

STATS 205: Homework Assignment 4

Brian Liu

6/3/2019

Solution to Problem 1

```
allergics = c(1651.0, 1112.0, 102.4, 100.0, 67.6, 65.9, 64.7, 39.6, 31.0)
nonallergics = c(48.1, 48.0, 45.5, 41.7, 35.4, 34.3, 32.4, 29.1, 27.3, 18.9, 6.6, 5.2, 4.7)
allergics; nonallergics
```

```
## [1] 1651.0 1112.0 102.4 100.0 67.6 65.9 64.7 39.6 31.0
```

```
## [1] 48.1 48.0 45.5 41.7 35.4 34.3 32.4 29.1 27.3 18.9 6.6 5.2 4.7
```

The null hypothesis is that allergic smokers have the same sputum histamine levels as nonallergic smokers. That is,

$$H_0 : p_a = p_n$$

The alternative hypothesis is that allergic smokers have higher sputum histamine levels than nonallergic smokers. That is,

$$H_0 : p_a > p_n$$

To test the null hypothesis against the alternative hypothesis, we will use the Mann-Whitney-Wilcoxin test, since the two samples are independent.

Two data samples are independent if they come from distinct populations and the samples do not affect each other.

– Mann-Whitney-Wilcoxin Test

```
wilcox.test(x = allergics, y = nonallergics, alternative = "greater")
```

```
##
```

```
## Wilcoxon rank sum test
```

```
##
```

```
## data: allergics and nonallergics
```

```
## W = 106, p-value = 0.000386
```

```
## alternative hypothesis: true location shift is greater than 0
```

The p -value is 0.000386, which is significant at the $\alpha = 0.05$ level. There is strong evidence that allergic smokers have higher sputum histamine levels than nonallergic smokers.

Solution to Problem 2

Original problem statement

```
karate = c(37, 39, 30, 7, 13, 139, 45, 25, 16, 146, 94, 16, 23, 1, 290, 169, 62, 145, 36, 20, 13)
olympics = c(12, 44, 34, 14, 9, 19, 156, 23, 13, 11, 47, 26, 14, 33, 15, 62, 5, 8, 0, 154, 146)
karate; olympics
```

```
## [1] 37 39 30 7 13 139 45 25 16 146 94 16 23 1 290 169 62
## [18] 145 36 20 13

## [1] 12 44 34 14 9 19 156 23 13 11 47 26 14 33 15 62 5
## [18] 8 0 154 146
```

The null hypothesis is that children who viewed the violent TV take the same amount of time to seek help (were as tolerant) as the children who viewed the nonviolent sports-action TV. That is,

$$H_0 : t_k = t_o$$

The alternative hypothesis is that children who viewed the violent TV take longer to seek help (were more tolerant) than the children who viewed the nonviolent sports-action TV. That is,

$$H_0 : t_k > t_o$$

```
wilcox.test(x = karate, y = olympics, alternative = "greater")
```

```
## Warning in wilcox.test.default(x = karate, y = olympics, alternative =
## "greater"): cannot compute exact p-value with ties

##
## Wilcoxon rank sum test with continuity correction
##
## data: karate and olympics
## W = 276.5, p-value = 0.08126
## alternative hypothesis: true location shift is greater than 0
```

The p -value is 0.08126, which is *not* significant at the $\alpha = 0.05$ level. There is *not enough* evidence that children who viewed the violent TV take longer to seek help (were more tolerant) than the children who viewed the nonviolent sports-action TV.

Solution to Problem 3

Let X be the nonallergics and Y be the allergics.

$$\delta = P(X < Y)$$

```
allergics = c(1651.0, 1112.0, 102.4, 100.0, 67.6, 65.9, 64.7, 39.6, 31.0)
nonallergics = c(48.1, 48.0, 45.5, 41.7, 35.4, 34.3, 32.4, 29.1, 27.3, 18.9, 6.6, 5.2, 4.7)
allergics; nonallergics
```

```
## [1] 1651.0 1112.0 102.4 100.0 67.6 65.9 64.7 39.6 31.0
## [1] 48.1 48.0 45.5 41.7 35.4 34.3 32.4 29.1 27.3 18.9 6.6 5.2 4.7
wilcox.test(x = allergics, y = nonallergics, conf.int=TRUE, conf.level=.90)
```

```
##
## Wilcoxon rank sum test
##
## data: allergics and nonallergics
## W = 106, p-value = 0.000772
## alternative hypothesis: true location shift is not equal to 0
## 90 percent confidence interval:
```

```
## 25.9 81.1
## sample estimates:
## difference in location
## 54.3
```

The estimate for δ is

$$\hat{\delta} = P(X < Y) = 54.3$$

and the 90 confidence interval for δ is

$$\hat{\delta} = P(X < Y) = (25.9, 81.1)$$

Solution to Problem 4

```
term = c(0.80, 0.83, 1.89, 1.04, 1.45, 1.38, 1.91, 1.64, 0.73, 1.46)
gest = c(1.15, 0.88, 0.90, 0.74, 1.21)

wilcox.test(x = term, y = gest, alternative = "greater")
```

```
##
## Wilcoxon rank sum test
##
## data: term and gest
## W = 35, p-value = 0.1272
## alternative hypothesis: true location shift is greater than 0
# ks.test(x = term, y = gest, alternative = "greater")
# ks.test(x = term, y = gest)
```

The p -value for the Wilcoxon ranked test is 0.1272.

The p -value for the one-sided Two-sample Kolmogorov-Smirnov test is 0.9355, which is larger than the p -value for the Wilcoxon-ranked test.

The p -value for the two-sided Two-sample Kolmogorov-Smirnov test is 0.1658, which is similar to that of the Wilcoxon-ranked test.

```
# install.packages("npsm", dependencies = TRUE, repos = "http://cran.us.r-project.org")
# install.packages("RVAideMemoire", dependencies = TRUE, repos = "http://cran.us.r-project.org")
# install.packages("robustrank", dependencies = TRUE, repos = "http://cran.us.r-project.org")
# fp.test(x = term, y = gest, alternative = 'two.sided')
# fp.test(x = term, y = gest)
```

Unfortunately, the package for `fp.test()` doesn't seem to be within reach, despite installing the following packages:

- `npsm` package
- `RVAideMemoire` package, with `fp.test()` demonstrated here
- R Documentation of `npsm` package including `fp.test()`
- `robustrank` package

Solution to Problem 5

Null hypothesis: “equal dispersions”

$$H_0 : p_{term} = p_{gest}$$

Alternative hypothesis: “the variation in tritiated water diffusion across human chorioamnion is different at term than at 12–26 weeks gestational age”

$$H_A : p_{term} \neq p_{gest}$$

```
term = c(0.80, 0.83, 1.89, 1.04, 1.45, 1.38, 1.91, 1.64, 0.73, 1.46)
gest = c(1.15, 0.88, 0.90, 0.74, 1.21)

ansari.test(x = term, y = gest, alternative = "t")
```

```
##
##  Ansari-Bradley test
##
## data:  term and gest
## AB = 36, p-value = 0.1372
## alternative hypothesis: true ratio of scales is not equal to 1
```

The p -value is 0.1372, which is *not* significant at the $\alpha = 0.05$ level. There is *not enough* evidence that the variation in tritiated water diffusion across human chorioamnion is different at term than at 12–26 weeks gestational age.

