# Using Interactive Jupyter Notebooks with R

Earl F Glynn

Kansas City R Users Group

2015-12-05

http://earlglynn.github.io/kc-r-users-jupyter/

# Using Interactive Jupyter Notebooks with R

- What is Jupyter?
- R User Interface Evolution
  - Command Line
  - RStudio
  - RStudio with Markdown
  - Jupyter Notebook
- Jupyter Markdown Cells
- Jupyter Code Cells
- Installation of Jupyter

# What is Jupyter?

- http://jupyter.org/
- Language-agnostic parts of IPython ("Interactive Python") http://ipython.org/
- Provides interactive data science and scientific computing across ~40 programming languages
- **Ju**lia – **Pyt**hon – **R**

# R User Interface Evolution

- R Command Line
- RStudio
- RStudio with Markdown
- Jupyter Notebook

Comparisons using `?lm` help example

# R Command Line

?lm

```
## Annette Dobson (1990) "An Introduction to Generalized Linear Models".
## Page 9: Plant Weight Data.
ctl <- c(4.17,5.58,5.18,6.11,4.50,4.61,5.17,4.53,5.33,5.14)
trt <- c(4.81,4.17,4.41,3.59,5.87,3.83,6.03,4.89,4.32,4.69)
group <- gl(2, 10, 20, labels = c("Ctl","Trt"))
weight <- c(ctl, trt)
lm.D9 <- lm(weight ~ group)
lm.D90 <- lm(weight ~ group - 1) # omitting intercept

anova(lm.D9)
summary(lm.D90)

opar <- par(mfrow = c(2,2), oma = c(0, 0, 1.1, 0))
plot(lm.D9, las = 1)        # Residuals, Fitted, ...
par(opar)
```

