

Project1

2023-11-28

R markdown project

This project analyses data from Wooldridge Source: M. Blackburn and D. Neumark (1992), “Unobserved Ability, Efficiency Wages, and Interindustry Wage Differentials,” Quarterly Journal of Economics 107, 1421-1436.

In this exploration, we aim to unpack the hypothesis suggesting that married women command higher wages. As we navigate the intricacies of a linear regression model, our focus extends beyond scrutinizing the “married” variable alone. We delve into not only the coefficient and significance of marital status but also other pertinent correlations with wage, presenting a holistic perspective on the myriad factors influencing earnings.

A data.frame with 935 observations on 17 variables:

wage: monthly earnings

hours: average weekly hours

IQ: IQ score

KWW: knowledge of world work score

educ: years of education

exper: years of work experience

tenure: years with current employer

age: age in years

married: =1 if married

black: =1 if black

south: =1 if live in south

urban: =1 if live in SMSA (Standard Metropolitan Statistical Area)

sibs: number of siblings

brthord: birth order

meduc: mother's education

feduc: father's education

lwage: natural log of wage

```
library(wooldridge)
library(Hmisc)

data('wage2')
# Removing missing values, total 663 observations
```

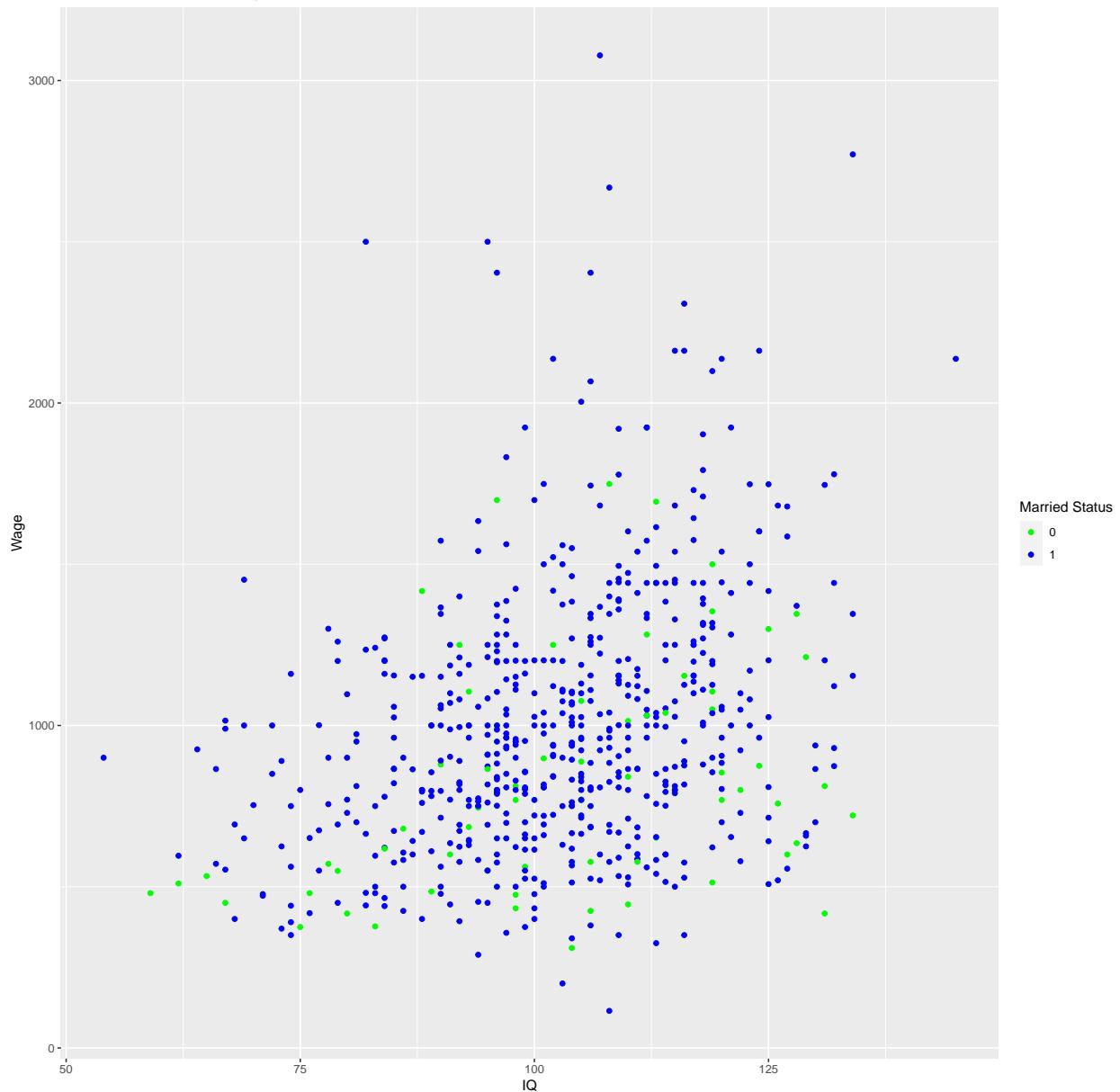
```
wage2=na.omit(wage2)
#describe(wage2)
head(wage2)
```

```
##   wage hours  IQ KWW educ exper tenure age married black south urban sibs
## 1 769    40  93  35   12   11      2  31      1     0     0     1     1
## 3 825    40 108  46   14   11      9  33      1     0     0     1     1
## 4 650    40  96  32   12   13      7  32      1     0     0     1     4
## 5 562    40  74  27   11   14      5  34      1     0     0     1    10
## 7 600    40  91  24   10   13      0  30      0     0     0     1     1
## 9 1154   45 111  37   15   13      1  36      1     0     0     0     2
##   brthord meduc feduc     lwage
## 1          2     8     8 6.645091
## 3          2    14    14 6.715384
## 4          3    12    12 6.476973
## 5          6     6    11 6.331502
## 7          2     8     8 6.396930
## 9          3    14    5 7.050990
```

