

Σημειώσεις για τις συναρτήσεις που χρησιμοποιήθηκαν στο εργαστήριο.

Πίνακας περιεχομένων

I) Δημιουργία δέγματος με βάση γνωστή κατανομή:.....	1
II) Έλεγχος υποθέσεων.....	2
A) T-τεστ : Σύγκριση δύο group δεδομένων.....	2
B) Z-τεστ.....	4

I) Δημιουργία δέγματος με βάση γνωστή κατανομή:

(στοχαστικής μεταβλητής με γνωστή κατανομή (παράδειγμα κανονικής κατανομής))

Βήματα :

- Δημιουργία στοχαστικής μεταβλητής που να παίρνει ισοπίθανα τιμές στο διάστημα [0,1)

=RAND()

- Δημιουργία στοχαστικής μεταβλητής με βάση την κανονική κατανομή (με δεδομένη μέση τιμή μ και τυπική απόκλιση σ)

=NORM.INV(RAND(); μ ; σ)

- Παράδειγμα δημιουργίας 9 δηγμάτων με $\mu = 200$ και $\sigma = 40$

II) Έλεγχος υποθέσεων

A) T-τεστ : Σύγκριση δύο group δεδομένων.

Ερώτημα:

Υπάρχουν στατιστικά σημαντικά στοιχεία (statistical significant evidences) για να καταρρίψουν την υπόθεση ότι ο μέσος όρος δύο δηγμάτων είναι

- ίδιος
- ο ένας δεν είναι μεγαλύτερος (ή δέν είναι μικρότερος)

Συμπληρωματικό Ερώτημα I : ένα T τεστ δύο δηγμάτων που έχουν δημιουργηθεί με βάση την ίδια κατανομή με τη πιθανότητα θα απορρίψουν την μηδενική υπόθεση (ενώ αυτή αληθεύει , αφού εμείς επιλέξαμε να τα φτιάξουμε από την ίδια).

Συμπληρωματικό Ερώτημα II : ένα T τεστ δύο δηγμάτων που έχουν δημιουργηθεί με βάση διαφορετική κατανομή με τη πιθανότητα δέν θα απορρίψουν την μηδενική υπόθεση (ενώ αυτή αληθεύει , αφού εμείς επιλέξαμε να τα φτιάξουμε από διαφορετική).

Δημιουργία 2 δηγμάτων με βάση την ίδια κατανομή (κανονική κατανομή με $\mu = 200$ και $\sigma = 40$)

Εφαρμογή T τεστ δύο δηγμάτων

Δυνατές επιλογών που έχουμε σε ένα **T τεστ** :

1.

- τα δήγματα είναι paired ή όχι
(π.χ. αντιστοιχούν στο ίδιο “άτομο” υπό δύο διαφορετικές συνθήκες ή σε διαφορετικά “άτομα” ?)
- Αναμένω τα δύο δήγματα να έχουν την ίδια διακύμανση
- Αναμένω τα δύο δήγματα να έχουν την διαφορετική διακύμανση (αν δεν το ξέρω επιλέγω συνήθως αυτό)

2. το τέστ είναι one or two tail ?

η Μηδενική υπόθεση H_0 αφορά

α) (two tails) ισότητα μέσων όρων $\mu_x =$ (άρα η εναλλακτική υπόθεση H_1 σημαίνει ότι είτε $\mu_x >$ είτε $\mu_x <$)

β) (one tail) ανισότητα μέσων όρων $\mu_x >$, ή $\mu_x <$ (ένα από τα δύο όχι είτε το ένα είτε το άλλο)

Σε περίπτωση που πρόκριτε για paired, το libreoffice έχει αυτόματο υπολογισμό μέσω των επιλογών

Δεδομένα → Στατιστικά → Δοκιμή T με ζεύγη.

Όπου αναφέρονται όχι μόνο το p-value αλλά και οι κρίσιμες τιμές.