Copy and paste to R console window

# R Command Line

```
> ## Annette Dobson (1990) "An Introduction to Generalized Linear Models".
> ## Page 9: Plant Weight Data.
> ctl <- c(4.17,5.58,5.18,6.11,4.50,4.61,5.17,4.53,5.33,5.14)
> trt <- c(4.81,4.17,4.41,3.59,5.87,3.83,6.03,4.89,4.32,4.69)
> group <- gl(2, 10, 20, labels = c("Ctl","Trt"))
> weight <- c(ctl, trt)
> lm.D9 <- lm(weight ~ group)
> lm.D90 <- lm(weight ~ group - 1) # omitting intercept
>
> anova(lm.D9)
Analysis of Variance Table

Response: weight
          Df Sum Sq Mean Sq F value Pr(>F)
group      1 0.6882 0.68820  1.4191  0.249
Residuals 18 8.7292 0.48496
> summary(lm.D90)

Call:
lm(formula = weight ~ group - 1)

Residuals:
    Min      1Q  Median      3Q     Max
-1.0710 -0.4938  0.0685  0.2462  1.3690

Coefficients:
         Estimate Std. Error t value Pr(>|t|)
groupCtl   5.0320     0.2202   22.85 9.55e-15 ***
groupTrt   4.6610     0.2202   21.16 3.62e-14 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.6964 on 18 degrees of freedom
Multiple R-squared:  0.9818,    Adjusted R-squared:  0.9798
F-statistic: 485.1 on 2 and 18 DF,  p-value: < 2.2e-16

>
> opar <- par(mfrow = c(2,2), oma = c(0, 0, 1.1, 0))
> plot(lm.D9, las = 1)        # Residuals, Fitted, ...
> par(opar)
```
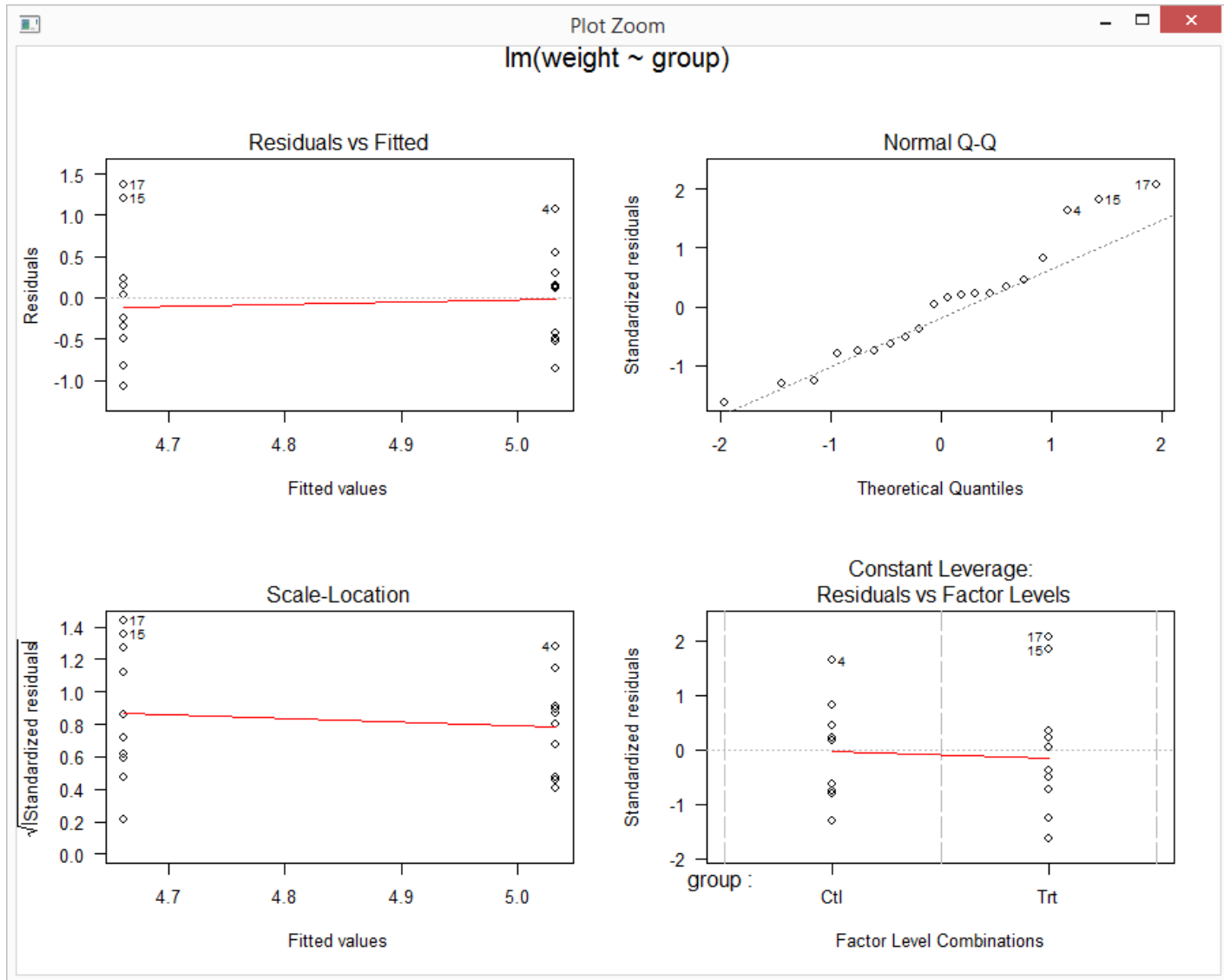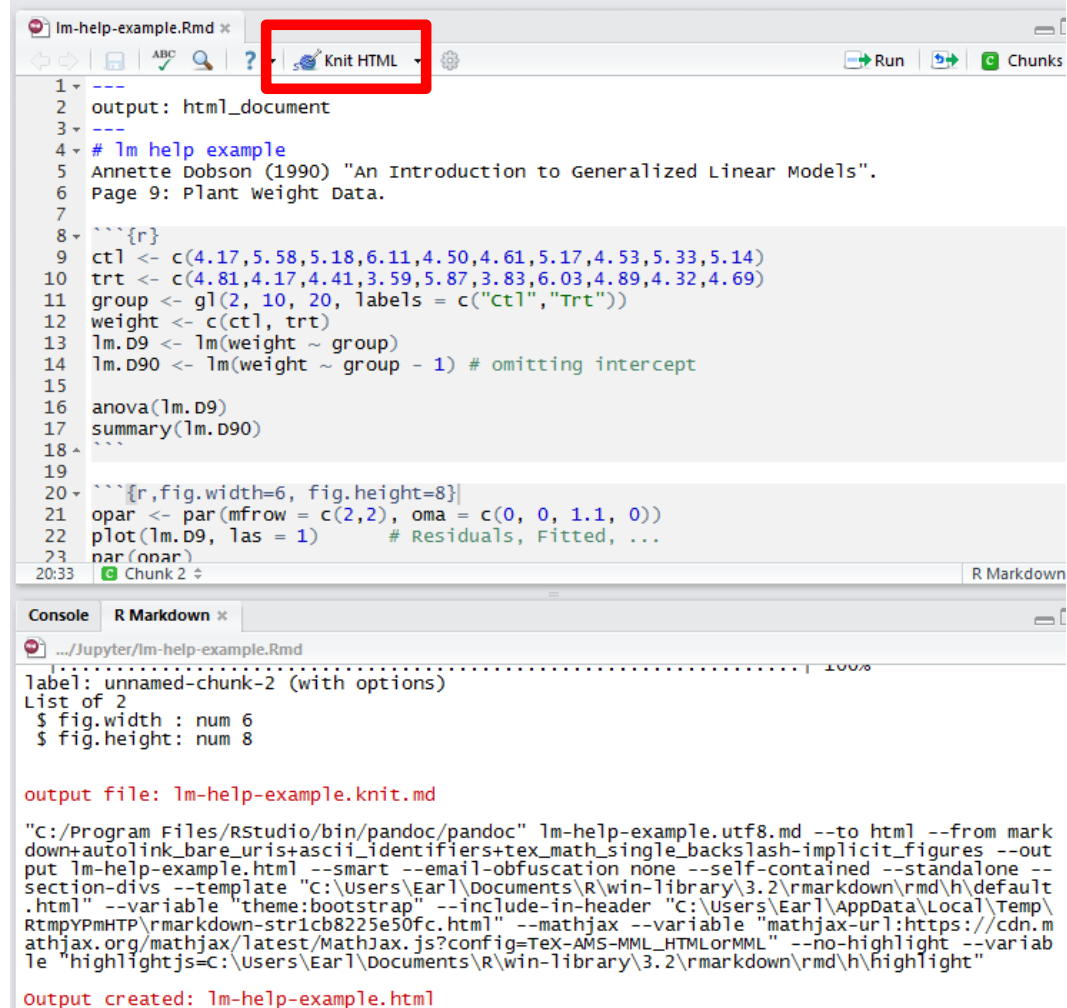
# RStudio