Solution to Problem 6

```
a = c(3.6, 2.6, 4.7, 8.0, 3.1, 8.8, 4.6, 5.8, 4.0, 4.6)
b = c(16.2, 17.4, 8.5, 15.6, 5.4, 9.8, 14.9, 16.6, 15.9, 5.3, 10.5)
```

The null hypothesis is that Type A subjects have the same Peak Levels of Human Plasma Growth Hormone after Arginine Hydrochloride Infusion as Type B subjects. That is,

$$H_0 : l_a = l_b$$

The alternative hypothesis is that Type A subjects have different Peak Levels of Human Plasma Growth Hormone after Arginine Hydrochloride Infusion as Type B subjects. That is,

$$H_0 : l_a > l_b$$

To test the null hypothesis against the alternative hypothesis, we will use the Mann-Whitney-Wilcoxin test, since the two samples are independent.

```
wilcox.test(x = a, y = b, alternative = "two.sided")

## Warning in wilcox.test.default(x = a, y = b, alternative = "two.sided"):
## cannot compute exact p-value with ties
##
##  Wilcoxon rank sum test with continuity correction
##
## data:  a and b
## W = 7, p-value = 0.0008201
## alternative hypothesis: true location shift is not equal to 0
```

The p -value is 0.0008201, which is significant at the $\alpha = 0.05$ level. There is *strong evidence* that Type A subjects have different Peak Levels of Human Plasma Growth Hormone after Arginine Hydrochloride Infusion as Type B subjects.

Solution to Problem 7

```
Darwin.data = data.frame(pair = seq(1, 15), pot = c(rep(1, times=3), rep(2, times = 3), rep(3, times = 3)),
saveRDS(Darwin.data, "Darwin_data.rds"); Darwin.data
```

| ## | pair | pot | cross.height | self.height |
|-------|------|-----|--------------|-------------|
| ## 1 | 1 | 1 | 23.500 | 17.375 |
| ## 2 | 2 | 1 | 12.000 | 20.375 |
| ## 3 | 3 | 1 | 21.000 | 20.000 |
| ## 4 | 4 | 2 | 22.000 | 20.000 |
| ## 5 | 5 | 2 | 19.125 | 18.375 |
| ## 6 | 6 | 2 | 21.500 | 18.625 |
| ## 7 | 7 | 3 | 22.125 | 18.625 |
| ## 8 | 8 | 3 | 20.375 | 15.250 |
| ## 9 | 9 | 3 | 18.250 | 16.500 |
| ## 10 | 10 | 3 | 21.625 | 18.000 |
| ## 11 | 11 | 3 | 23.250 | 16.250 |
| ## 12 | 12 | 4 | 21.000 | 18.000 |
| ## 13 | 13 | 4 | 22.125 | 12.750 |
| ## 14 | 14 | 4 | 23.000 | 15.500 |
| ## 15 | 15 | 4 | 12.000 | 18.000 |

(i)

The null hypothesis is that there is *no* difference between heights of crossed and self-fertilized plants. That is,

$$H_0 : h_c = h_s$$

The alternative hypothesis is that there *is* a difference between heights of crossed and self-fertilized plants. That is

$$H_0 : h_c \neq h_s$$

```
t.test(x = Darwin.data$cross.height, y = Darwin.data$self.height, alternative = "two.sided")
```

```
##
## Welch Two Sample t-test
##
## data: Darwin.data$cross.height and Darwin.data$self.height
## t = 2.4371, df = 22.164, p-value = 0.02328
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## 0.3909566 4.8423767
## sample estimates:
## mean of x mean of y
## 20.19167 17.57500
```

The p -value is 0.02328, which is significant at the $\alpha = 0.05$ level. There is *strong evidence* that there is a difference between heights of crossed and self-fertilized plants.

The first assumption made regarding t-tests concerns the scale of measurement. The assumption for a t-test is that the scale of measurement applied to the data collected follows a continuous or ordinal scale, such as the scores for an IQ test.

The second assumption made is that of a simple random sample, that the data is collected from a representative, randomly selected portion of the total population.

The third assumption is the data, when plotted, results in a normal distribution, bell-shaped distribution curve.

The fourth assumption is a reasonably large sample size is used. A larger sample size means the distribution of results should approach a normal bell-shaped curve.

The final assumption is homogeneity of variance. Homogeneous, or equal, variance exists when the standard deviations of samples are approximately equal.

– What assumptions are made when conducting a t-test?

(ii)

```
# install.packages("coin", dependencies=TRUE, repos='http://cran.us.r-project.org')
# use permutation test here
```

(iii)

The null hypothesis is that there is *no* difference between heights of crossed and self-fertilized plants. That is,

$$H_0 : h_c = h_s$$

The alternative hypothesis is that there *is* a difference between heights of crossed and self-fertilized plants. That is

$$H_0 : h_c \neq h_s$$

```
wilcox.test(x = Darwin.data$cross.height, y = Darwin.data$self.height, alternative = "two.sided")
```

```
## Warning in wilcox.test.default(x = Darwin.data$cross.height, y =
## Darwin.data$self.height, : cannot compute exact p-value with ties
##
## Wilcoxon rank sum test with continuity correction
##
## data: Darwin.data$cross.height and Darwin.data$self.height
## W = 185.5, p-value = 0.002608
## alternative hypothesis: true location shift is not equal to 0
```

The p -value is 0.002608, which is significant at the $\alpha = 0.05$ level. There is *strong evidence* that there is a difference between heights of crossed and self-fertilized plants.

