

STAT 2005 – PROGRAMMING LANGUAGES FOR STATISTICS

TUTORIAL 3 STATISTICAL GRAPHICS

2020

LIU Ran

Department of Statistics, The Chinese University of Hong Kong

1 Statistical Graphics

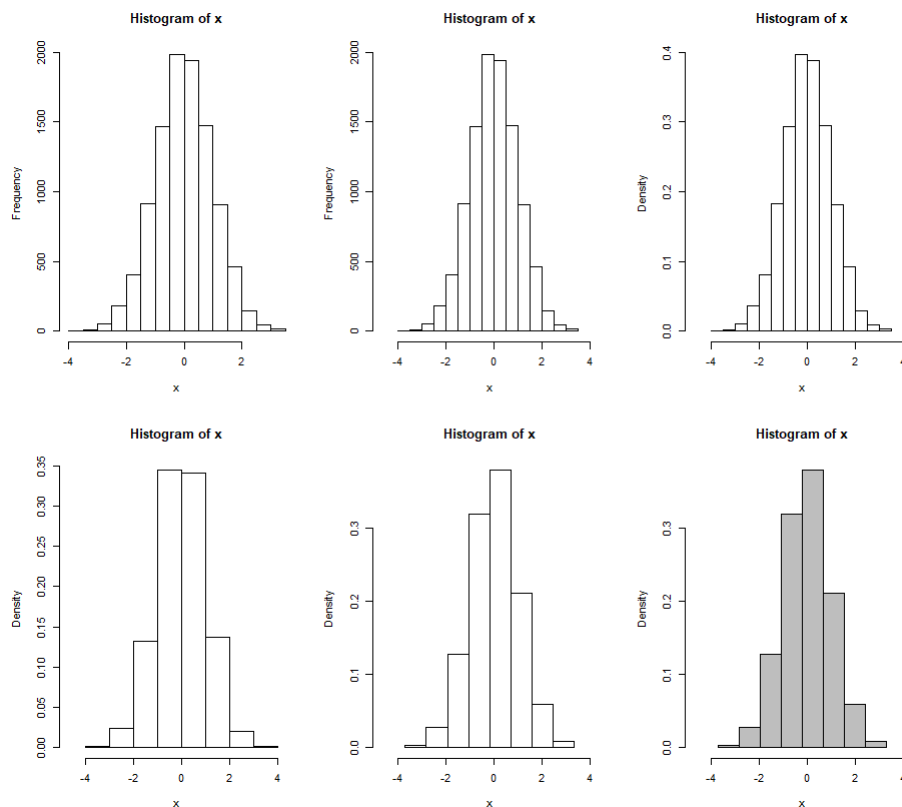
We are going to review the concepts about the histogram, qqnorm and boxplot commands. And after that `matplot` command will also be mentioned.

1.1 histogram

A histogram is an approximate representation of the distribution of numerical data. When we want to check the density function of a dataset, it is helpful to draw a histogram.

```

1 set.seed(2005)
2 x = rnorm(10000)
3 par(mfrow=c(2,3))
4 hist(x)
5 hist(x,xlim = c(min(x)-1,max(x)+1))
6 # xlim means the value range of xaxis
7 hist(x,xlim = c(min(x)-1,max(x)+1), freq = F)
8 # freq represents the type for yaxis(density or frequency)
9 hist(x,xlim = c(min(x)-1,max(x)+1), freq = F, nclass = 10)
10 # nclass means the number of bins
11 hist(x,xlim = c(min(x)-1,max(x)+1), freq = F, breaks = seq(min(x),max(x),length.out = 9))
12 # You can also use 'breaks' argument to divide the dataset
13 hist(x,xlim = c(min(x)-1,max(x)+1), freq = F, breaks = seq(min(x),max(x),length.out = 9),
14 col='grey')
```



Remark 1.1. `freq` if TRUE, the histogram graphic is a representation of frequencies, the counts component of the result; if FALSE, density values, are plotted.

