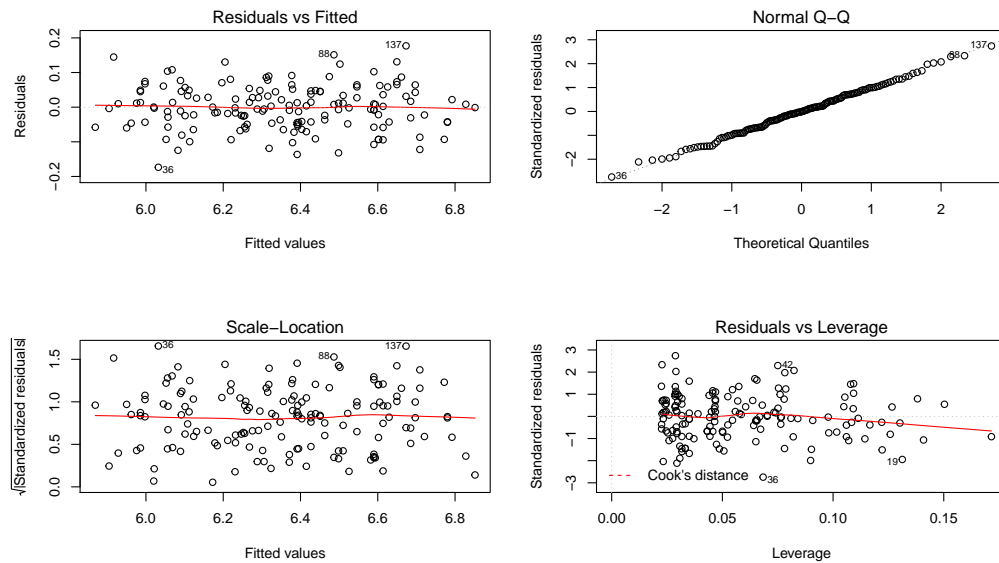






Stratification:

```
##
## Call:
## lm(formula = sal9495 ~ exper + dept + clin + cert, data = male_data)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.173529 -0.043958 -0.000524  0.045339  0.176938
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    6.0738401   0.0222569  272.897 < 2e-16 ***
## exper          0.0118970   0.0008179   14.546 < 2e-16 ***
## deptGenetics    0.1132455   0.0262915    4.307 3.03e-05 ***
## deptMedicine    0.2580777   0.0191627   13.468 < 2e-16 ***
## deptPediatrics  0.0984634   0.0261797    3.761 0.000245 ***
## deptPhysiology -0.0692456   0.0193867   -3.572 0.000482 ***
## deptSurgery     0.4566300   0.0214087   21.329 < 2e-16 ***
## clinResearch   -0.0796016   0.0156645   -5.082 1.14e-06 ***
## certNot certified -0.1156107  0.0140047   -8.255 8.55e-14 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.06552 on 145 degrees of freedom
## Multiple R-squared:  0.9315, Adjusted R-squared:  0.9277
## F-statistic: 246.5 on 8 and 145 DF,  p-value: < 2.2e-16
```



```
##
```

```
## Call:
```

```
## lm(formula = sal9495 ~ exper + dept + clin + cert, data = female_data)
```

```
##
```

```
## Residuals:
```

```
##      Min      1Q   Median      3Q      Max
## -0.175835 -0.039097  0.005635  0.047279  0.167929
```

```
##
```

```
## Coefficients:
```

```
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    5.976891   0.030766 194.267 < 2e-16 ***
## exper          0.020835   0.001774  11.748 < 2e-16 ***
## deptGenetics    0.085369   0.027196   3.139 0.002246 **
## deptMedicine    0.268917   0.025088  10.719 < 2e-16 ***
## deptPediatrics  0.113209   0.027126   4.173 6.55e-05 ***
## deptPhysiology -0.082054   0.023029  -3.563 0.000571 ***
## deptSurgery     0.478852   0.039957  11.984 < 2e-16 ***
## clinResearch   -0.138052   0.016656  -8.289 6.51e-13 ***
## certNot certified -0.068639  0.017731  -3.871 0.000197 ***
```

```
## ---
```