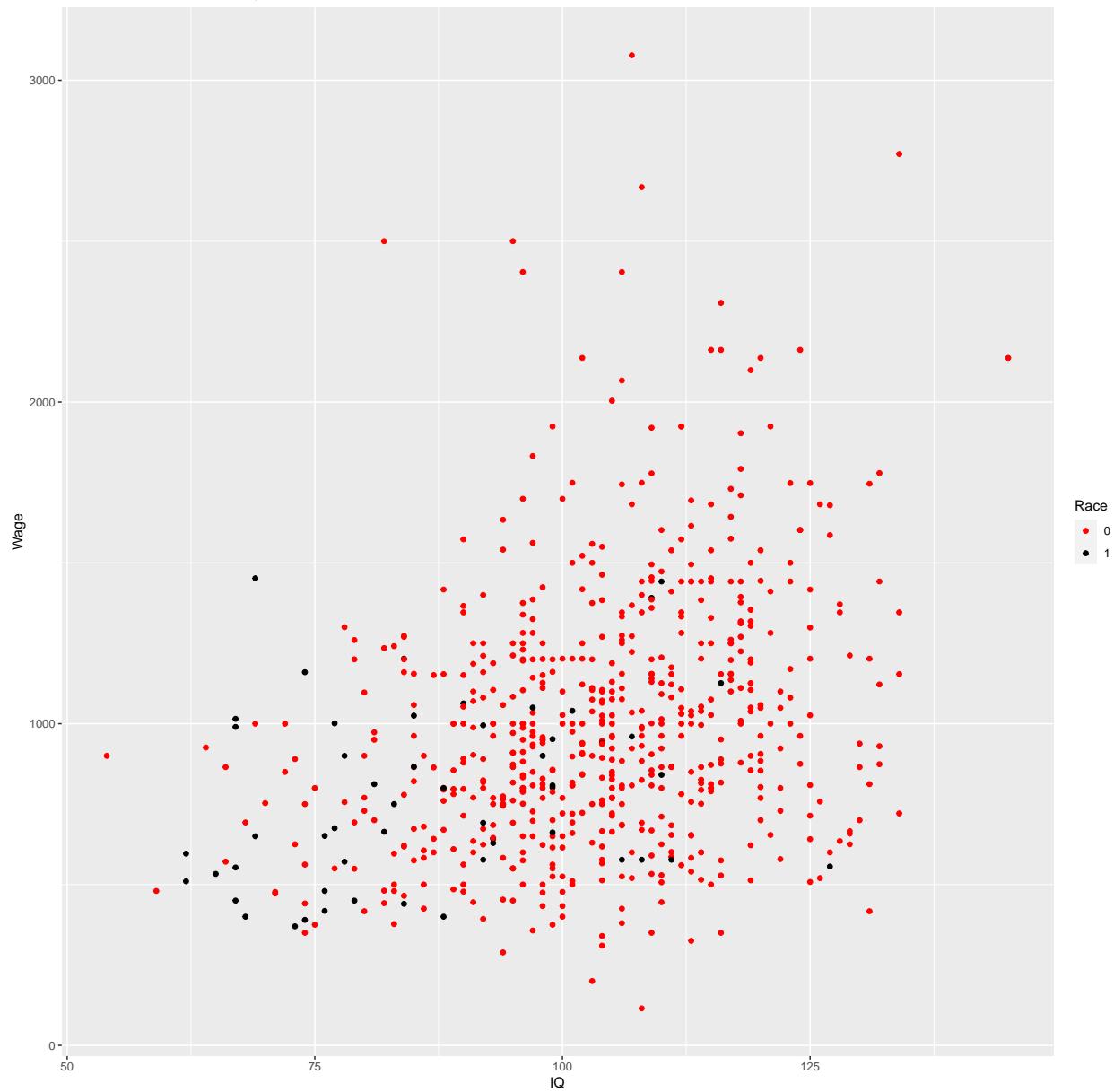
Data visualisation

Creating a scatter matrix plot of data with correlation coefficients. The scatter plot matrix did not highlight