Remark 1.2. The density value 0.4 may not be the proportion of this bin, which means $2000/10000 \neq 0.4$. Because $\int_a^b f(x)dx = P(X \in [a, b])$, the density value multiplied by the width of this bin roughly equals to the proportion $0.4 * width \approx \frac{2000}{10000}$.

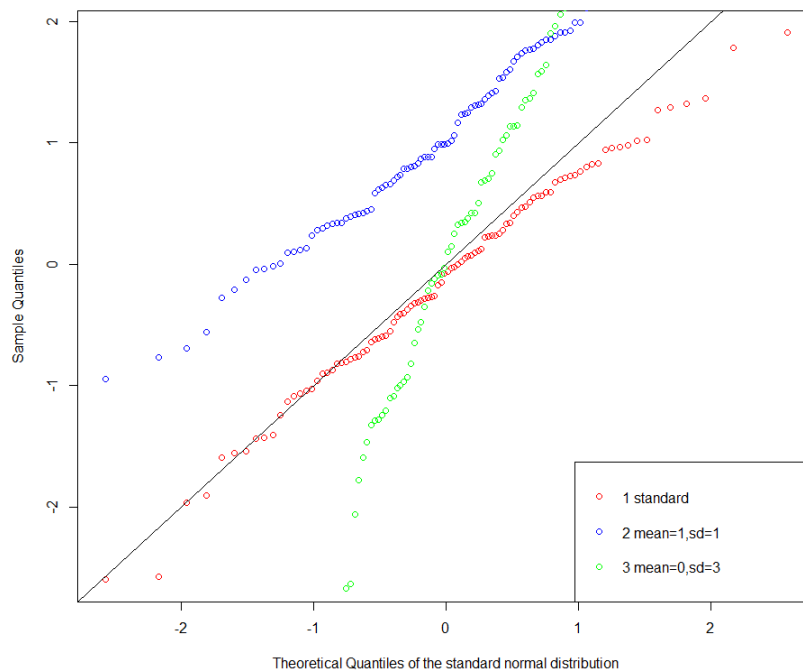
1.2 qqnorm

A Q-Q (quantile-quantile) plot is a probability plot, which is a graphical method for comparing two probability distributions by plotting their quantiles against each other. `qqnorm` compares the data with the standard normal distribution.

```

1 set.seed(2005)
2 y1 = rnorm(100, mean = 0, sd = 1)
3 y2 = rnorm(100, mean = 1, sd = 1)
4 y3 = rnorm(100, mean = 0, sd = 3)
5 x1 = qqnorm(y1, plot.it = F) # not plot it individually
6 x2 = qqnorm(y2, plot.it = F)
7 x3 = qqnorm(y3, plot.it = F)
8 plot(x1, col = 'red',
9       xlab = 'Theoretical Quantiles of the standard normal distribution',
10      ylab = 'Sample Quantiles')
11 points(x2, col = 'blue')
12 points(x3, col = 'green')
13 abline(a=0, b=1) # y=x
14 # If your data satisfy the defined distribution, the fitted line will be y = x
15 legend('bottomright', legend = paste(1:3, c('standard', 'mean=1, sd=1', 'mean=0, sd=3')),
16        pch = 1, col = c('red', 'blue', 'green'))

```



Remark 1.3. The `qqnorm` function will give the standard normal distribution quantiles. If the data satisfy the standard normal distribution, points generated by `qqnorm` will be around the line $y = x$. If data satisfy a normal distribution but it has not standard mean or variance, the slope and the intercept will be influenced. And from the figure we can see the variance affects the slope and the mean affects the intercept. But no matter what mean and variance are, as long as your data is generated by a normal distribution, it will still be a straight line. Therefore, it is appropriate to use `qqnorm` to see if the data satisfy a normal distribution.