https://www.rstudio.com/products/RStudio/

# RStudio

# RStudio

# RStudio

# RStudio with Markdown



Markdown Basics:  http://rmarkdown.rstudio.com/authoring_basics.html

# RStudio with Markdown

Output to HTML, PDF, Word.
Graphics output included.

lm-help-example.html | Open in Browser | Find |

## lm help example

Annette Dobson (1990) "An Introduction to Generalized Linear Models". Page 9: Plant Weight Data.

```
ctl <- c(4.17,5.58,5.18,6.11,4.50,4.61,5.17,4.53,5.33,5.14)
trt <- c(4.81,4.17,4.41,3.59,5.87,3.83,6.03,4.89,4.32,4.69)
group <- gl(2, 10, 20, labels = c("Ctl","Trt"))
weight <- c(ctl, trt)
lm.D9 <- lm(weight ~ group)
lm.D90 <- lm(weight ~ group - 1) # omitting intercept

anova(lm.D9)
```
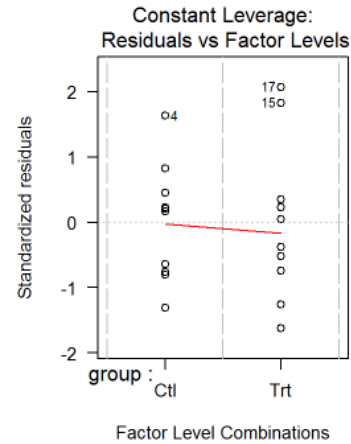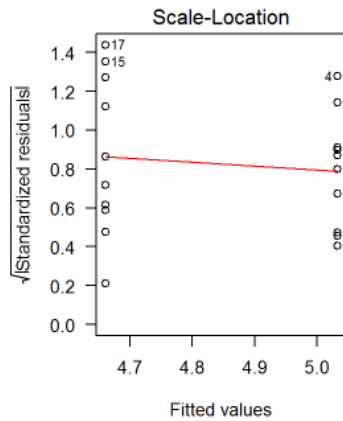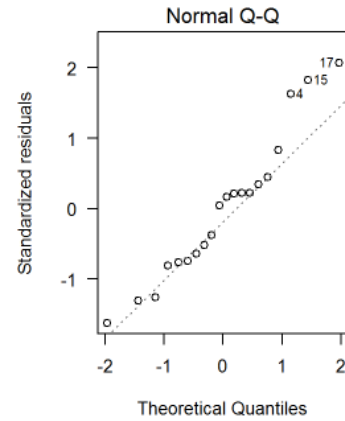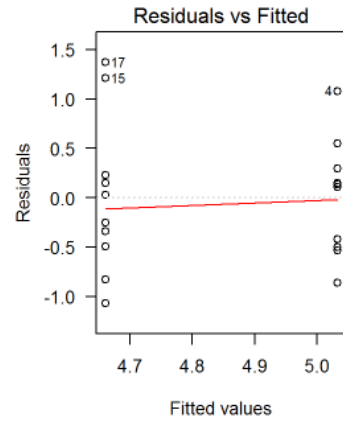
```
## Analysis of Variance Table
##
## Response: weight
##           Df Sum Sq Mean Sq F value Pr(>F)
## group      1 0.6882 0.68820  1.4191  0.249
## Residuals 18 8.7292 0.48496
```
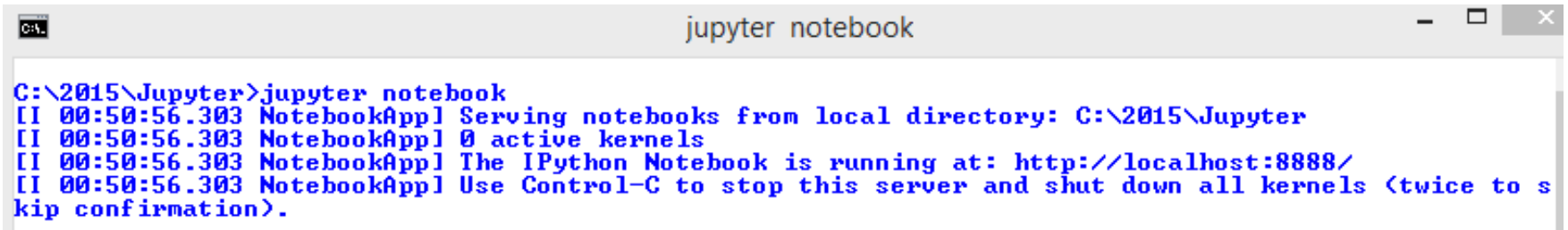
```
summary(lm.D90)
```

```
##
## C-11
```

# RStudio with Markdown

# Jupyter Notebook

From command window in working directory, start Jupyter notebook server:

```
jupyter notebook
```

# Jupyter Notebook

# Jupyter Notebook



Add Cell

Each Jupyter cell contains Markdown or the equivalent of a Code "chunk" in RStudio

# Jupyter Notebook



Markdown →

Annette Dobson (1990) "An Introduction to Generalized Linear Models". Page 9: Plant Weight Data.