best correlations, that is why providing a separate plot for correlations.

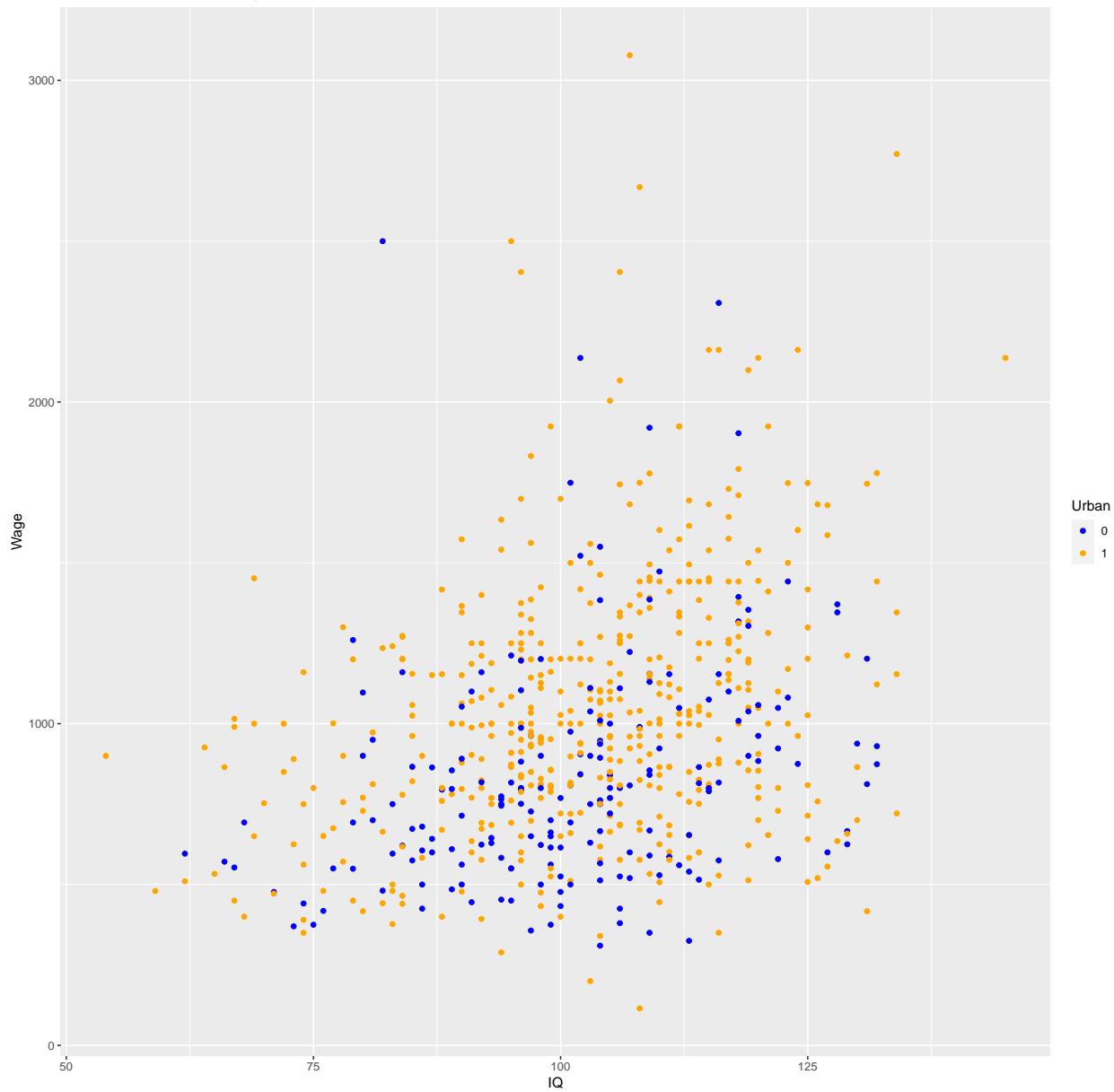
Scatter Plot of IQ and Wage Colored By Married Status