1.3 boxplot

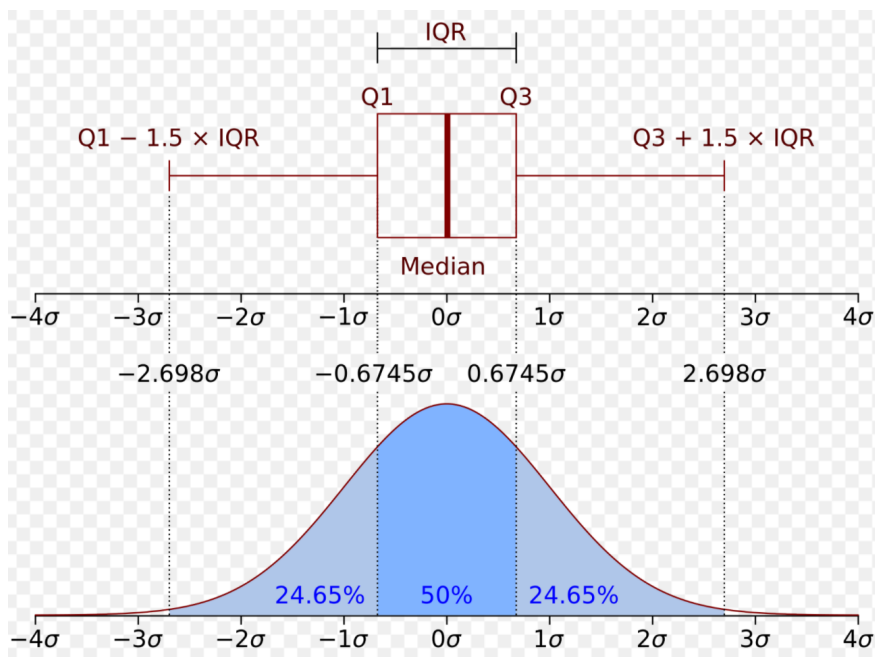
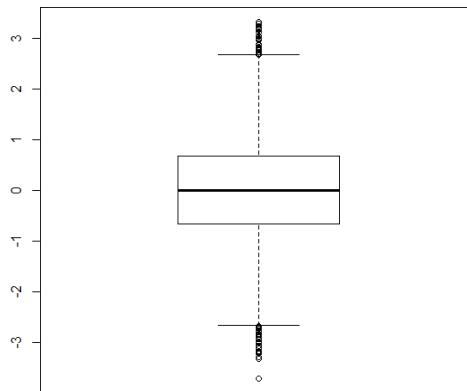
A box plot is a method for graphically depicting groups of numerical data through their quartile . A quartile is a type of quantile which divides the number of data points into four parts.

1. Median (Q2 / 50th percentile)
2. First quartile (Q1 / 25th percentile)
3. Third quartile (Q3 / 75th percentile)
4. Interquartile range (IQR) = $Q3 - Q1$

```

1 boxplot(x) # Box Plot
2 boxplot(x)$out # the outliers
3 index = which(x %in% boxplot(x)$out) # find the index of outliers

```



Remark 1.4. Although the histogram seems to be more intuitive, the box plot needs much less space and you can draw different groups' box plots in one figure. Meanwhile, it does not need to choose the number of bins or the width between breaks, which heavily affect the appearance of a histogram.