Code →

```
In [1]:  ctl <- c(4.17,5.58,5.18,6.11,4.50,4.61,5.17,4.53,5.33,5.14)
         trt <- c(4.81,4.17,4.41,3.59,5.87,3.83,6.03,4.89,4.32,4.69)
         group <- gl(2, 10, 20, labels = c("Ctl","Trt"))
         weight <- c(ctl, trt)
         lm.D9 <- lm(weight ~ group)
         lm.D9
```

```
Out[1]:
         Call:
         lm(formula = weight ~ group)

         Coefficients:
         (Intercept)      groupTrt
               5.032        -0.371
```
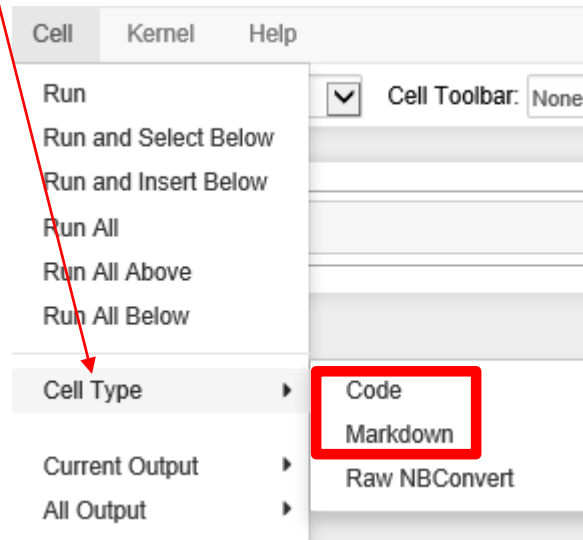
Unlike RStudio/knittr, no special syntax for code chunk.
Enter "Ctrl-Enter" to execute code in cell interactively.
Out[1] is the R output here from cell In[1].

# Jupyter Notebook

```
In [2]: lm.D90 <- lm(weight ~ group - 1) # omitting intercept

anova(lm.D9)
```

Out[2]:

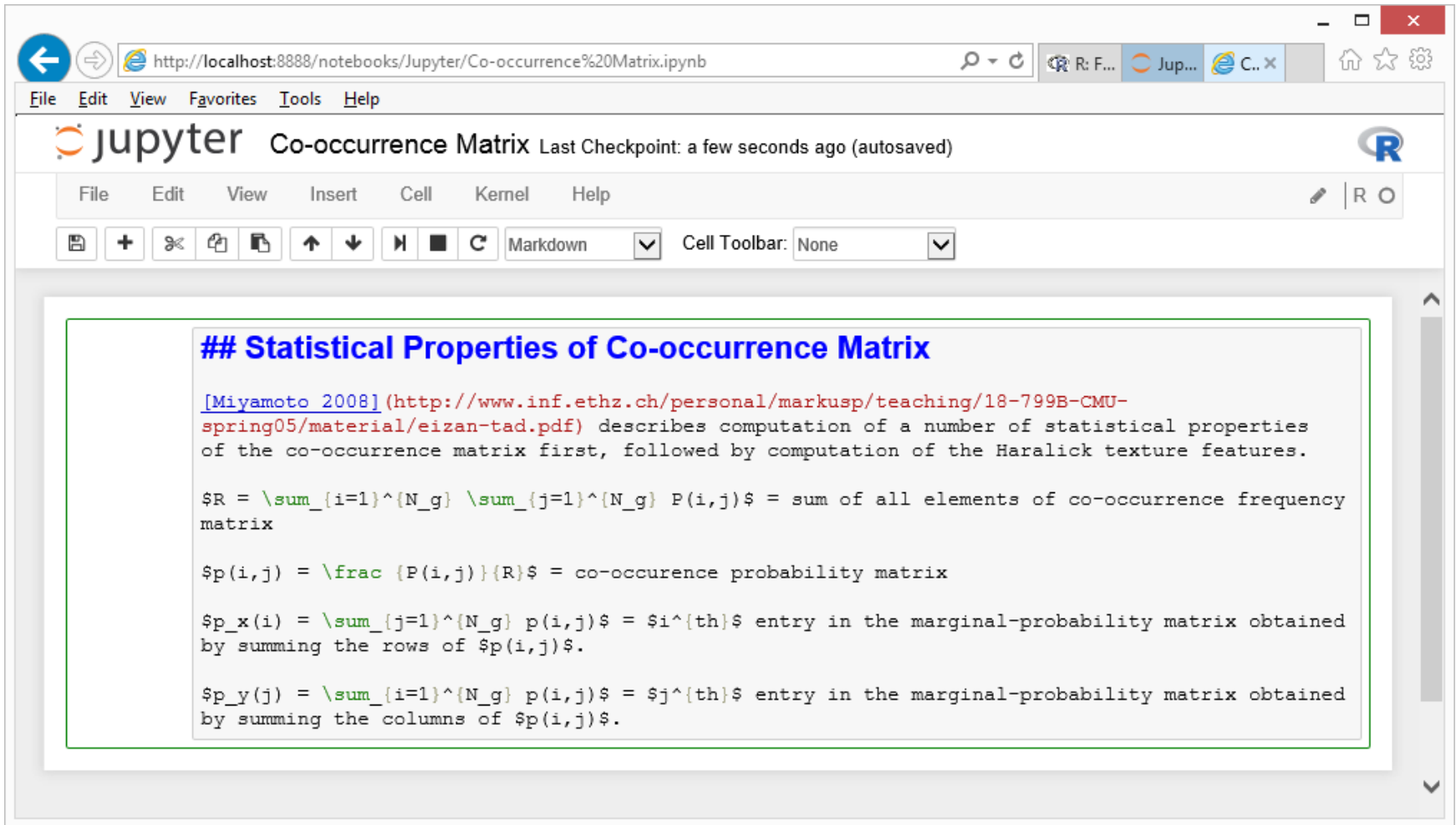|  | Df | Sum Sq | Mean Sq | F value | Pr(>F) |
|---|---|---|---|---|---|
| **group** | 1 | 0.688205 | 0.688205 | 1.419101 | 0.2490232 |
| **Residuals** | 18 | 8.72925 | 0.4849583 | NA | NA |