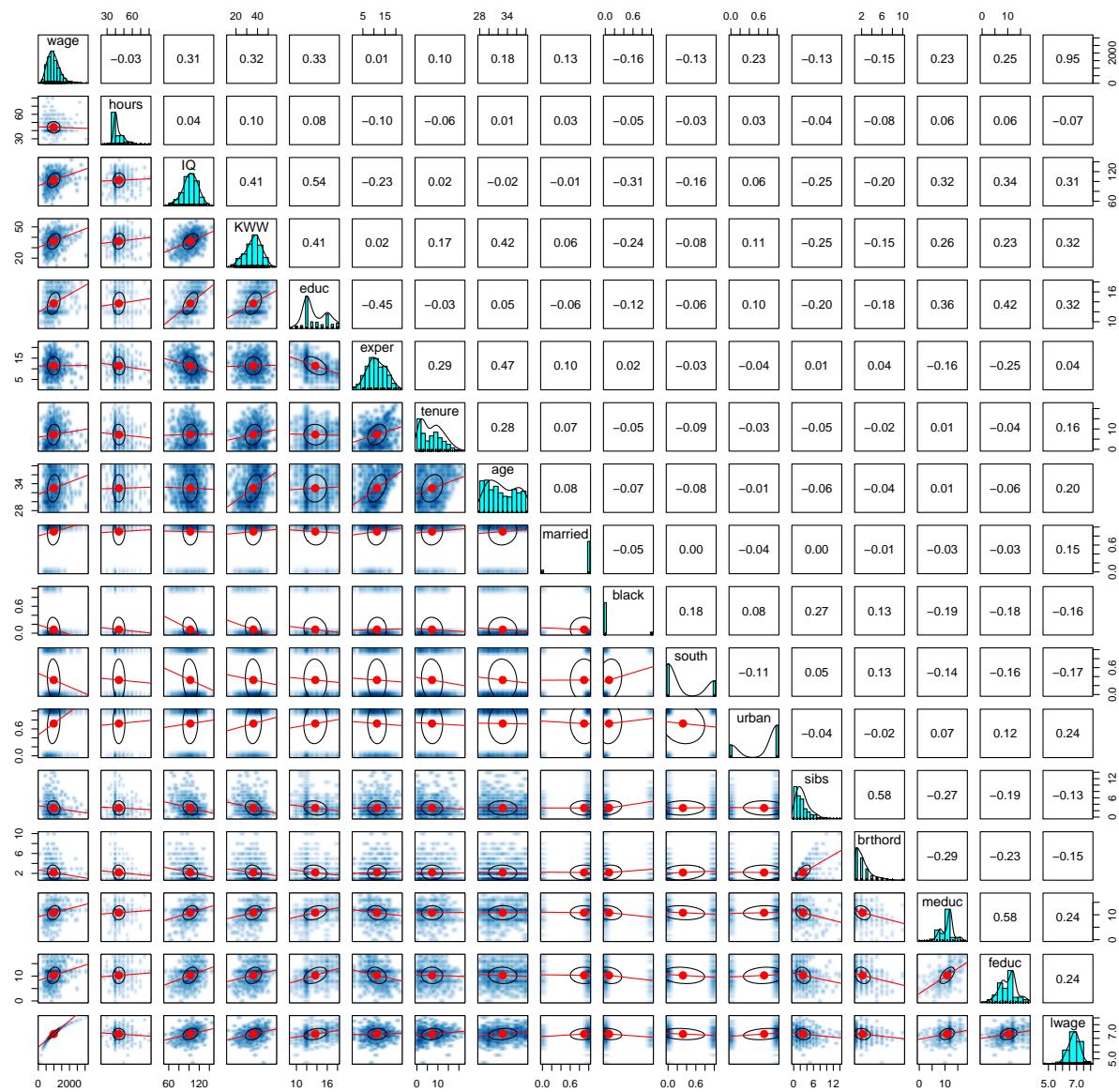
Scatter Plot of IQ and Wage Colored By Race