1.4 Add Elements To Existing Graphs

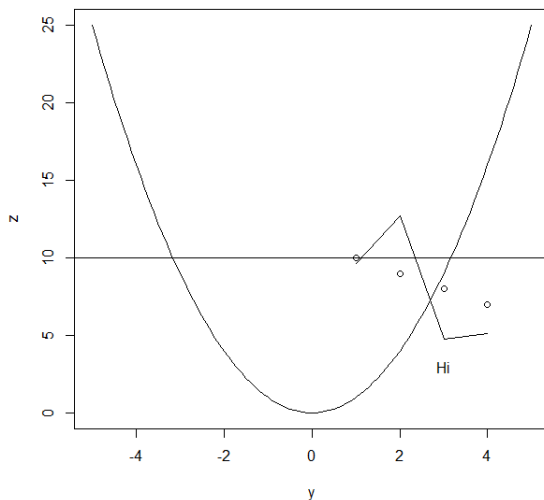
| Function | Description |
|--------------------------------|--|
| <code>abline</code> (a, b) | A line with intercept a and slope b |
| <code>points</code> (x, y) | Points (x, y) |
| <code>lines</code> (x, y) | Several line segments pass through points (x, y) |
| <code>text</code> (x, y, text) | Text at (x, y) |

```

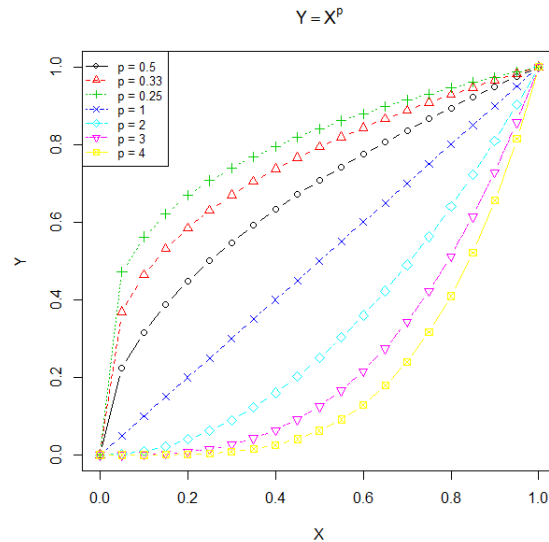
1 set.seed(2005)
2 y = seq(-5,5,length.out=10)
3 z = y**2
4 plot(y, z, type="l")
5 abline(h=10) # add a horizontal line
6 # use v = 0 to draw a vertical line
7 points(1:4, 10:7) # Add points
8 # x coordinates c(1,2,3,4)
9 # y coordinates c(10,9,8,7)
10 # point coordinates (1,10), (2,9), (3,8), (4,7)
11 lines(1:4, 10*abs(rnorm(4))) # Add three line segments
12 text(3, 3, "Hi") # Text

```

Add Elements



matplot



2 matplot

- `matplot`: Plot several lines in one graph at the same time.

```

1 # plot the figure for f(x) = x^p
2 x = seq(0, 1, by=0.05)
3 n = 4
4 p = c(1/(2:n), 1:n)
5 N = length(p)
6 X = matrix(rep(x, N), ncol=N) # each column of X is the vector x
7 Y = t(t(X)^p) # Y[,j] = X[,j]^p[j] calculate the function values by columns
8 col = pch = lty = 1:N
9 matplot(X, Y, main=expression(Y==X^p), pch=pch, col=col, lty=lty, type="b")
10 # We can use the expression command to show the mathematical formula
11 legend("topleft", legend = paste0("p = ", round(p, 2)),
12       pch=pch, col=col, lty=lty, cex=0.8)
13 # round(x,4): rounds the values in its first argument to the specified number of decimal places

```

```

1 > p
2 [1] 0.5000000 0.3333333 0.2500000 1.0000000 2.0000000 3.0000000 4.0000000