# Jupyter Notebook

```
In [4]:  options(repr.plot.width=6, repr.plot.height=6)
         opar <- par(mfrow = c(2,2), oma = c(0, 0, 1.1, 0))
         plot(lm.D9, las = 1)        # Residuals, Fitted, ...
         par(opar)
```



lm(weight ~ group)

# Jupyter Markdown Cells



Markdown example including inline LaTeX equations. *Ctrl-Enter* to render.

# Jupyter Markdown Cells

# Jupyter Code Cells

Online Examples:
[http://earlglynn.github.io/kc-r-users-jupyter/](http://earlglynn.github.io/kc-r-users-jupyter/)

- Jupyter First Look

- lm help example
- Co-occurrence Matrix

- Exploring Kaggle Facial Keypoints Detection Data

# Installation of Jupyter

Perhaps easiest:

Install Anaconda Python from Continuum Analytics

https://www.continuum.io/downloads

- Python 3.5, Windows 64-bit graphical installer
- Package List: http://docs.continuum.io/anaconda/pkg-docs
  – Includes:  numpy, scipy, scikit-learn, matplotlib, …

# Installation of Jupyter

From command prompt:

- **Conda:** `conda update conda`

- Jupyter: `conda install jupyter`

- R Essentials:
  `conda install -c r r-essentials`

- R Kernel:
  `conda install -c r ipython-notebook r-irkernel`
  http://irkernel.github.io/installation/
  https://www.continuum.io/blog/developer/jupyter-and-conda-r

# R Packages Used by Jupyter

```
In [1]:  .libPaths()

Out[1]:       "C:/Users/Earl/Documents/R/win-library/3.1"   "C:/Anaconda3/R/library"
```

```
In [2]:  library()
```

```
Packages in library 'C:/Anaconda3/R/library':

base                    The R Base Package
base64enc               Tools for base64 encoding
boot                    Bootstrap Functions (Originally by Angelo Canty
                        for S)
class                   Functions for Classification
cluster                 Cluster Analysis Extended Rousseeuw et al.
codetools               Code Analysis Tools for R
compiler                The R Compiler Package
datasets                The R Datasets Package
```

. . .

# Installation of Jupyter

## Kernels for other languages:

https://github.com/ipython/ipython/wiki/IPython-kernels-for-other-languages

# Take Home Message

Jupyter is a great way to use R interactively to document the steps in a data analysis project.

Jupyter's interactive approach is better (IMHO) than the batch processing by RStudio/knitr to document reproducible results.