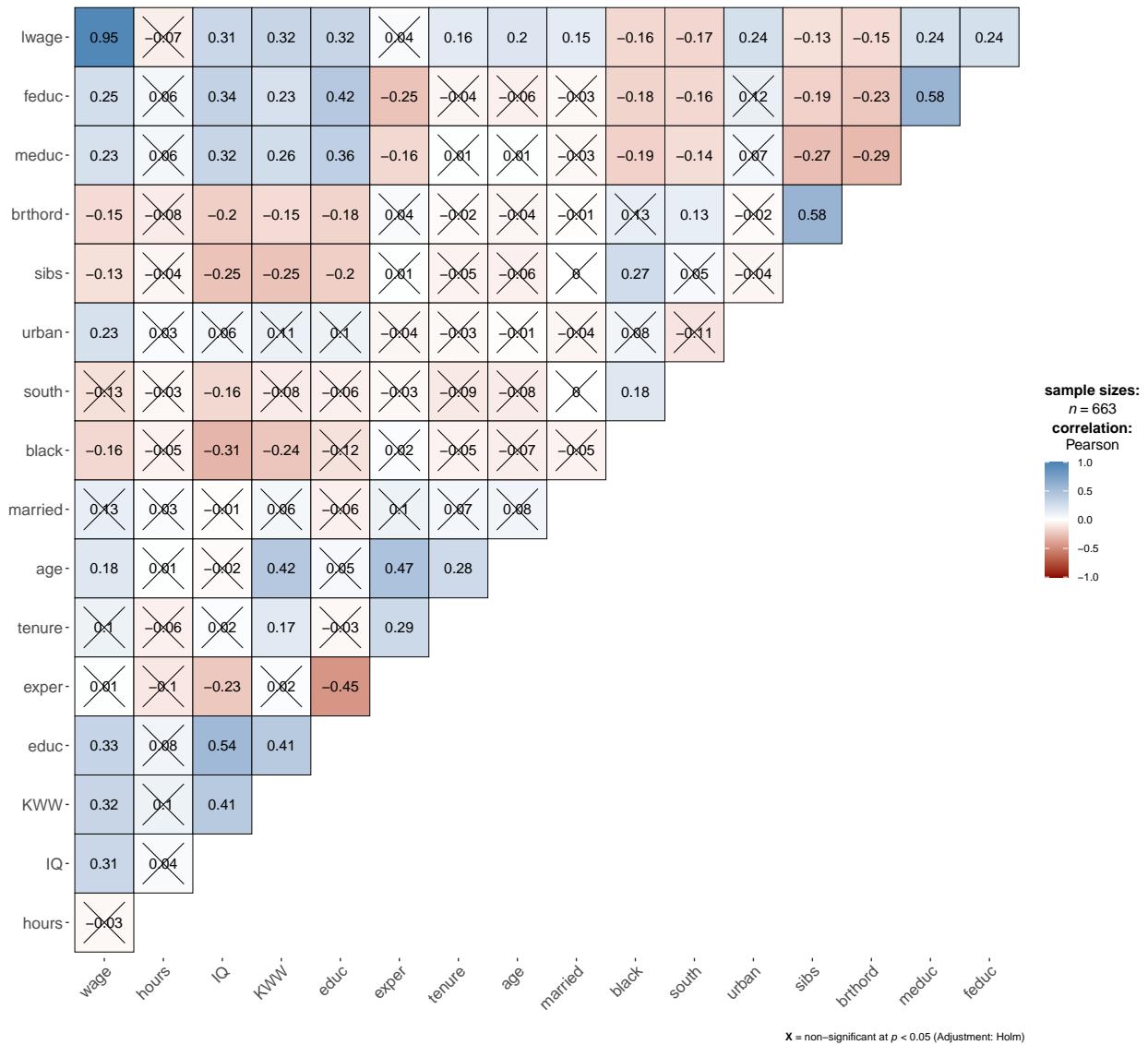


Scatter Plot of IQ and Wage Colored By Urban



Scatter Plot Matrix for Wage2 Dataset





Normality

Checking the normality of the data. This block of code checks for normality of the data. Wage is skewed to left as expected, but the log of wage is more towards normal distribution

```
library('ggpubr')
par(mfrow=c(2,2))
shapiro.test(wage2$wage)

## 
##  Shapiro-Wilk normality test
##
```

```

## data: wage2$wage
## W = 0.93802, p-value = 5.575e-16

shapiro.test(wage2$lwage)