```

```

3 > X
4      [,1] [,2] [,3] [,4] [,5] [,6] [,7]
5 [1,] 0.00 0.00 0.00 0.00 0.00 0.00 0.00
6 [2,] 0.05 0.05 0.05 0.05 0.05 0.05 0.05
7 [3,] 0.10 0.10 0.10 0.10 0.10 0.10 0.10
8 [4,] 0.15 0.15 0.15 0.15 0.15 0.15 0.15
9 [5,] 0.20 0.20 0.20 0.20 0.20 0.20 0.20
10 [6,] 0.25 0.25 0.25 0.25 0.25 0.25 0.25
11 [7,] 0.30 0.30 0.30 0.30 0.30 0.30 0.30
12 [8,] 0.35 0.35 0.35 0.35 0.35 0.35 0.35
13 [9,] 0.40 0.40 0.40 0.40 0.40 0.40 0.40
14 [10,] 0.45 0.45 0.45 0.45 0.45 0.45 0.45
15 [11,] 0.50 0.50 0.50 0.50 0.50 0.50 0.50
16 [12,] 0.55 0.55 0.55 0.55 0.55 0.55 0.55
17 [13,] 0.60 0.60 0.60 0.60 0.60 0.60 0.60
18 [14,] 0.65 0.65 0.65 0.65 0.65 0.65 0.65
19 [15,] 0.70 0.70 0.70 0.70 0.70 0.70 0.70
20 [16,] 0.75 0.75 0.75 0.75 0.75 0.75 0.75
21 [17,] 0.80 0.80 0.80 0.80 0.80 0.80 0.80
22 [18,] 0.85 0.85 0.85 0.85 0.85 0.85 0.85
23 [19,] 0.90 0.90 0.90 0.90 0.90 0.90 0.90
24 [20,] 0.95 0.95 0.95 0.95 0.95 0.95 0.95
25 [21,] 1.00 1.00 1.00 1.00 1.00 1.00 1.00
26 > Y
27      [,1]      [,2]      [,3] [,4]      [,5]      [,6]      [,7]
28 [1,] 0.0000000 0.0000000 0.0000000 0.00 0.0000 0.0000000 0.0000000
29 [2,] 0.2236068 0.3684031 0.4728708 0.05 0.0025 0.000125 0.00000625
30 [3,] 0.3162278 0.4641589 0.5623413 0.10 0.0100 0.001000 0.00010000
31 [4,] 0.3872983 0.5313293 0.6223330 0.15 0.0225 0.003375 0.00050625
32 [5,] 0.4472136 0.5848035 0.6687403 0.20 0.0400 0.008000 0.00160000
33 [6,] 0.5000000 0.6299605 0.7071068 0.25 0.0625 0.015625 0.00390625
34 [7,] 0.5477226 0.6694330 0.7400828 0.30 0.0900 0.027000 0.00810000
35 [8,] 0.5916080 0.7047299 0.7691606 0.35 0.1225 0.042875 0.01500625
36 [9,] 0.6324555 0.7368063 0.7952707 0.40 0.1600 0.064000 0.02560000
37 [10,] 0.6708204 0.7663094 0.8190363 0.45 0.2025 0.091125 0.04100625
38 [11,] 0.7071068 0.7937005 0.8408964 0.50 0.2500 0.125000 0.06250000
39 [12,] 0.7416198 0.8193213 0.8611735 0.55 0.3025 0.166375 0.09150625
40 [13,] 0.7745967 0.8434327 0.8801117 0.60 0.3600 0.216000 0.12960000
41 [14,] 0.8062258 0.8662391 0.8979008 0.65 0.4225 0.274625 0.17850625
42 [15,] 0.8366600 0.8879040 0.9146912 0.70 0.4900 0.343000 0.24010000
43 [16,] 0.8660254 0.9085603 0.9306049 0.75 0.5625 0.421875 0.31640625
44 [17,] 0.8944272 0.9283178 0.9457416 0.80 0.6400 0.512000 0.40960000
45 [18,] 0.9219544 0.9472682 0.9601846 0.85 0.7225 0.614125 0.52200625
46 [19,] 0.9486833 0.9654894 0.9740037 0.90 0.8100 0.729000 0.65610000
47 [20,] 0.9746794 0.9830476 0.9872585 0.95 0.9025 0.857375 0.81450625
48 [21,] 1.0000000 1.0000000 1.0000000 1.00 1.0000 1.000000 1.00000000

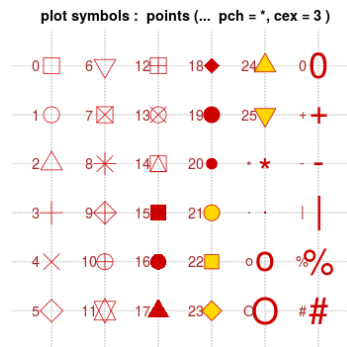
```

Remark 2.1. Suitable legend should be added for understanding and interpretation.

