Subseting elements:
selected_elements ← my-matrix[my-matrix >10]

- cessing elements
matrix[, 3] 3rd column | matrix[2,3] second row, 3rd volum
- transpose
t(my-matrix)
- multiplication, addition
my-matrix %*% another matix, matrix1 + matrix2
- Combine Vertically
combined_matrix ← rbind (matrix1, matrix2)
- Selecting elements based on sequence
matrix [, m[1, ]>2] # select columns for which 1it row is >2.
- diag(A), diag(n).
- Reshape matrix : dim(matrix) ← c(2,3)

### Arrays                    (Multi dimensional data structures)
- Create array : my_array ← array (1:18, dim=c(3,3,2))
- accessing elements: element ← my-array (2,3,1)
                                        row col depth.
- subarray: Sub_array ← my-array[1:2, 1:2,1]
- transpose: transposed-array ← aperm (my-array, c(3,2,1))
- dimension set: dim (array) ← c(2,3,3).
★ - names to array: dimnames(my-array) ←list(c ("Row1", "Row2", "Row3",
                                        c("col1", "col2", "col3")
- reshaping: reshaped-array ← aperm (array, c(3,2,1)) c("Depth1","Depth2"))

### Data frames

- data1 ← data[data.# age >25, ] : Select column with condition.
- columns ← data [, c('age','name')] : Select multiple columns
- sorted_data ← data [order (data.# age, decreasing = TRUE), ]
  το ξιτόριοη            , decreasing = FALSE,)]
Subset - subset_data ← subset (data, age >25)
★ - Προσθήκη στήλης: new_data ← transform (data, new_col = age * 2)

- Δύο περιορισμοί, σαν διάστημα:
selected_rows ← data[(data.$ Income <100) &
                        (data.$ Income >1000),]


### Basics

- rm (list = ls()) = καθαρισμός της R
- save image (file = "df.RData")


- Recycle R: εάν προσθέσω έναν vector με 10 στοιχεία σε ένα διάνυσμα με 5, τότε θα προστεθούν μόνο τα 5 πρώτα.
- data[-c(X,Y)] παίρνω τα μη NAN. στοιχεία
- για αρνητικά seq βάζω -0,5.


### Apply family

apply: | operate functions in specific parts of an array.
Suppose we want to find median for every column
- apply (state.x77, 2, median)
- apply (data, 2, median, c(0,25, 0,75))
    Λύση
- column_sums ← apply (data, 2, sum)
- new_data ← rbind (data, column_sums)


lapply: | applies a function to every element in a list and returns a list.
- my_list ← list (a = 1:3, b = 4:6, c = 7:9)
result ← lapply (my_list, mean)
result ← lapply (ataframe, sum)  # Same to dataframe


sapply: | returns a vector or matrix. 'S' → simplify       Αλλάζει μόνο η
- data_means ← sapply (data, means)                         εμφάνιση


✳ | apply (state.x77, 2, quantile, c(0,25, 0,75))

Σε κάθε αλλαγή ή υπολογισμε μετά ή πριν τον παρένθεση βάζουμε ","
για να προσδιορίσουμε αλλαγή σε γραμμή ή στήλες

* Διαγραφή low correlated variables

```
cors <- sapply (pred, cor, y = response)  # calculate correlations
mask <- (rank (- abs( cors )) <= 10)
```

* - Εύρεση μεγάλης συσχέτισης:  cor-values <- sapply (data, function(col),
    selected_cols <- names( cor_values[cor...]) cor(col, data $Y) >0.75)

tapply: Ομαδοποιώ και υπολογίζω συναρτήσεις σε data.frames βάση ενός
    factor. πχ.  tapply (X, f  fun)  f=factor στήλη, X=vector
    - chickwts :  weight = numeric, feed = factor.
    tapply ( chickwts $ weight, chickwts $ feed, mean, sd, length etc. )

### Αφαίρεση στηλών

1ος τρόπος:
```
remove_col <- c("X1", "X2" ...)
new_df <- df[, !(names(df) %in% remove_col)]
```