Τα T-τεστ μπορούν να υπολογιστούν με την χρήση μίας συνάρτησης.
=TTEST(**Array1** ; **Array2** ; **Tails** ; **Type**)

όπου

Array1 ; τα κελιά με τα δεδομένα του πρώτου group

Array2 ; τα κελιά με τα δεδομένα του δεύτερου group

Tails ; 1 ή 2 (one tail ή two tails)

Type :

1 τα δῆγματα είναι paired

2 τα δῆγματα αναμένετε να ἔχουν την ίδια διακύμανση

3 δῆγματα δεν αναμένετε να ἔχουν την ίδια διακύμανση (ή απαλά δεν ξέρω)

Παράδειγμα

p value :

paired , one tail

θεωρώντας ίδια διακῆμανση, one tail

μή θεωρώντας ίδια διακῆμανση, one tail

(η διαφορά του να χρησιμοποιήσεις 2 tails αντί για 1 είναι ότι διπλασιάζετε η τιμή του p δηλαδή :

=TTEST(**Array1** ;**Array2** ; **2** ; **Type**) = 2* TTEST(**Array1** ;**Array2** ; **1** ; **Type**)

έτσι ώστε το p να συγκρίνετε πάντα με το επίπεδο σημαντικότητας α (π.χ. το 0.05) και να μην χρειάζετε να χωρίσεις το α σε αριστερό και δεξιό μέρος.)

Link για διαφορές του ttest

B) Z-τεστ

A

Δεδομένα : ένας μέσος όρος $\langle X \rangle$, και η διακύμανση της κατανομής σ_x .

Σύγκριση του $\langle X \rangle$ με την υπόθεση μας για μ_x

(μ_x η μέση κατανομή του πληθυσμού π.χ. τι θα παίρναμε εάν μπορούσαμε να κάνουμε άπειρα πειράματα ελέγχοντας όλα τα ενδεχόμενα)

Ερώτημα : Ο μέσος όρος που παρατήρησα $\langle X \rangle$ στις μετρήσεις μου υποστηρίζουν ή εναντιώνετε στην υπόθεση ότι η μέση τιμή του πληθυσμού είναι μ_x (είτε γνωρίζοντάς το σ_x είτε εκτιμώντας το με βάση την διακύμανση του δείγματος (link).

Παραδείγματα στο αρχείο Hypothesis_testing_examples_1.xlsx στο eclass.

B

Δεδομένα : Ένα σέτ μετρήσεων, (και όχι υποχρεωτικά η διακύμανση της κατανομής σ_x , μιας και μπορεί να εκτιμηθεί από το σέτ των μετρήσεων)

Ερώτημα : Το σέτ μετρήσεων μας προσφέρουν στατιστικά σημαντικά στοιχεία για να απορίψουμε την υπόθεση ότι η μέση τιμή του πληθυσμού είναι μ_x (το όπου το σ_x μπορεί να εκτιμάτε το με βάση την διακύμανση του δείγματος ή να δοθεί σαν δεδομέν(link).

=ztest(**Array1**, μ_x , σ_x)

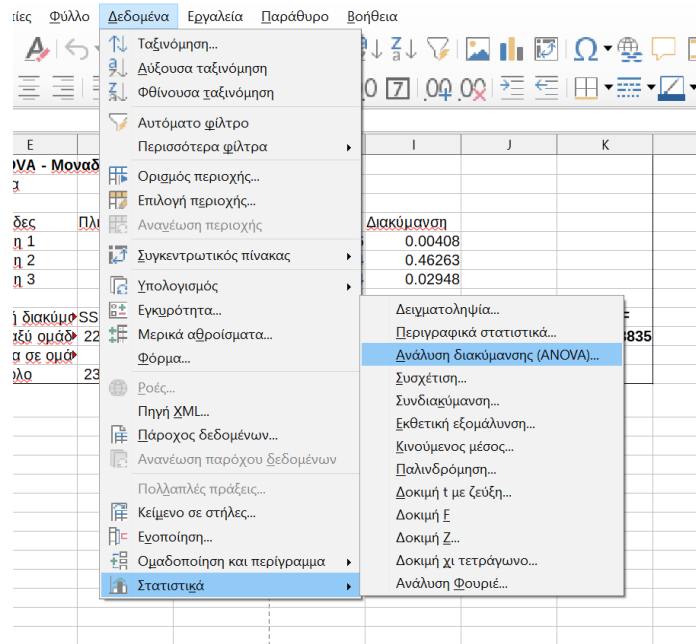
Zτεστ θεωρεί κανονική κατανομή και κατά συνέπεια θα πρέπει να εφαρμόζετε για μέγεθος δείγματος μεγαλύτερου του 30.