The Wilcoxon Sign test makes four important assumptions:

1. Dependent samples – the two samples need to be dependent observations of the cases. The Wilcoxon sign test assess for differences between a before and after measurement, while accounting for individual differences in the baseline.
2. Independence – The Wilcoxon sign test assumes independence, meaning that the paired observations are randomly and independently drawn.
3. Continuous dependent variable – Although the Wilcoxon signed rank test ranks the differences according to their size and is therefore a non-parametric test, it assumes that the measurements are continuous in theoretical nature. To account for the fact that in most cases the dependent variable is binominal distributed, a continuity correction is applied.
4. Ordinal level of measurement – The Wilcoxon sign test needs both dependent measurements to be at least of ordinal scale. This is necessary to ensure that the two values can be compared, and for each pair, it can be said if one value is greater, equal, or less than the other.

Furthermore, in order for the differences between measures to be rankable, the observations must be comparable. For every difference of observations, it must be clear which one is greater or if both observations are equal.

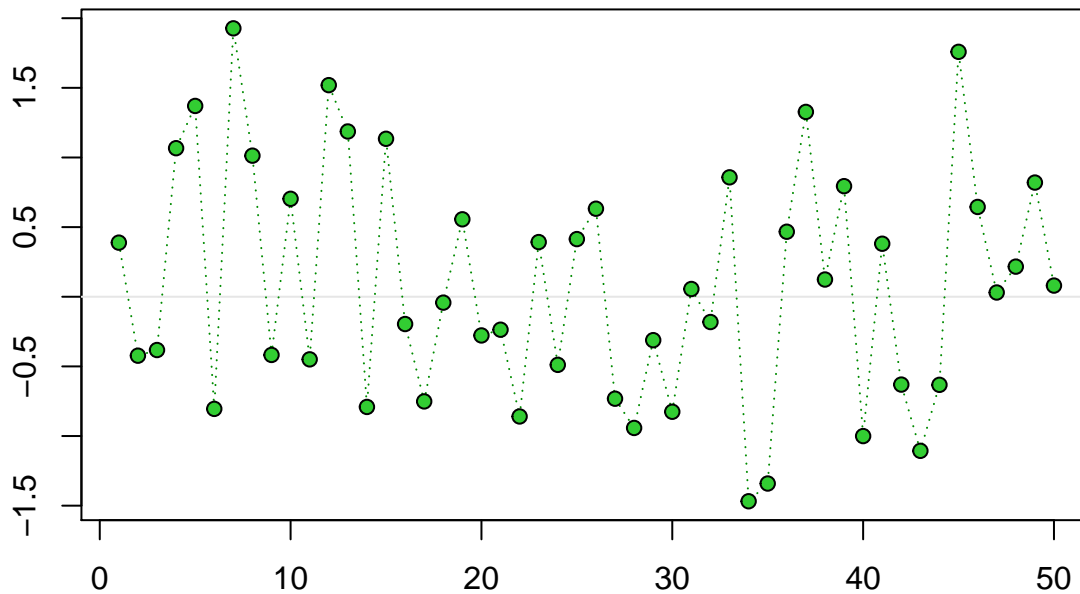
The test of significance of the Wilcoxon test further assumes that both samples have a continuous distribution function. This implies that tied ranks cannot occur. However, if tied ranks exist in the sample a continuity correction can be calculated. It is also possible to use an exact test that relies on permutation testing.

– Assumptions of the Wilcoxon Sign Test

```
demo(graphics)

##
##
## demo(graphics)
## ---- ~~~~~
##
## > # Copyright (C) 1997-2009 The R Core Team
## >
## > require(datasets)
##
## > require(grDevices); require(graphics)
##
## > ## Here is some code which illustrates some of the differences between
## > ## R and S graphics capabilities. Note that colors are generally specified
## > ## by a character string name (taken from the X11 rgb.txt file) and that line
## > ## textures are given similarly. The parameter "bg" sets the background
## > ## parameter for the plot and there is also an "fg" parameter which sets
## > ## the foreground color.
## >
## >
## > x <- stats::rnorm(50)
##
## > opar <- par(bg = "white")
##
## > plot(x, ann = FALSE, type = "n")
```

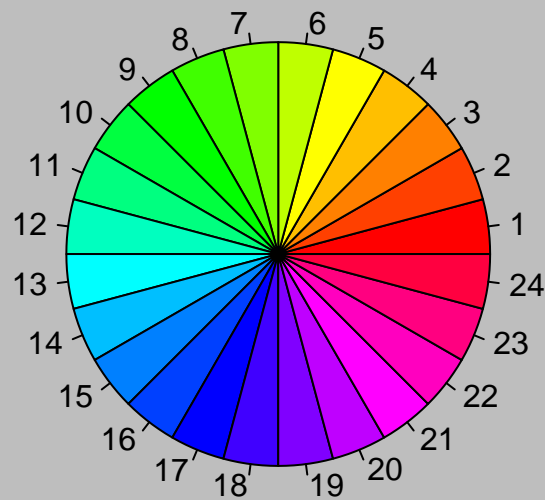
Simple Use of Color In a Plot



Just a Whisper of a Label

```
##
## > abline(h = 0, col = gray(.90))
##
## > lines(x, col = "green4", lty = "dotted")
##
## > points(x, bg = "limegreen", pch = 21)
##
## > title(main = "Simple Use of Color In a Plot",
## +       xlab = "Just a Whisper of a Label",
## +       col.main = "blue", col.lab = gray(.8),
## +       cex.main = 1.2, cex.lab = 1.0, font.main = 4, font.lab = 3)
##
## > ## A little color wheel.    This code just plots equally spaced hues in
## > ## a pie chart.    If you have a cheap SVGA monitor (like me) you will
## > ## probably find that numerically equispaced does not mean visually
## > ## equispaced.  On my display at home, these colors tend to cluster at
## > ## the RGB primaries.  On the other hand on the SGI Indy at work the
## > ## effect is near perfect.
## >
## > par(bg = "gray")
##
## > pie(rep(1,24), col = rainbow(24), radius = 0.9)
```