2ος τρόπος.
```
clear_data <- subset (data, select = - c("col1", "col2"))
```

### Υπολογισμός statistics με apply

```
descriptives_df <- data.frame (
 min = sapply (df, function(X), min(X, na.rm=TRUE)),
 quantile = sapply (df, function(X). quantile (X, 0,25, na.rm = TRUE)),
 median = sapply (df....., median (X, na.rm = TRUE)) )
```

### Sequences και κατανομής

- rnorm (n, mu, sigma) (generate)      rbinom (n, n, p) - rπιθ
- pnorm (X, mu, sigma) (ισούται Φ(x))   pbinom (X, n, p)
- dnorm (X, mu, sigma) (shape)         dbinom (X, n, p)
- qnorm (0.95, mu, sigma) (given 95%)   qbinom (0.95, n, p)
                                        όταν θα ψάχνουμε critical value

```
seq (from = value, to = value, by = value, length out = value)
- seq (0, 20, 4)
- seq (20, 1, -1)
- replicated_seq ← rep (1:3, times = 3)

- paste ("no", 1:5) → πέντε no με αριθμό διπλά
  paste ("no", 1:5, rep = "") χωρίς κενό στα strings.
```

### Διαστήματα Εμπιστοσύνης.

$$\bar{X} \pm Z_{1-a/2} \, 6/\sqrt{n} \quad , \quad \bar{X} - \bar{Y} \pm Z_{1-a/2} \sqrt{6_x^2/n + 6_y^2/m}$$

Παράδειγμα: έστω ένα data sample.
```
conf_interval ← t.test (sample_data) $ conf.Int   default 95%
conf_int_90 ← t.test (sample_data, conf.level = 0.90) $ conf.Int
print (conf_int)
```

Παράδειγμα2: δηλαδή η δ'ετω εχω τα μέτρα.
```
n = 175
mean ← mean (data) , sd_value ← (sqrt (sum (x - mean)^2)/(n-1))
t_value ← qt (Z_{1-a/2}, df = n-1)
ci_lower ← mean - t_value * (sd_value / sqrt (n))
ci_upper ← mean + t_value * (sd_value / sqrt (n))
print (paste (...., ci_lower, ci_upper))
```

[Άσκηση]

Έστω το UKDriverdeaths dataset. Να βρω το 95% CI για τον
μέσο των θανάτων ανα μήνα για κάθε χρονιά.

Hint: χρήση της apply ονομάζω data το set.
```
Xbars ← apply (data, 2, mean) , sds ← apply (data, 2, sd)
n ← nrow (data) ,   Z ← 1.96
standar_error ← sds / sqrt (n)
```

Για να υπολογίσω μία πιθανότητα
βάζω από qnorm (0.975) = 1,96

```
lower_CI ← Xbars - Z · Standard-error.
upper_CI ← Xbars + Z · Standard-error.
print (paste ( ...., lower_CI, upper_CI))
αλλώς  Z ← qnorm (0.975) αντί να ξέρουμε ήδη τον αριθμό
```

[Άσκηση]

Έστω ένα dataset με 1 numeric και factor. Να βρεθεί το διάστημα
95% για τους μέσους των 9 ομάδων.

1ος:
```
Xbars ← tapply (data $ scores,  data $ teams, mean)
sds ← tapply (...)   n ← tapply (....)
Z ← qnorm (0.975)
standard-error ← sds / sqrt (n)
lower_CI ← xbars - (z · st..)  Upper_CI ← xbars+ (z + st..).
print(paste( lower_CI, upper_CI))
```

2ος:
```
                   numeric  factor
result ← tapply ( ...., ......, function(x) t.test(x) $ conf.int)
print (result)
```

• Διάβασμα αρχείων:  dec = ',' εάν έχει κόμμα μέσα το κάνει τελεία

## Έλεγχοι Υποθέσεων test

```
t.test :    t.test (sample_data, mu = expected_mean)  #1 sample
            t.test (group1, group2)                    #2 sample
► t.test(x, y=null, alternative = c("two sided", "less", "greater"), mu = 0,
      paired = FALSE, var.equal = FALSE, conf.level = 0.95)
```
- Από default δεν θεωρεί το t.test ότι έχουμε ίσες διακυμάνσεις, αλλά:
```
var.test (group1, group2) ελέγχουμε με F-test.
t.test (group1, group2, var equal = TRUE)
```

```
paired = FALSE τα δείγματα είναι ανεξάρτητα
paired = TRUE είναι συσχετιζόμενα ή paired.
```

