

## Agenda for today

1. Descriptive Data Analysis
2. Graphics



## Prerequisites

```
1 library("xplore")
2 library("stats")
3 setenv("outputstringformat", "%s")
```

Avoid quotes in text output.



## Descriptive Data Analysis

- typically the first part of statistical modeling
- evaluation of data
- All routines from libraries xplore (basic routines) and stats (basic statistical methods)



## Data Matrices

- $z = \#(x_1, x_2, \dots, x_n)$

creates a column vector  $z$  from scalar numbers  $x_1, x_2, \dots, x_n$

- $z = x|y$

concatenates two arrays  $x$  and  $y$  rowwise

- $z = x \sim y$

concatenates two arrays  $x$  and  $y$  columnwise



## Reading Data

```
1 x = read("file")  
2   reads numeric data from file.dat  
3 x = readm("file")  
4   reads mixed text and numeric data from file.dat
```



## Dimensions of a Dataset

```
1 d = dim(x)
2     shows the dimension of an array x
3 n = rows(x)
4     shows the number of rows of an array x
5 p = cols(x)
6     shows the number of columns of an array x
7 y = x[i,j] or y = x[i,] or y = x[:,j]
8     extracts element i,j or row i or column j from x
9 z = x[k:l,m:n]
10    extracts rows k to l and columns m to n
```



## Minimum and Maximum

```
1 mx = min(x {,d})  
2     computes the minimum of an array x, optionally  
3     with respect to dimension d  
3 mx = max(x {,d})  
4     computes the maximum of an array x, optionally  
4     with respect to dimension d
```



## Mean, Variance and other Moments

```
1 mx = mean(x {,d})
2     computes the mean of an array x, optionally with
3     respect to dimension d
3 vx = var(x {,d})
4     computes the variance of an array x, optionally
5     with respect to dimension d
5 kx = kurtosis(x)
6     computes the (columnwise) kurtosis of an array x
7 sx = skewness(x)
8     computes the (columnwise) skewness of an array x
```





## Median and Quantiles

```
1 mx = median(x)
2     computes the (columnwise) median of an array x
3 qx = quantile(x, alpha)
4     computes the (columnwise) quantile of an array x
     at level alpha
```



## Covariance and Correlation

```
1 cx = cov(x)
2     computes the covariance matrix of a data matrix
   x
3 rx = corr(x)
4     computes the correlation matrix of a data matrix
   x
```



## Categorical Data

```
1 {xr, r} = discrete(x {,y})  
2     reduces a matrix to its distinct rows  
3     and gives the number of replications of each row  
4     in the original data set
```



## Categorical Data

```
1 setenv("outputstringformat", "%s")
2 library("xplore")
3 library("stats")
4 earn=read("cps85")
5 earn=earn[,1|2|5|8|10|11|12]
6 {cat, freq}=discrete(earn[,3])
7 cat
8 freq
```



## Missing Data

```
1 nx = countNaN(x)
2     counts missing values in an array x
3 nx = countNotNumber(x)
4     counts missing and infinite values in an array x
5 ix = isNaN(x)
6     determines whether the elements of an array x
7     are missing values
7 ix = isInf(x)
8     determines whether the elements of an array x
8     are infinite values
```



## Missing Data

```
1 ix = isNumber (x)
2     determines whether the elements of an array x
   are regular numeric values
3 y = paf(x, i)
4     deletes all rows of x for which the
   corresponding element of i equals 0
5 y = replace(x, w, b)
6     replaces all elements of x which equal w by the
   value b
```



## Summarizing Information

```
1 s = summarize(x {,xvars})  
2   computes a short summary of descriptive  
   statistics  
3 s = fivenum(x {,xvars})  
4   computes the five number summary for each column  
   of a matrix x  
5 s = descriptive(x {,xvars})  
6   computes detailed descriptive statistics for  
   each column of a matrix x
```

optionally a vector of variable names `xvars` can be given



## Summarizing Metric Data

```
1 s = summarize(x {,xvars})  
2   computes a short summary of descriptive  
   statistics  
3 s = fivenum(x {,xvars})  
4   computes the five number summary for each column  
   of a matrix x  
5 s = descriptive(x {,xvars})  
6   computes detailed descriptive statistics for  
   each column of a matrix x
```

optionally a vector of variable names `xvars` can be given





## Summarizing Categorical Data

```
1 s = frequency(x {, xvars {, outwidth}})
2   computes frequency table for each column of a
   matrix x
3 s = crosstable(x{,xvars})
4   computes pairwise cross tables from all columns
5   computes the result of a  $\chi^2$  independence
   test
```

optionally a vector of variable names `xvars` can be given



## Graphics Overview

- the graphical tools (high-level): plot\*
- the graphical primitives (low-level) gr\*
- the graphical commands