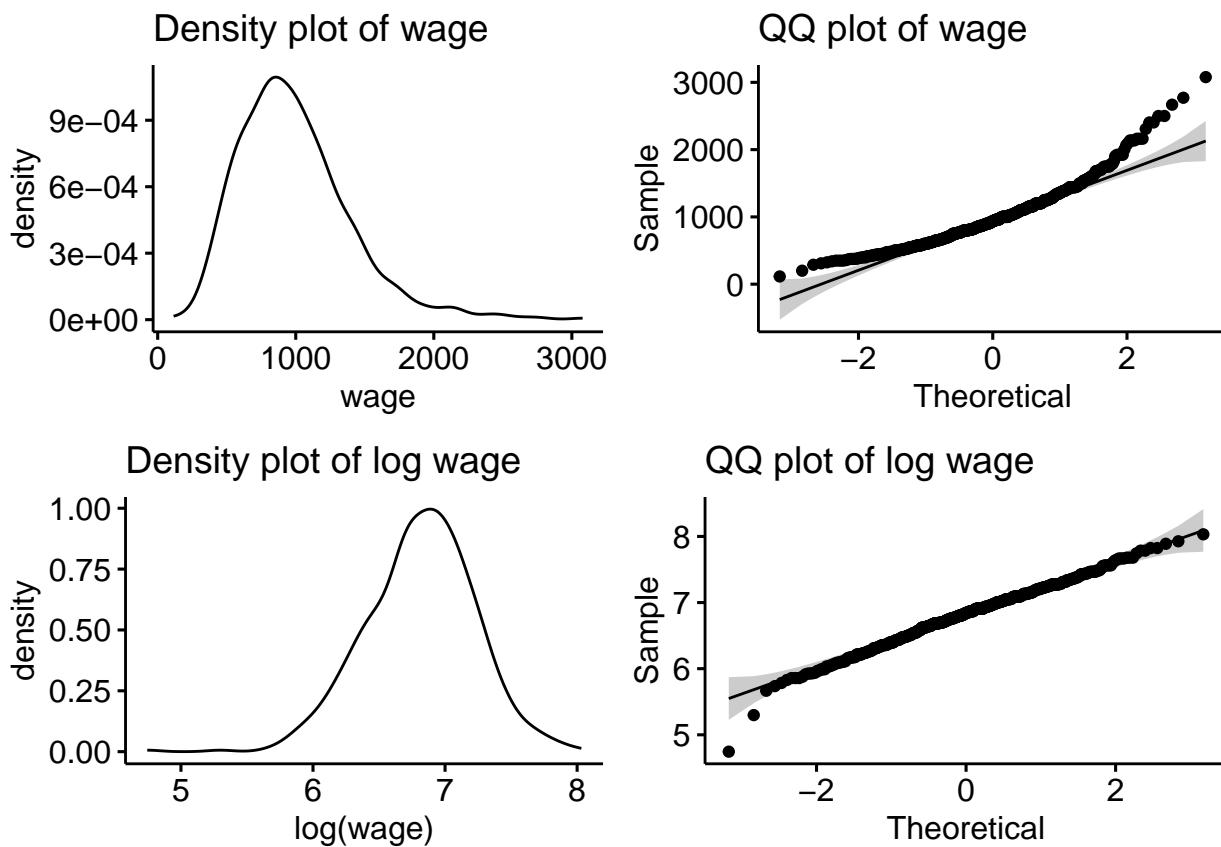
##
## Shapiro-Wilk normality test
##
## data: wage2$lwage
## W = 0.99131, p-value = 0.0006163

plot0 <- ggdensity(wage2$wage, main="Density plot of wage", xlab="wage")
plot1 <- ggqqplot(wage2$wage, title="QQ plot of wage")

plot2 <- ggdensity(wage2$lwage, main="Density plot of log wage", xlab="log(wage)")
plot3 <- ggqqplot(wage2$lwage, title="QQ plot of log wage")

ggarrange(plot0, plot1, plot2, plot3, ncol=2, nrow=2)