Παραδείγματα στο αρχείο Hypothesis_testing_examples_1.xlsx στο eclass.

Άλλα παραδείγματα χρήσης στατιστικών τεστ στο LibreOffice .

Anova Analysis Of Variance

Εφαρμογή χρησιμοποιώντας τις επιλογές :



Αποτελέσματα :

P value , το F είναι μεγαλύτερο ή μικρότερο από το Fcrit

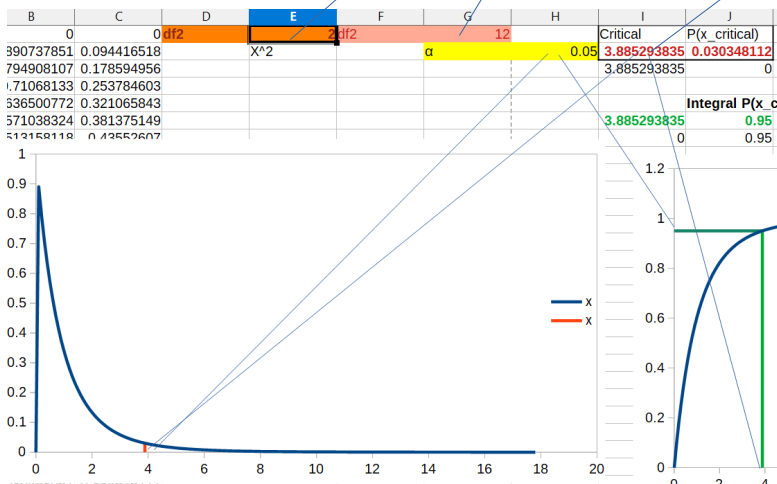
Το Fcrit εξαρτ'ατε από 2 βαθμούς ελευθερίας

τον αριθμό των γκρουπς - 1

και τον αριθμό των μετρίσεων - τον αριθμό των γκρουπς.

Αναλυτικό παράδειγμα υπολογισμού ανονα στο αρχείο :anova_pH (2019_04_12 12_59_09 UTC)