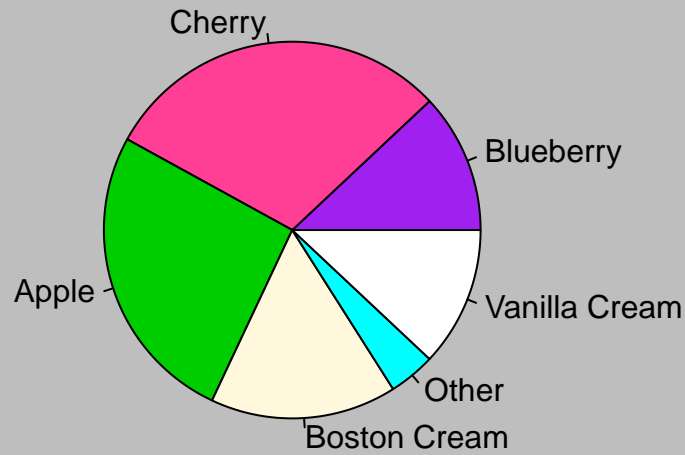

A Sample Color Wheel



(Use this as a test of monitor linearity)

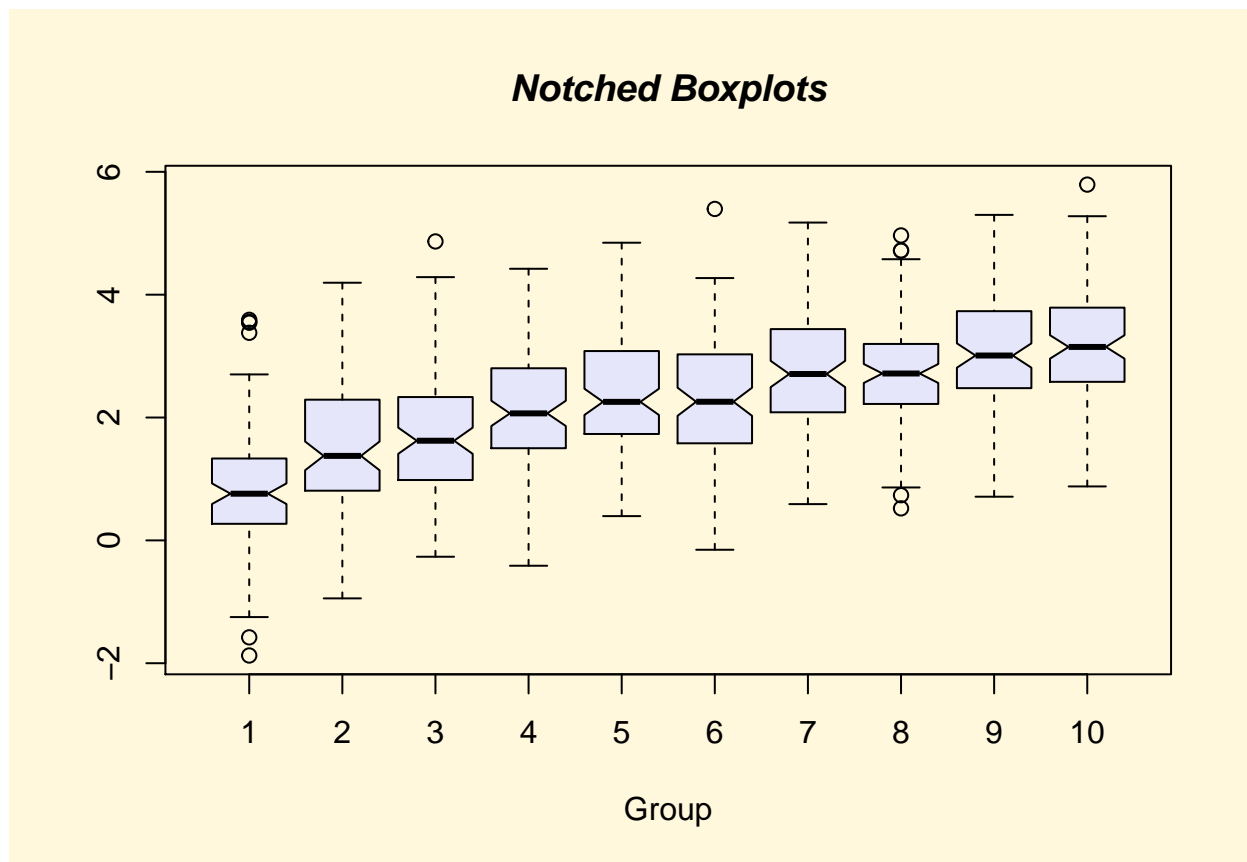
```
##
## > title(main = "A Sample Color Wheel", cex.main = 1.4, font.main = 3)
##
## > title(xlab = "(Use this as a test of monitor linearity)",
## +       cex.lab = 0.8, font.lab = 3)
##
## > ## We have already confessed to having these. This is just showing off X11
## > ## color names (and the example (from the postscript manual) is pretty "cute".
## >
## > pie.sales <- c(0.12, 0.3, 0.26, 0.16, 0.04, 0.12)
##
## > names(pie.sales) <- c("Blueberry", "Cherry",
## +                       "Apple", "Boston Cream", "Other", "Vanilla Cream")
##
## > pie(pie.sales,
## +     col = c("purple", "violetred1", "green3", "cornsilk", "cyan", "white"))
```