3 Useful things

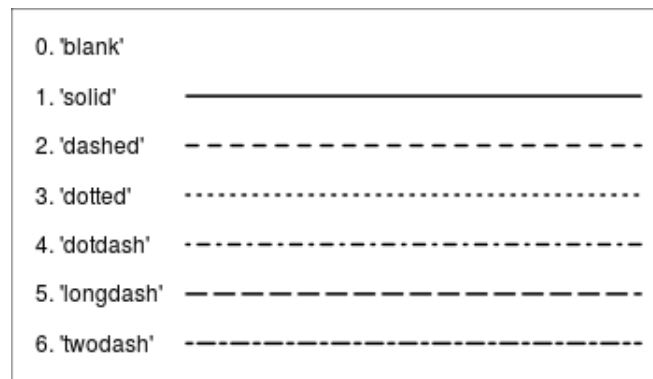
3.1 pch

Use the `pch=` option to specify symbols to use when plotting points. For symbols 21 through 25, specify border color (`col=`) and fill color (`bg=`).



3.2 lty

You can use `lty=` to specify the type of lines. You can choose the type you want by the integer index or the character.



3.3 Color

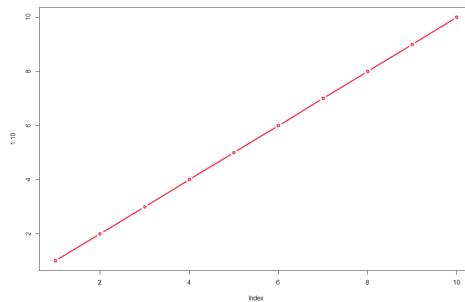
You can specify colors in R by index, name, hexadecimal, or RGB. For example `col=1`, `col="white"`, and `col="#FFFFFF"` are equivalent. Because the colors in R are too many to choose, sometimes we select the optimal color from the Crayola Crayon colors(optional). The following are the Crayola Crayon colors:



Because the names in the table do not have the same names in R, for example, you cannot use `col='Scarlet'` directly. You should load the package 'broman' to get the corresponding hexadecimal of the color you want.

```
1 library(broman) # contains the crayons color
2 plot_crayons() # display the above figure
3 colors = brocolors("crayons") # return a list which contains all the colors by their hexadecimal
4 sel_color = colors['Scarlet'] # choose the color you want to plot
5 plot(1:10,type='b',lwd = 3,col = sel_color) # use the color
```

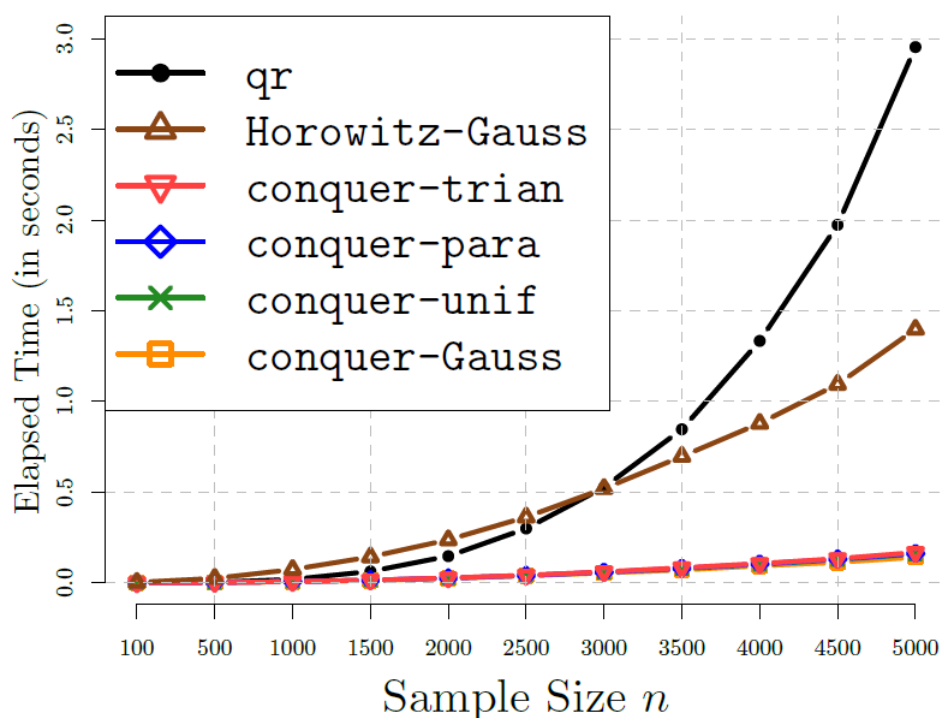
Remark 3.1. Go to wiki to check more detail about Crayola Crayon colors.