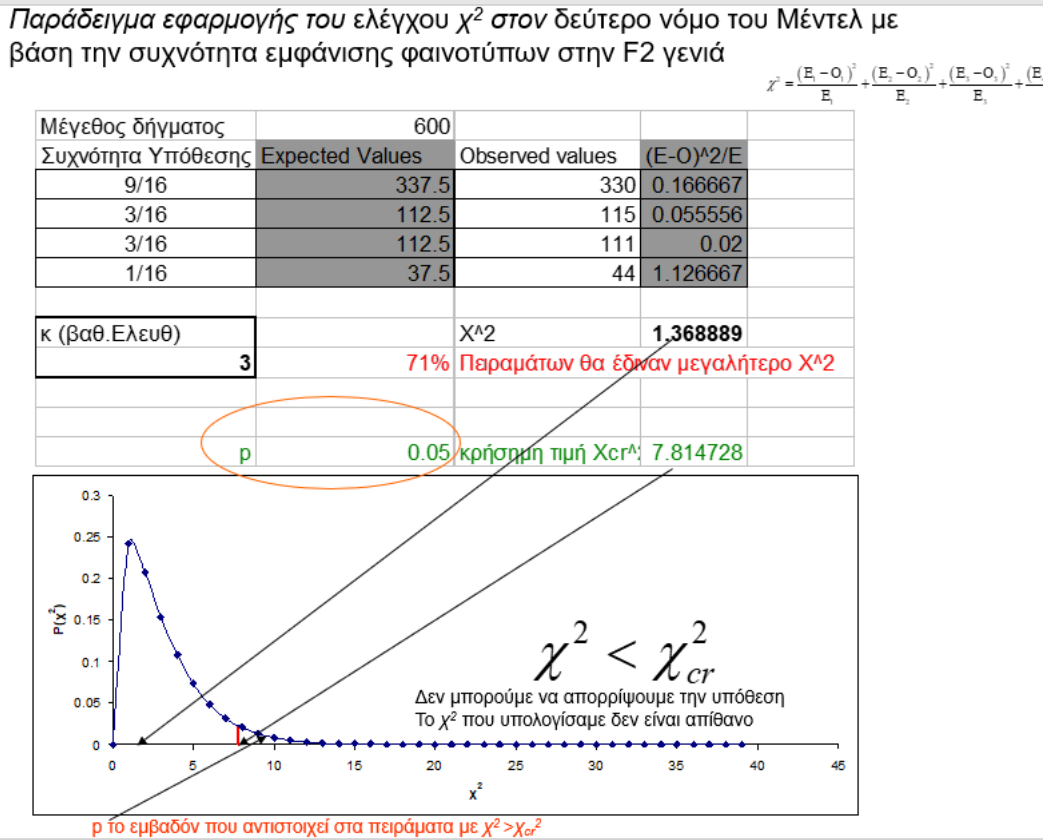
Μέτρηση 1	Μέτρηση 2	Μέτρηση 3	ANOVA - Μοναδικός παράγοντας						
1.16	3.57	1.16	Άλφα	0.05					
1.24	3.57	1.15	Ομάδες	Πλήθος	Άθροισμα	Μέσος	Διακύμανση		
1.06	4.04	1.06	Στήλη 1	5	5.78	1.156	0.00408		
1.16	4.74	1.3	Στήλη 2	5	18.82	3.764	0.46263		
1.16	2.9	1.5	Στήλη 3	5	6.17	1.234	0.02948		
			Πηγή διακύμ	SS	df	MS	F	Τιμή-P	Κρίσιμο F
			Μεταξύ ομάδ	22.01441333	2	11.00720667	66.5503537	3.19941E-07	3.885293835
			Μέσα σε ομά	1.98476	12	0.165396667			
			Σύνολο	23.99917333	14				



Τέστ χ^2

Παράδειγμα από Γενετική.

1) Σύγκριση προβλεπόμενης αναλογίας φαινοτύπων με βάση τον δεύτερο νόμο του Μεντελ για πεπερασμένο δῆγμα.



Επισήμανση της διαφοράς σε περίπτωση που χρησιμοποιήτε το τεστ χ^2 για τον έλεγχο γενετικής ισοροπίας : οι βαθμοί ελευθερίας δεν είναι ο αριθμός των κατηγοριών-1 όπως σε αυτό το παράδειγμα αλλά ίσος με τον αριθμό των παραμέτρων του μοντέλου γενετικής ισοροπίας (λυγότεροι)

Κατανομές που χρησιμοποιούνται στα τεστ υποθέσεων και η σημασία των βαθμών ελευθερίας.