January Pie Sales



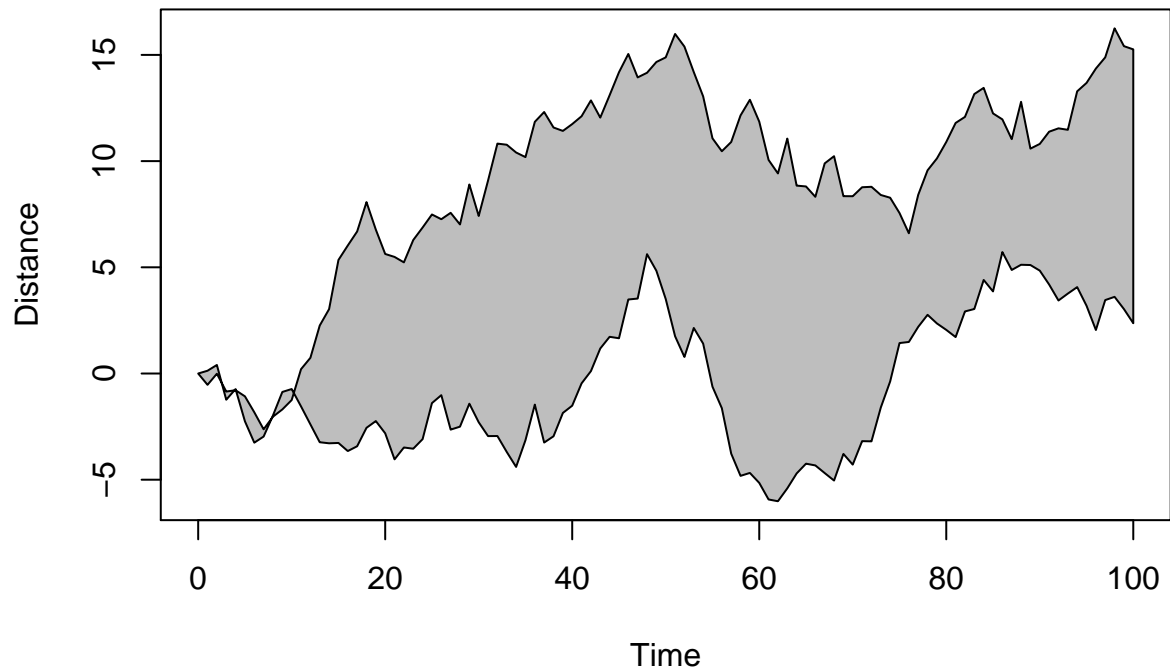
(Don't try this at home kids)

```
##  
## > title(main = "January Pie Sales", cex.main = 1.8, font.main = 1)  
##  
## > title(xlab = "(Don't try this at home kids)", cex.lab = 0.8, font.lab = 3)  
##  
## > ## Boxplots: I couldn't resist the capability for filling the "box".  
## > ## The use of color seems like a useful addition, it focuses attention  
## > ## on the central bulk of the data.  
## >  
## > par(bg="cornsilk")  
##  
## > n <- 10  
##  
## > g <- gl(n, 100, n*100)  
##  
## > x <- rnorm(n*100) + sqrt(as.numeric(g))  
##  
## > boxplot(split(x,g), col="lavender", notch=TRUE)
```

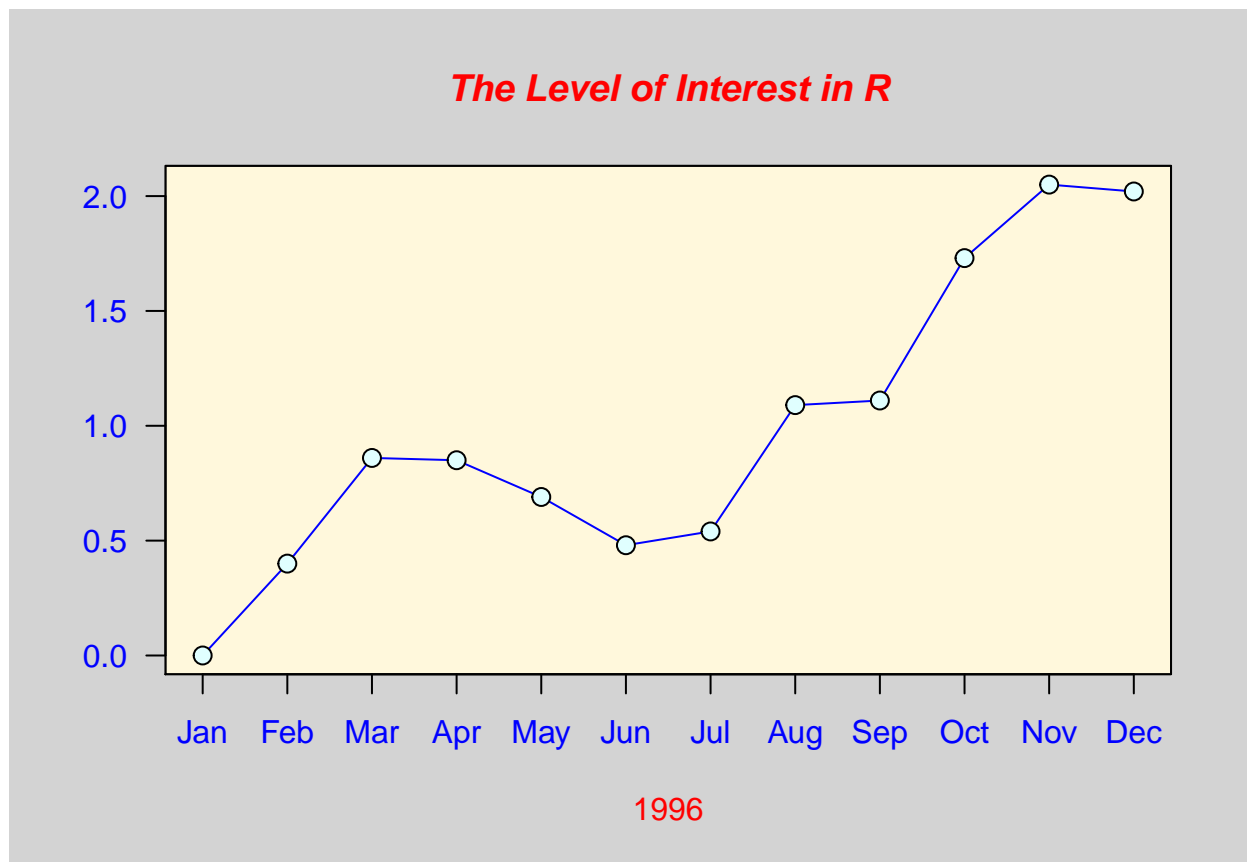


```
##
## > title(main="Notched Boxplots", xlab="Group", font.main=4, font.lab=1)
##
## > ## An example showing how to fill between curves.
## >
## > par(bg="white")
##
## > n <- 100
##
## > x <- c(0,cumsum(rnorm(n)))
##
## > y <- c(0,cumsum(rnorm(n)))
##
## > xx <- c(0:n, n:0)
##
## > yy <- c(x, rev(y))
##
## > plot(xx, yy, type="n", xlab="Time", ylab="Distance")
```