## Basic Plotting

```
1 plot(x1 {, x2 {,...{x5}}})  
2     plots the data sets x1,...,x5  
3 line(x1 {, x2 {, ... {x5}}})  
4     plots the lines sets x1, ..., x5  
5 y = setmask (x, opt1 {, opt2 {, ... {opt9}}})  
6     modifies a data set for plotting  
7 disp = createdisplay (r,c)  
8     creates a display disp  
9 show(disp, i, j, x1 {, x2 {, ... {, xn}}})  
10    plots the data sets x1, ..., xn in the display  
     disp d
```

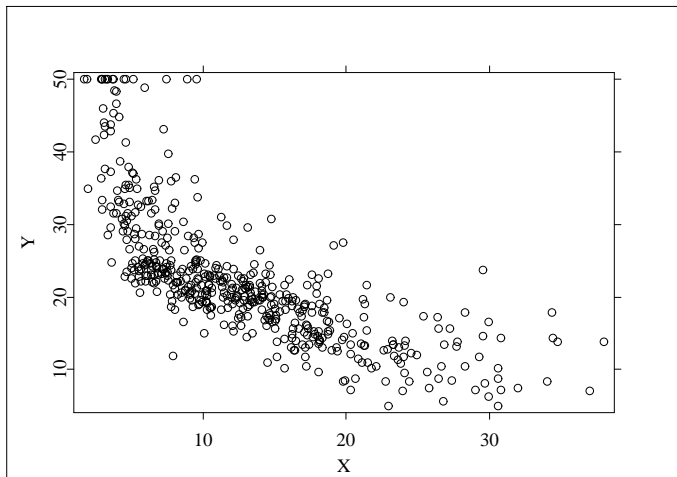


## Basic Plotting

```
1 library ("plot")           ; loads library plot
2 data = read ("bostonh")    ; reads Boston Housing data
3 x = data[,13:14]          ; selects columns 13 and 14
4 plot(x)                   ; plots data set
```



## Example plot

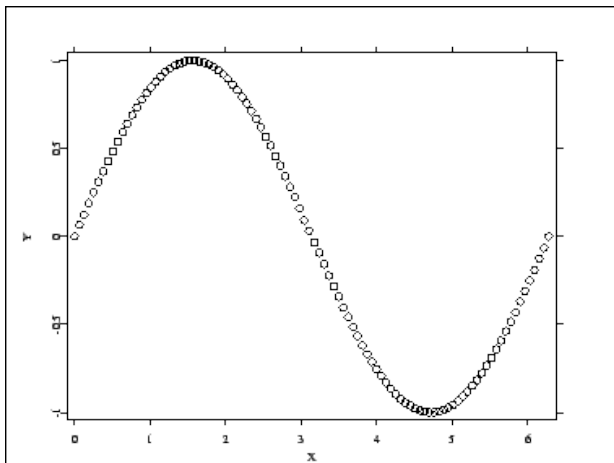


## Basic Plotting

```
1 library("plot")           ; loads library plot
2 xmin = 0                  ; grid minimum
3 xmax = 2*pi               ; grid maximum
4 n     = 100                ; number of grid
   points
5 x = xmin + (xmax-xmin)/(n-1) .* (0:n-1)
6 ; generates grid
7 y = sin(x)                ; computes sin(x)
8 plot(x~y)                 ; plots data set
```



## Example plot



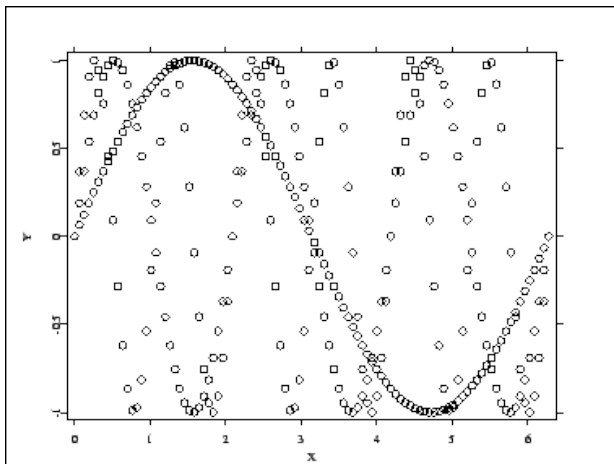
## Multiple Plots

```
1  library("plot")
2  xmin = 0
3  xmax = 2*pi
4  n     = 100
5  x    = xmin + (xmax-xmin)/(n-1) .* (0:n-1)
6  y1   = sin(x)
7  y2   = sin(3.*x)
8  y3   = sin(6.*x)
9  plot(x~y1, x~y2, x~y3)
```