Στο Αρχείο distributions_df.xlsx

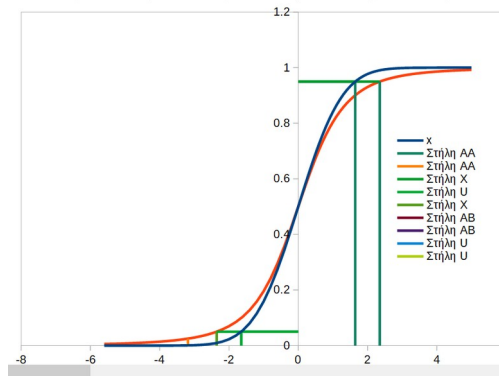
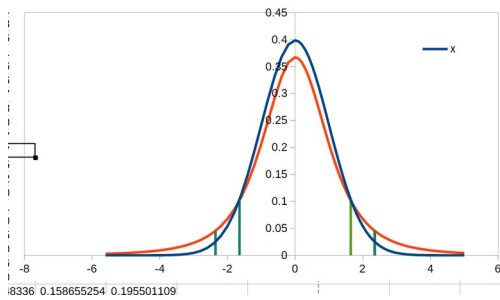
θα βρείτε παραδείγματα όπου δείχνουν την Κατανομές που χρησιμοποιούνται στα τεστ που αναφέραμε (μαζί με τα διαγράμματα των αθριστικών κατανομών που μας δίνουν τα εμβαδά που εμπλέκονται στους υπολογισμούς του επιπέδου σημαντικότητας σε σχέση με τις κρίσιμες τιμές και το p value):

- a) κανονική κατανομή (με $\mu_x=0$, $\sigma_x=1$)
- b) κατανομή t (σαν συνάρτηση των βαθμών ελευθερίας)
- c) κατανομή χ^2 (σαν συνάρτηση των βαθμών ελευθερίας)
- d) κατανομή F που χρησιμοποιείτε στο ANOVA (σαν συνάρτηση 2 βαθμών ελευθερίας)

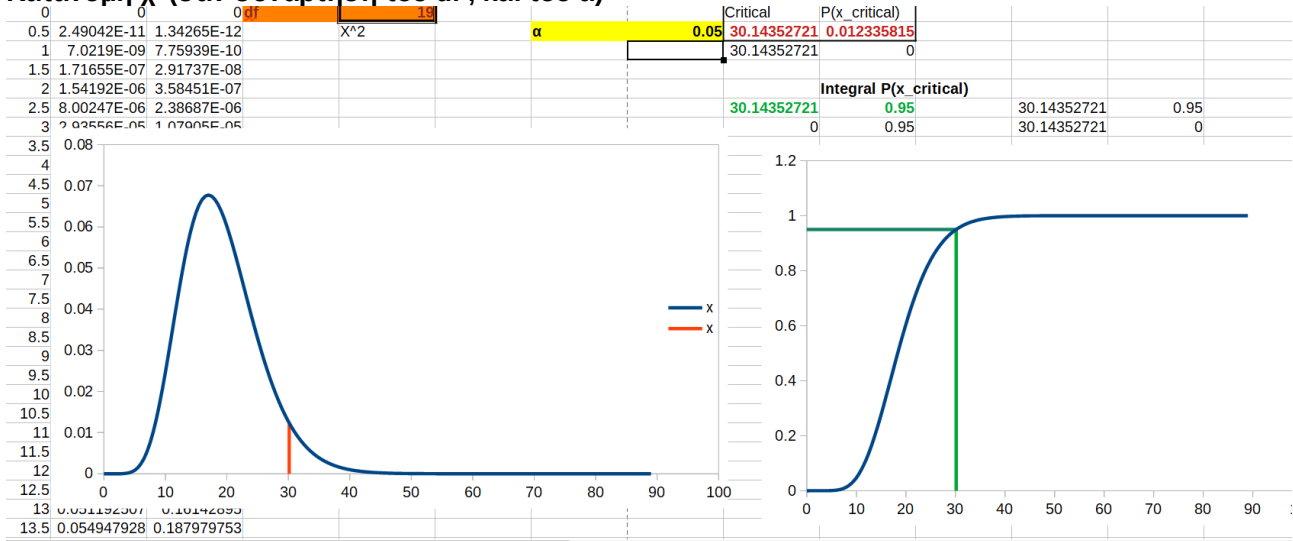
Ερωτήματα που μπορούν να κατανοηθούν εξετάζοντας τα διαγράμματα του αρχείου.