Distance Between Brownian Motions

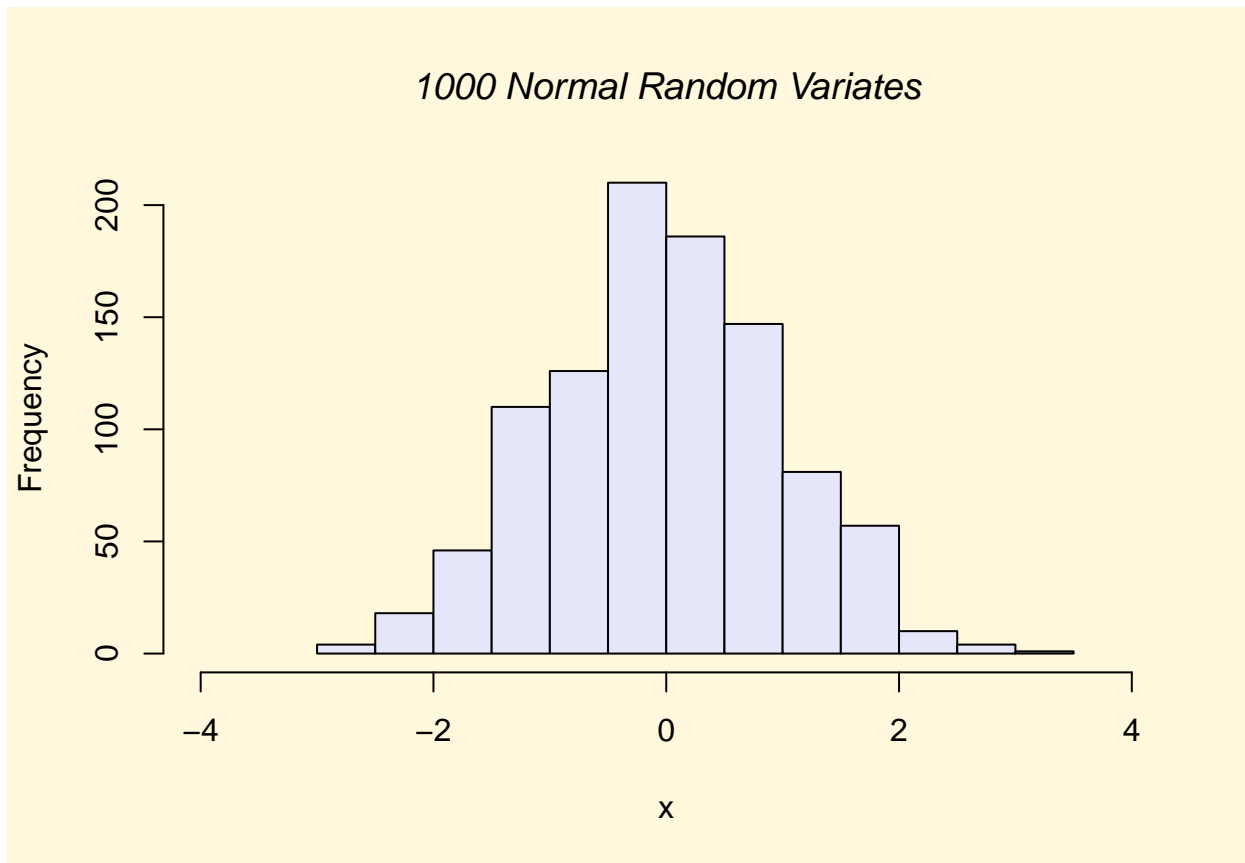


```
##  
## > polygon(xx, yy, col="gray")  
##  
## > title("Distance Between Brownian Motions")  
##  
## > ## Colored plot margins, axis labels and titles.    You do need to be  
## > ## careful with these kinds of effects.    It's easy to go completely  
## > ## over the top and you can end up with your lunch all over the keyboard.  
## > ## On the other hand, my market research clients love it.  
## >  
## > x <- c(0.00, 0.40, 0.86, 0.85, 0.69, 0.48, 0.54, 1.09, 1.11, 1.73, 2.05, 2.02)  
##  
## > par(bg="lightgray")  
##  
## > plot(x, type="n", axes=FALSE, ann=FALSE)
```

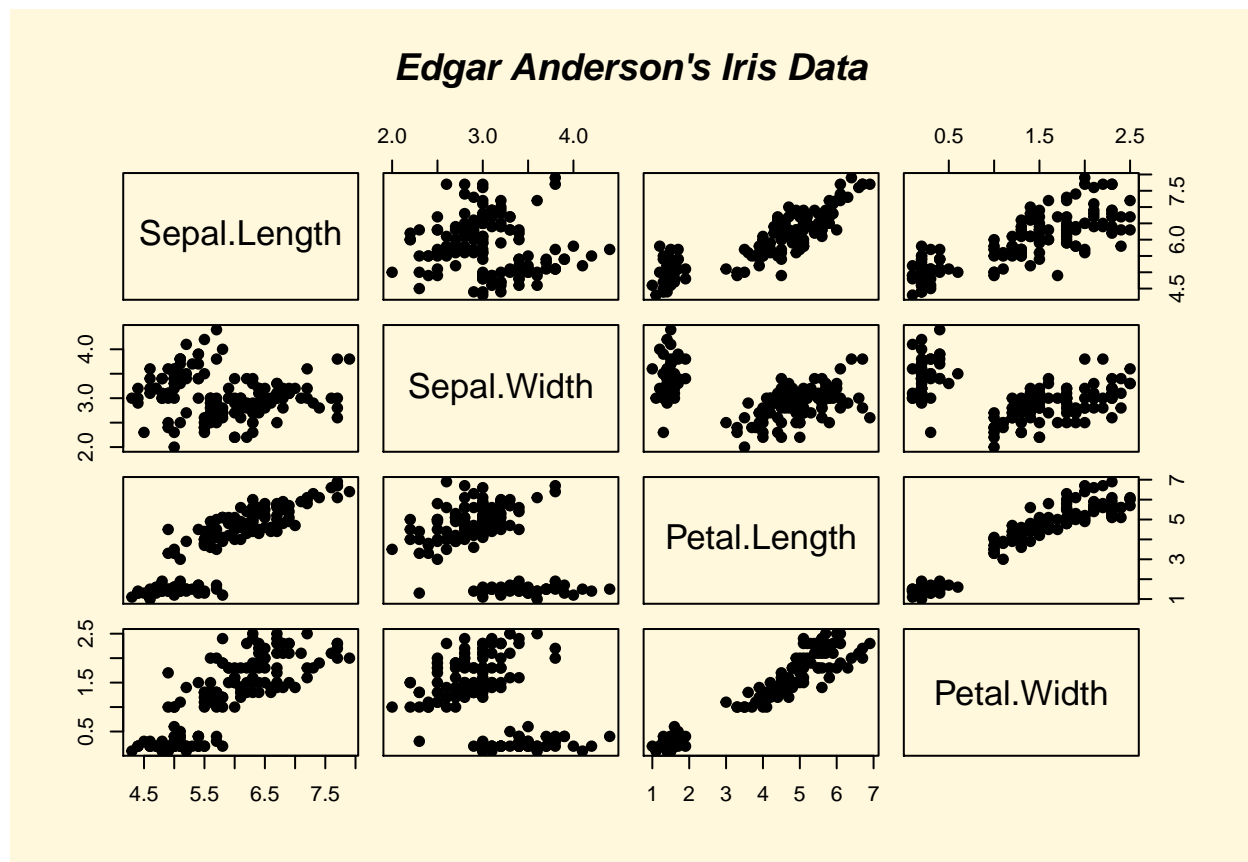


```
##
## > usr <- par("usr")
##
## > rect(usr[1], usr[3], usr[2], usr[4], col="cornsilk", border="black")
##
## > lines(x, col="blue")
##
## > points(x, pch=21, bg="lightcyan", cex=1.25)
##
## > axis(2, col.axis="blue", las=1)
##
## > axis(1, at=1:12, lab=month.abb, col.axis="blue")
##
## > box()
##
## > title(main= "The Level of Interest in R", font.main=4, col.main="red")
##
## > title(xlab= "1996", col.lab="red")
##
## > ## A filled histogram, showing how to change the font used for the
## > ## main title without changing the other annotation.
## >
## > par(bg="cornsilk")
##
## > x <- rnorm(1000)
##
```

```
## > hist(x, xlim=range(-4, 4, x), col="lavender", main="")
```