```



```

## Linear regression

# Remove the wage from the model
wage2 <- subset(wage2, select = -wage)

# Standard model

```

```

model0 <- lm(wage2$lwage~, data = wage2)

library(MASS)
# step select smaller model
fit <- lm(wage2$lwage~, data = wage2)
model1 <- stepAIC(fit, direction='both', trace=0)

summary(model0)

##
## Call:
## lm(formula = wage2$lwage ~ ., data = wage2)
##
## Residuals:
##      Min       1Q   Median       3Q      Max 
## -1.96887 -0.19460  0.00923  0.22401  1.34185 
##
## Coefficients:
##             Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 5.156439  0.225286 22.888 < 2e-16 ***
## hours      -0.006548  0.001934 -3.385 0.000754 *** 
## IQ          0.003186  0.001223  2.604 0.009425 **  
## KWW         0.003735  0.002390  1.562 0.118662    
## educ        0.041267  0.008942  4.615 4.74e-06 ***
## exper       0.010749  0.004435  2.424 0.015629 *   
## tenure      0.007102  0.002894  2.454 0.014401 *  
## age          0.009107  0.005977  1.524 0.128058    
## married     0.200760  0.045998  4.365 1.48e-05 *** 
## black       -0.105141  0.055667 -1.889 0.059373 .  
## south       -0.049076  0.030753 -1.596 0.111019    
## urban        0.195658  0.031240  6.263 6.88e-10 *** 
## sibs         0.009619  0.007876  1.221 0.222423    
## brthord     -0.018465  0.011569 -1.596 0.110975    
## meduc        0.009633  0.006167  1.562 0.118753    
## feduc        0.005590  0.005398  1.036 0.300804    
## ---        
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.3507 on 647 degrees of freedom
## Multiple R-squared:  0.2925, Adjusted R-squared:  0.2761 
## F-statistic: 17.84 on 15 and 647 DF,  p-value: < 2.2e-16

summary(model1)

```

```

##
## Call:
## lm(formula = wage2$lwage ~ hours + IQ + KWW + educ + exper +
##     tenure + age + married + black + south + urban + meduc, data = wage2)
##
## Residuals:
##      Min       1Q   Median       3Q      Max 
## -2.01294 -0.20059  0.01195  0.22524  1.32195 

```

```

## 
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 5.134460  0.217379 23.620 < 2e-16 ***
## hours      -0.006371  0.001932 -3.298 0.00103 **  
## IQ          0.003273  0.001221  2.681 0.00753 **  
## KWW         0.003546  0.002377  1.492 0.13612    
## educ        0.042825  0.008799  4.867 1.42e-06 ***
## exper       0.010333  0.004424  2.336 0.01981 *   
## tenure      0.006943  0.002895  2.398 0.01676 *   
## age          0.009068  0.005958  1.522 0.12849    
## married     0.203743  0.046006  4.429 1.11e-05 *** 
## black        -0.096556  0.054384 -1.775 0.07629 .  
## south        -0.057885  0.030457 -1.901 0.05781 .  
## urban        0.196127  0.031141  6.298 5.55e-10 *** 
## meduc        0.013480  0.005348  2.521 0.01195 *  
## --- 
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 0.351 on 650 degrees of freedom
## Multiple R-squared:  0.2882, Adjusted R-squared:  0.275 
## F-statistic: 21.93 on 12 and 650 DF,  p-value: < 2.2e-16

```

Collinearity

if higher than 5 indicates collinearity. Here we do not have this issue, but if we did, we would remove one of the variables with high VIF from the model.

```

library(car)
vif(model0)

##    hours      IQ      KWW      educ      exper      tenure      age      married
## 1.032256 1.737404 1.743251 2.142799 1.919560 1.152359 1.803937 1.022335
##    black      south      urban      sibs      brthord      meduc      feduc
## 1.249702 1.114358 1.061829 1.676644 1.594288 1.631398 1.695625

vif(model1)

##    hours      IQ      KWW      educ      exper      tenure      age      married
## 1.027821 1.727252 1.720755 2.071801 1.907037 1.151330 1.789792 1.021123
##    black      south      urban      meduc
## 1.190886 1.091371 1.053447 1.224933