- 1) Πόσο διαφέρει η κατανομή t από την κανονική, και πως όταν έχουμε μεγαλύτερο δείγμα (περισσότερους βαθμούς ελευθερίας) οι δύο κατανομές πλησιάζουν)
- 2) Τι αλλάζει όταν αλλάζουμε το διάστημα εμπιστοσύνης.
- 3) Πως αλλάζουν οι κρίσιμες τιμές των τεστ όταν αλλάζουν οι βαθμοί ελευθερίας.

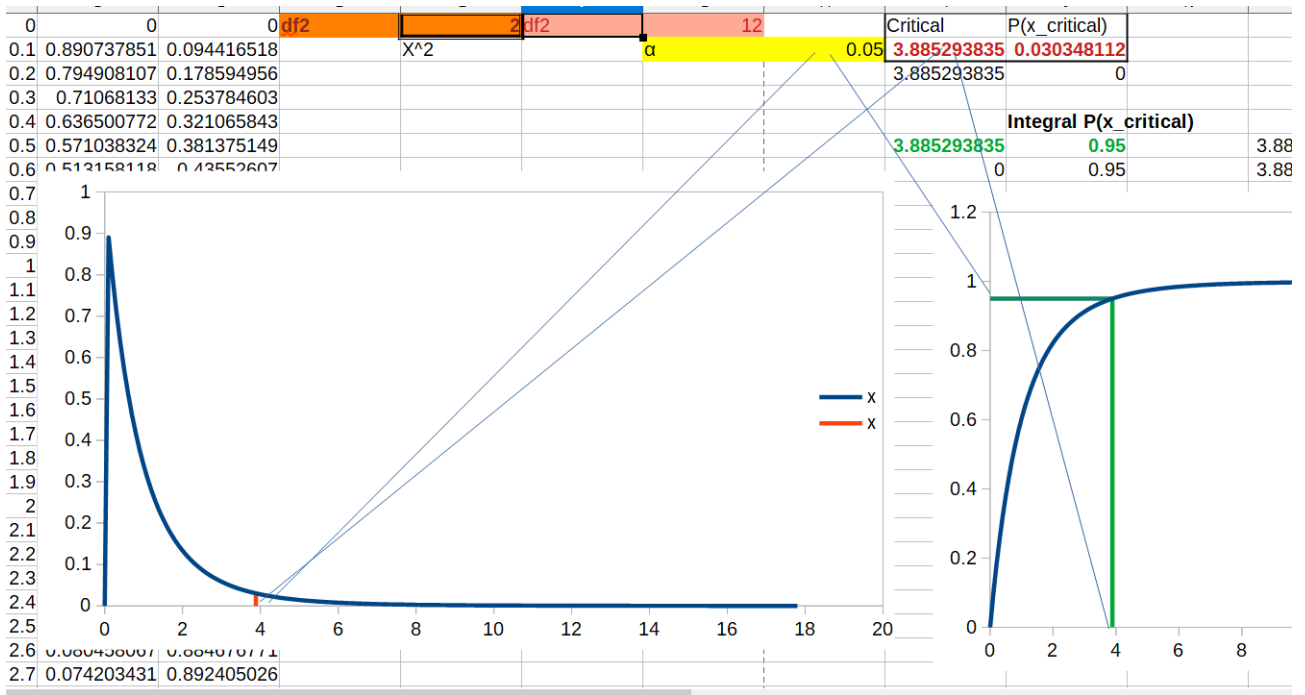
T-κατανομή με 3 βαθμούς ελευθερίας (κελί H2 στο αρχείο distributions_df.xlsx) και επίπεδα εμπιστοσύνης $\alpha = 0.05$ (κελί K2 στο αρχείο distributions_df.xlsx)



Κατανομή χ^2 (σαν συνάρτηση του df , και του α)



Κατανομή F.