4 Example(Optional)

How to draw a nice figure:

1. Be familiar with commands for drawing, such as the parameter commands for `plot(base)`. Try to know the effect of each command.
2. Watch many wonderful figures in others' paper(Nature, Science...), and try to find the advantages of their setting.
3. Try to reproduce the figures by yourself.
4. Create your own setting for next time use.



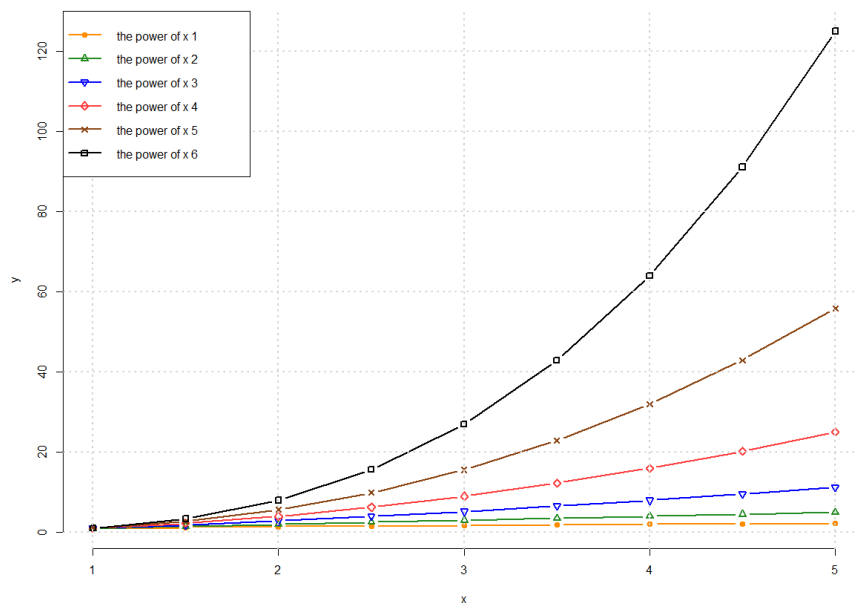
1. Generate similar data, here we use $f(x) = x^{p/2}$ power function to generate similar data.
2. Notice the `par()` setting, the isolated axis means we should use `btty = 'n'`.
3. Use the same `pch` and `color` parameters.
4. Use the same type of lines and proper line width, here we use the solid type and `lwd = 2`.
5. Notice the grid setting, here we use `grid(lwd=2)`.
6. Make a nice legend.
7. Detail adjustment.(lab title, axis range, main title)

Remark 4.1. Go to related websites to pick the colors you want from images: [imagecolorpic](#).

```

1 par(bty = 'n')
2 index = 1:6
3 x_range = seq(1,5,by=0.5)
4 sines <- outer(x_range, index, function(x, p) x^(p/2))
5 # generate a matrix: sines[i,j] = func(x_range[i],index[j])
6 matplot(x_range,sines)
7 my_shape = as.integer(c('19','2','6','5','4','0'))
8 # black, brown, red, blue, green, yellow
9 # many similar colors, col = 'green' is not the one we want
10 my_color = c('#000000','#8B4513','#FF4040','#0000FF','#228B22','#FF8C00')
11 my_color = my_color[seq(length(my_color),1)]
12 matplot(x_range,sines, type = "b",
13         pch = my_shape[index],
14         lty = 1,
15         col = my_color[index],
16         lwd = 2,
17         xlab = 'x',
18         ylab = 'y',
19         xlim = c(min(x_range),max(x_range)),
20         ylim = c(min(sines),max(sines)),
21         )
22 grid(lwd=2)
23 legend('topleft', legend = paste('the power of x',index),
24        col = my_color[index],
25        pch = my_shape[index],
26        lty = 1,
27        lwd = 2)

```



After reproducing others' figures, you can store the setting, such as `pch` and `color`, and next time, use these settings to plot your own figures.