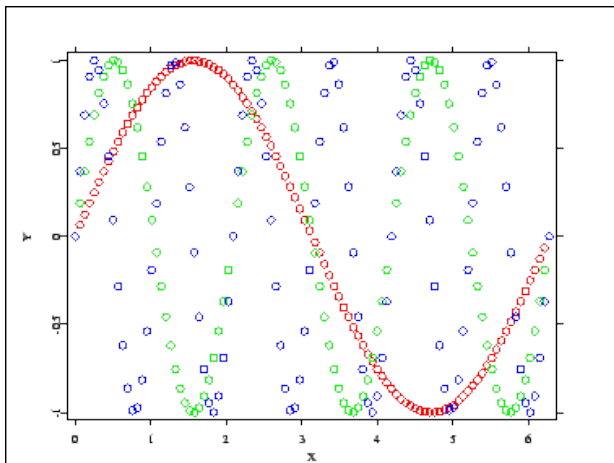
## Example plot



# Colors

```
1  library("plot")
2  xmin = 0
3  xmax = 2*pi
4  n = 100
5  x = xmin + (xmax-xmin)/(n-1) .* (0:n-1)
6  y1 = sin(x)
7  y2 = sin(3.*x)
8  y3 = sin(6.*x)
9  z1 = setmask(x~y1, "red")
10 z2 = setmask(x~y2, "green")
11 z3 = setmask(x~y3, "blue")
12 plot(z1, z2, z3)
```

## Example plot

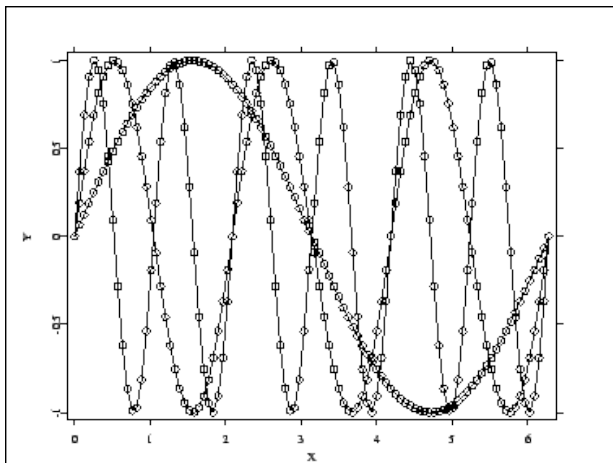


## Plotting Lines

```
1  library("plot")
2  xmin = 0
3  xmax = 2*pi
4  n     = 100
5  x    = xmin + (xmax-xmin)/(n-1) .* (0:n-1)
6  y1   = sin(x)
7  y2   = sin(3.*x)
8  y3   = sin(6.*x)
9  line(x~y1, x~y2, x~y3)
```



## Example plot

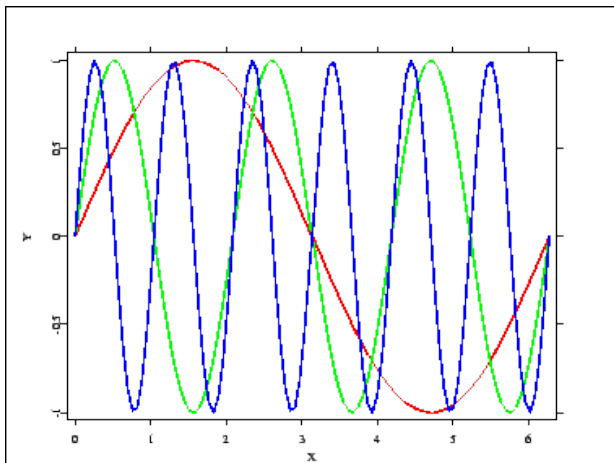


## Lines and Colors

```
1  library("plot")
2  xmin = 0
3  xmax = 2*pi
4  n     = 100
5  x    = xmin + (xmax-xmin)/(n-1) .* (0:n-1)
6  y1   = sin(x)
7  y2   = sin(3.*x)
8  y3   = sin(6.*x)
9  plot(x~y1, x~y2, x~y3)
10 z1 = setmask(x~y1, "line", "red")
11 z2 = setmask(x~y2, "line", "green")
12 z3 = setmask(x~y3, "line", "blue")
13 plot(z1, z2, z3)
```



## Example plot



## Multiple Plots

```
1 disp = createdisplay(rownum, colnum)
2
3 show(disp, <row>, <col>, what)
```