· Extract pvalue: p-value ← result$p.value

Ho : $\mu = 75$ , Hı : $\mu < 75$.
t_test_result ← t.test (sample, mu = 75, alternative = "less")
θγzoυμs a = 0.05
· if ( a > p-value) {
  cat (" reject null hypothesis)
} else {
  cat ("fail to reject null hypothesis) }

- Wilcox. test (group 1, group 2)

## Regression

### Stepwise AIC
lm ( Y ~ . , data = data)
- stepAIC (model, direction = "forward")
- stepAIC (model, direction = "backward")
- stepAIC (model, direction = "both")

| Influence.measures(model) |
| confint (fit) |

uαι zωρίs zo AIC

### Stepwise BIC
- step (model, K = log (length(x)), direction = "backward" )
- step (model, K = log (length (x)), direction = "forward" )
- . . . . . . direction = "both")

## Predictions, Plots and style
### Predictions:
confint (fit) = φτιάχνει διάστημα για τα coefficients.
- muhat ← predict (fit, interval = "confidence")
μέση τιμή των predictions.

- yhat ← predict (fit, inerval = "prediction")
y μαπέλο με διάστημα

✱ Προσδιορισμ νέων δεδομένων                    ) συνδυασμος με
επιλογίζωντας xbar=mu, sd =, και βάζοντας) apply για να το
rnorm (1,xbar, sd) → βάλε video μάνω για πολλές στηλές.

φτιάχνω ἐνα νέο data frame. ή σειρά-row και το βάζω
στα predict        apply (data, 2 mean, sd
explanatory-data ← data.frame(

ορίζω      X1 = rnorm (①) xbar[1], sds[1], .....)
τους αριθμ.
            predict (model, explanatory-data)  → το νέο ŷ.


Plots.
- Setting the layout.
     par(.... c(2,2), bg=".....")   pch=16 για fill

     Plot every explanatory vs fitted.        [fitted ← fitted(model)]

- plot (data$X1, model$residuals, main = "fitted vs predicted X1")
- abline (h=0, col= "red")


            fitted vs residuals : plot (model$residuals, fitted-values)
Διάστημα 95% για τα regression coefficients.
CI ← confint (model, level = 0.95)


✱✱ Διαγραφή ύποπτων τιμών influence ← dataframe, order → διαγραφω
                                        τα μεγαλύτερα coocks.distance
Extract R²:   summary (model) $ r.squared.
✱ New predictions με CI 95%:
     φτιάχνω τα νέα data: new-predictions ← predict (model, explanatory,
                    interval = "confidence", level = 0.95)
✱ New predictions με CI 95% για το response:
     observed-predictions ← predict (model, explanatory,
                    interval = "prediction", level = 0.95)


Ⓧ  cook ← influence $ cooks.distance
    threshold ← 4 / nrow (data) (το 4 το αλλάζω)
    influencial-obs ← which (cook > threshold)
    clean-data ← data [cook <= threshold, ]
    new-model ← lm (Y~ ., data = clean-data)

## Matrices    (two dimensional data structures)

- Construct matrix
  matrix (1:9, nrow=... , ncol=...)
  matrix (rnorm (6*4, mean=0, sd=1, nrow=... , ncol=... , byrow = True)

- Combine Vectors
  vect1 ← c(....)
  Vect 2 ← c(....)
  Vect 3 ← c(.)
  mymatrix ← c(vect1, vect2, vect3)                    Define as realikes
  matrix ← matrix (mymatrix, byrow = TRUE, nrow =3)
  mymatrix ← cbind (vect1, vect2, vect3) # 2 entries.

- Assign names to matrix
  colnames (my-matrix) ← c("var1", "var2", ....)
  rownames (my-matrix) ← c("x", "y", "z")