Μαθηματικό υπόβαθρο στην διαφορά του υπολογισμού της διακύμανσης, και των βαθμών ελευθερίας στα είδη των τεστ.

$$\mu_x = \mu_y \rightarrow t \rightarrow \in N(0, \sigma_t)$$

Type :

1 τα δέγματα είναι paired

$$\mu_x = \mu_y \rightarrow \bar{x} \approx \bar{y} \rightarrow \bar{x} - \bar{y} \rightarrow \in Norm(0, \sigma_{\bar{x} - \bar{y}})$$

$$\sigma_{\bar{x} - \bar{y}}^2 \rightarrow \sigma_x^2 + \sigma_y^2, \text{ αν } x \text{ και } y \text{ είναι ασυσχέτιστα}$$

$$\sigma_y^2 \rightarrow \sigma_y^2 / \sqrt{(M)}$$

$$\sigma_x^2 \rightarrow \sigma_x^2 / \sqrt{(N)},$$

$$\sigma_x^2 / \sqrt{(N)} \rightarrow S_x^2 = \frac{\sum (x_i - \bar{x})^2}{(N - 1)} \sigma_y^2 / \sqrt{(M)} \rightarrow S_y^2 = \frac{\sum (y_i - \bar{y})^2}{(M - 1)}$$

,

$$t = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{S_x^2}{N} + \frac{S_y^2}{M}}} \rightarrow \bar{x} - \bar{y} \rightarrow \in N\left(\mu_x - \mu_y, \sqrt{\frac{S_x^2}{N} + \frac{S_y^2}{M}}\right)$$

$$df = \frac{2}{\left(\frac{1}{N} + \frac{1}{M}\right)}$$

Βαθμοί ελευθερίας : $(N+M)/2-1$ ή

2 τα δέγματα αναμένετε να έχουν την ίδια διακύμανση

$$\mu_x = \mu_y \rightarrow \bar{x} \approx \bar{y} \rightarrow \bar{x} - \bar{y} \rightarrow \in Norm(0, \sigma_p)$$

$$t = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{s_x^2}{N} + \frac{s_y^2}{M}}} \rightarrow \frac{\bar{x} - \bar{y}}{s_p \sqrt{\frac{1}{N} + \frac{1}{M}}} \rightarrow \bar{x} - \bar{y} \rightarrow \in N\left(\mu_x - \mu_y, s_p \sqrt{\frac{1}{N} + \frac{1}{M}}\right)$$

$$s_p^2 = \frac{(N-1) * s_x^2 + (M-1) * s_y^2}{N+M-2}$$

$$DF = N+M-2$$

3 unpaired και τα δέγματα δεν αναμένετε να έχουν την ίδια διακύμανση

$$\mu_x = \mu_y \rightarrow \bar{x} \approx \bar{y} \rightarrow \bar{x} - \bar{y} \rightarrow \in Norm(\mu_x - \mu_y, \sigma_{\bar{x}-\bar{y}})$$

$$\sigma_y^2 \rightarrow \sigma_y^2 / \sqrt{(M)}$$

$$\sigma_x^2 \rightarrow \sigma_x^2 / \sqrt{(N)},$$

$$\sigma_x^2 / \sqrt{(N)} \rightarrow S_x^2 = \frac{\sum (x_i - \bar{x})^2}{(N-1)} \sigma_y^2 / \sqrt{(M)} \rightarrow S_y^2 = \frac{\sum (y_i - \bar{y})^2}{(M-1)}$$

,

$$t = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{s_x^2}{N} + \frac{s_y^2}{M}}} \rightarrow \frac{\bar{x} - \bar{y}}{s_p \sqrt{\frac{1}{N} + \frac{1}{M}}} \rightarrow \bar{x} - \bar{y} \rightarrow \in Norm\left(\mu_x - \mu_y, s_p \sqrt{\frac{1}{N} + \frac{1}{M}}\right)$$

$$s_p^2 = \frac{(N-1) * s_x^2 + (M-1) * s_y^2}{N+M-2}$$

$$DF = N+M-2$$