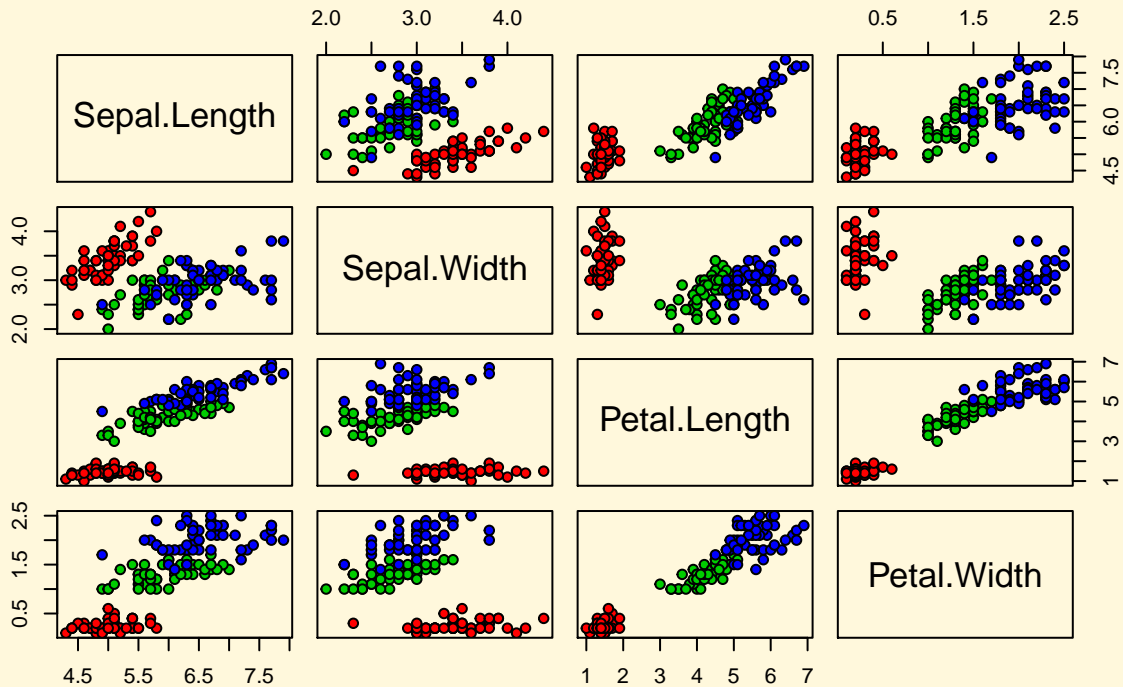


```
##  
## > title(main="1000 Normal Random Variates", font.main=3)  
##  
## > ## A scatterplot matrix  
## > ## The good old Iris data (yet again)  
## >  
## > pairs(iris[1:4], main="Edgar Anderson's Iris Data", font.main=4, pch=19)
```



```
##
## > pairs(iris[1:4], main="Edgar Anderson's Iris Data", pch=21,
## +       bg = c("red", "green3", "blue")[unclass(iris$Species)])
```