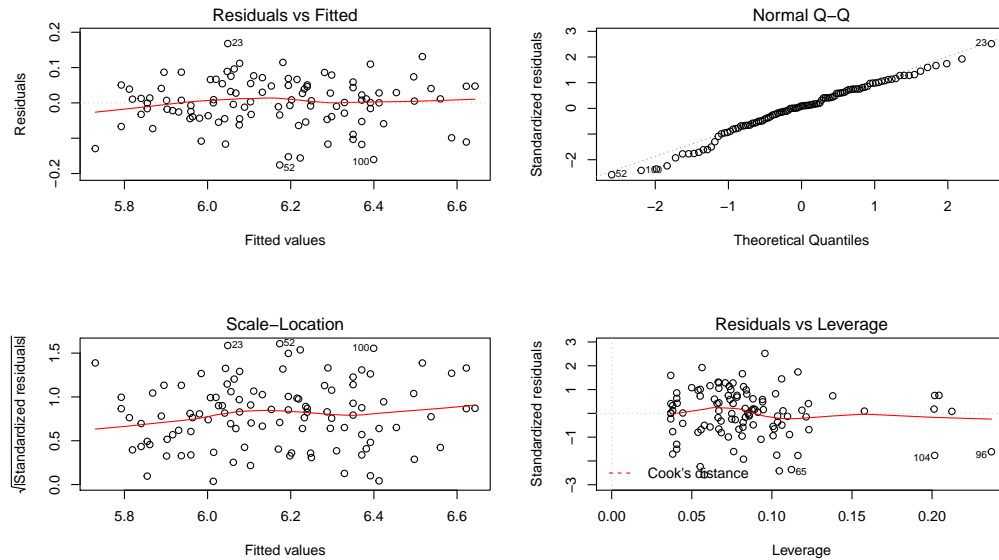
```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
```

```
## Residual standard error: 0.07011 on 97 degrees of freedom
```

```
## Multiple R-squared:  0.9109, Adjusted R-squared:  0.9035
```

```
## F-statistic: 123.9 on 8 and 97 DF,  p-value: < 2.2e-16
```



**Results**

**Conclusions**

**Discussion**

## Figures and tables

	Female (N=106)	Male (N=155)	Total (N=261)	p value
Department				< 0.001
- Biochemistry	20 (18.9%)	30 (19.4%)	50 (19.2%)	
- Genetics	11 (10.4%)	10 (6.5%)	21 (8.0%)	
- Medicine	30 (28.3%)	50 (32.3%)	80 (30.7%)	
- Pediatrics	20 (18.9%)	10 (6.5%)	30 (11.5%)	
- Physiology	20 (18.9%)	20 (12.9%)	40 (15.3%)	
- Surgery	5 (4.7%)	35 (22.6%)	40 (15.3%)	
- Missing	0	0	0	
Clinical				0.197
- Clinical	60 (56.6%)	100 (64.5%)	160 (61.3%)	
- Research	46 (43.4%)	55 (35.5%)	101 (38.7%)	
- Missing	0	0	0	
Certified				0.074
- Board certified	70 (66.0%)	118 (76.1%)	188 (72.0%)	
- Not certified	36 (34.0%)	37 (23.9%)	73 (28.0%)	
- Missing	0	0	0	
Publication Rate				0.004
- Mean (SD)	5.350 (1.886)	4.646 (1.938)	4.932 (1.944)	
- Median (IQR)	5.250 (3.725, 7.275)	4.000 (3.100, 6.700)	4.400 (3.200, 6.900)	
- Min - Max	2.400 - 8.700	1.300 - 8.600	1.300 - 8.700	
- Missing	0	0	0	
Years since obtaining				< 0.001
MD				
- Mean (SD)	7.491 (4.166)	12.103 (6.704)	10.230 (6.227)	
- Median (IQR)	7.000 (5.000, 10.000)	10.000 (7.000, 15.000)	9.000 (6.000, 14.000)	
- Min - Max	1.000 - 23.000	2.000 - 37.000	1.000 - 37.000	
- Missing	0	0	0	
Rank				< 0.001

(continued)

	Female (N=106)	Male (N=155)	Total (N=261)	p value
- Assistant	69 (65.1%)	43 (27.7%)	112 (42.9%)	
- Associate	21 (19.8%)	43 (27.7%)	64 (24.5%)	
- Full professor	16 (15.1%)	69 (44.5%)	85 (32.6%)	
- Missing	0	0	0	
Salary in academic year 1994				< 0.001
- Mean (SD)	118871.274 (56168.006)	177338.761 (85930.540)	153593.345 (80469.667)	
- Median (IQR)	108457.000 (75774.500, 143096.000)	155006.000 (109687.000, 231501.500)	133284.000 (90771.000, 200543.000)	
- Min - Max	34514.000 - 308081.000	52582.000 - 428876.000	34514.000 - 428876.000	
- Missing	0	0	0	
Salary after increment to Sal94				< 0.001
- Mean (SD)	130876.915 (62034.507)	194914.090 (94902.728)	168906.655 (88778.425)	
- Median (IQR)	119135.000 (82345.250, 154170.500)	170967.000 (119952.500, 257163.000)	148117.000 (99972.000, 218955.000)	
- Min - Max	38675.000 - 339664.000	58923.000 - 472589.000	38675.000 - 472589.000	
- Missing	0	0	0	

## References

## Appendix