```

Heteroscedasticity

```

bptest(model1)

## 
## studentized Breusch-Pagan test
## 
## data: model1
## BP = 37.256, df = 12, p-value = 0.000203

```

```

model1.fixed <- coeftest(model1, vcov. = vcovHC(model1, type = 'HC3'))
model1.fixed

## 
## t test of coefficients:
## 
##           Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 5.1344598  0.2194234 23.3998 < 2.2e-16 ***
## hours       -0.0063708  0.0023217 -2.7440  0.006237 **  
## IQ          0.0032725  0.0011007  2.9730  0.003058 **  
## KWW         0.0035465  0.0027280  1.3001  0.194043  
## educ        0.0428251  0.0094180  4.5472 6.487e-06 *** 
## exper       0.0103329  0.0047388  2.1805  0.029578 *   
## tenure      0.0069429  0.0030515  2.2752  0.023217 *  
## age         0.0090684  0.0062250  1.4568  0.145663  
## married     0.2037429  0.0480420  4.2409 2.549e-05 *** 
## black       -0.0965560  0.0504821 -1.9127  0.056229 .  
## south        -0.0578847  0.0329506 -1.7567  0.079438 .  
## urban        0.1961270  0.0311978  6.2866 5.953e-10 *** 
## meduc       0.0134802  0.0049654  2.7148  0.006807 **  
## ---        
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

bpptest(model0)

## 
## studentized Breusch-Pagan test
## 
## data: model0
## BP = 38.914, df = 15, p-value = 0.00066

model0.fixed <- model1.fixed <- coeftest(model1, vcov. = vcovHC(model1, type = 'HC3'))
model0.fixed

## 
## t test of coefficients:
## 
##           Estimate Std. Error t value Pr(>|t|)    
## (Intercept) 5.1344598  0.2194234 23.3998 < 2.2e-16 ***
## hours       -0.0063708  0.0023217 -2.7440  0.006237 **  
## IQ          0.0032725  0.0011007  2.9730  0.003058 **  
## KWW         0.0035465  0.0027280  1.3001  0.194043  
## educ        0.0428251  0.0094180  4.5472 6.487e-06 *** 
## exper       0.0103329  0.0047388  2.1805  0.029578 *   
## tenure      0.0069429  0.0030515  2.2752  0.023217 *  
## age         0.0090684  0.0062250  1.4568  0.145663  
## married     0.2037429  0.0480420  4.2409 2.549e-05 *** 
## black       -0.0965560  0.0504821 -1.9127  0.056229 .  
## south        -0.0578847  0.0329506 -1.7567  0.079438 .  
## urban        0.1961270  0.0311978  6.2866 5.953e-10 *** 
## meduc       0.0134802  0.0049654  2.7148  0.006807 **  
## ---        
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

Normality of residuals

```
library(olsrr)

## 
## Attaching package: 'olsrr'

## The following object is masked from 'package:MASS':
## 
##     cement

## The following object is masked from 'package:wooldridge':
## 
##     cement

## The following object is masked from 'package:datasets':
## 
##     rivers

ols_test_normality(model0)

## -----
##      Test       Statistic      pvalue
## -----
## Shapiro-Wilk      0.9804      0.0000
## Kolmogorov-Smirnov 0.0383      0.2838
## Cramer-von Mises   110.23      0.0000
## Anderson-Darling    1.5062      7e-04
## -----


ols_test_normality(model1)

## -----
##      Test       Statistic      pvalue
## -----
## Shapiro-Wilk      0.9804      0.0000
## Kolmogorov-Smirnov 0.0358      0.3624
## Cramer-von Mises   109.7335     0.0000
## Anderson-Darling    1.3792      0.0014
## -----
```

Conclusion

The coefficient for “married” stands at 0.0357696, suggesting a positive association between marital status and wages. Crucially, the p-value of 0.023591, less than the conventional significance level of 0.05, signifies statistical significance. This implies that, according to the model, being married is not merely incidental but significantly associated with higher earnings. Beyond the exploration of marital status, other variables like wage levels, working hours, tenure, and regional dynamics emerge as influential contributors to the earnings tapestry.

In conclusion, the analysis lends support to the hypothesis that married women do, indeed, earn higher wages. The positive coefficient and statistically significant p-value imply a meaningful connection between marital status and increased earnings, providing empirical evidence to corroborate the initial claim. It's imperative to note that this statistical evidence, while robust, does not establish causation, and other unconsidered factors may contribute to the observed relationship.