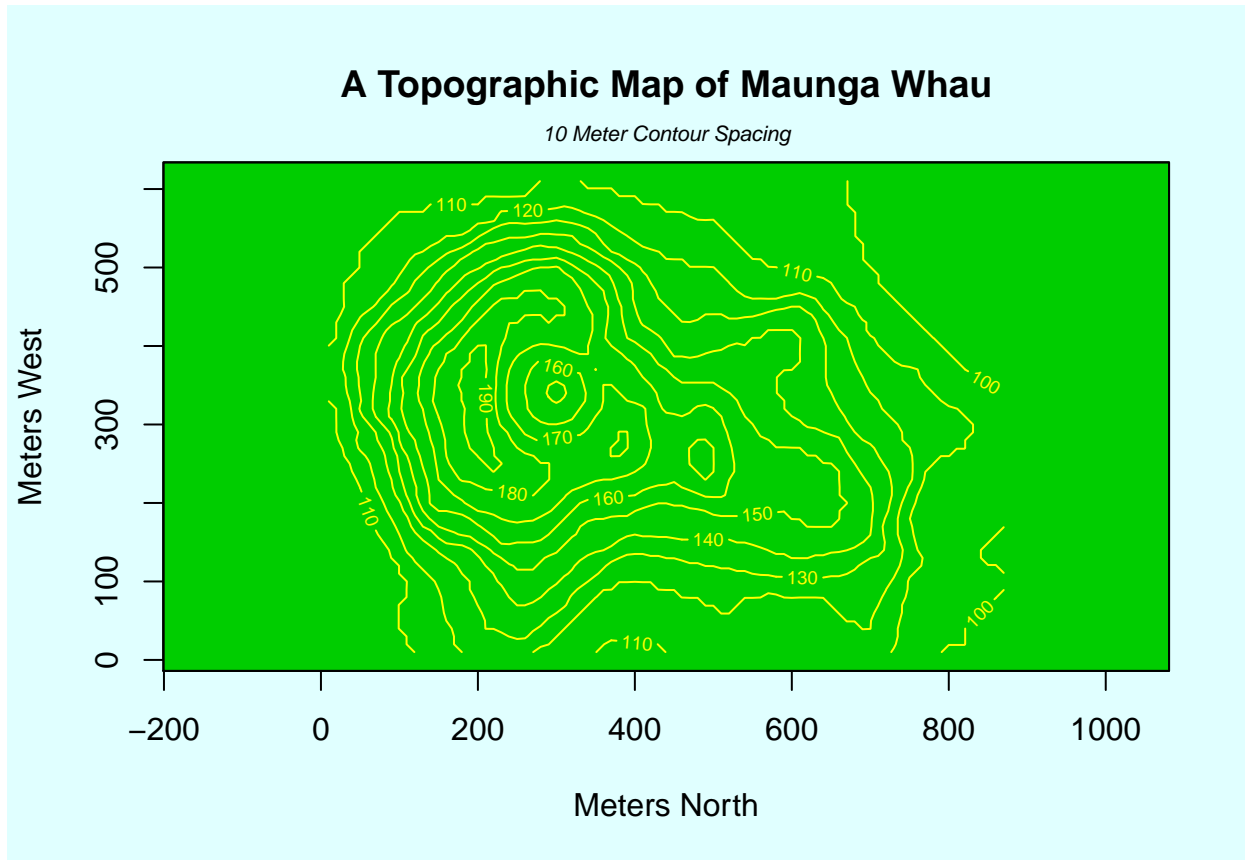
Edgar Anderson's Iris Data



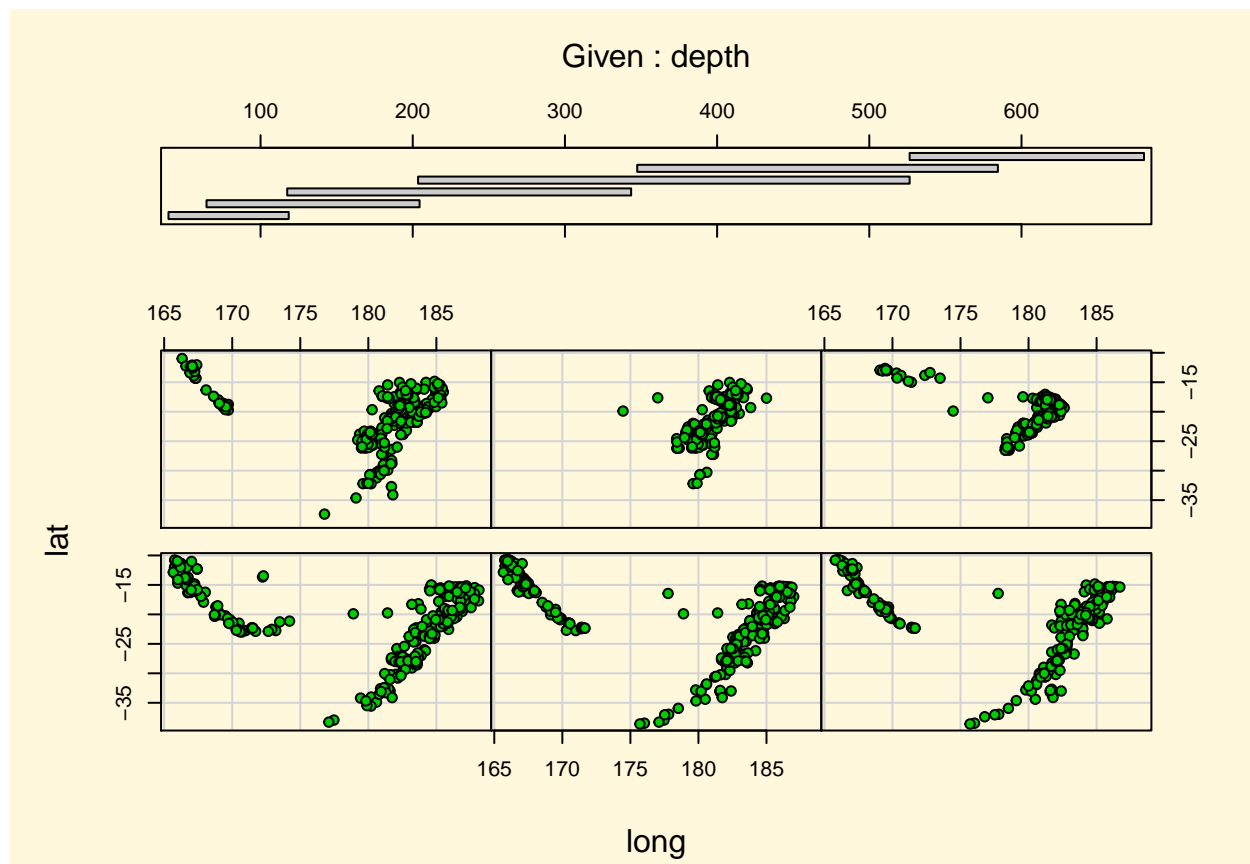
```
##
## > ## Contour plotting
## > ## This produces a topographic map of one of Auckland's many volcanic "peaks".
## >
## > x <- 10*1:nrow(volcano)
##
## > y <- 10*1:ncol(volcano)
##
## > lev <- pretty(range(volcano), 10)
##
## > par(bg = "lightcyan")
##
## > pin <- par("pin")
##
## > xdelta <- diff(range(x))
##
## > ydelta <- diff(range(y))
##
## > xscale <- pin[1]/xdelta
##
## > yscale <- pin[2]/ydelta
##
## > scale <- min(xscale, yscale)
##
## > xadd <- 0.5*(pin[1]/scale - xdelta)
##
```



```
## > yadd <- 0.5*(pin[2]/scale - ydelta)
##
## > plot(numeric(0), numeric(0),
## +      xlim = range(x)+c(-1,1)*xadd, ylim = range(y)+c(-1,1)*yadd,
## +      type = "n", ann = FALSE)
```



```
##
## > usr <- par("usr")
##
## > rect(usr[1], usr[3], usr[2], usr[4], col="green3")
##
## > contour(x, y, volcano, levels = lev, col="yellow", lty="solid", add=TRUE)
##
## > box()
##
## > title("A Topographic Map of Maunga Whau", font= 4)
##
## > title(xlab = "Meters North", ylab = "Meters West", font= 3)
##
## > mtext("10 Meter Contour Spacing", side=3, line=0.35, outer=FALSE,
## +      at = mean(par("usr")[1:2]), cex=0.7, font=3)
##
## > ## Conditioning plots
## >
## > par(bg="cornsilk")
##
## > coplot(lat ~ long | depth, data = quakes, pch = 21, bg = "green3")
```



```
##
## > par(opar)
```