

LECTURE 1

C10/02

- Introduce myself

- Παρουσίαση Μαθήματος. Τι είναι το μάθημα;

↳ Εφαρμοσμένη Στατιστική (επί προσδιορ: εφαρμοσμένη)

⇒ Εφαρμοχές σε πραγματικά προβλήματα από μια πληθώρα διαφορετικών πεδίων επιστημών

* Χρήση υπολογιστικών στατιστικών πακέτων (ΜΑΤΛΑΒ, SPSS, MINITAB, R ...)

Εργαστήριο: ΜΑΤΛΑΒ

* Μάθημα: Αρχικά 4 ώρες θεωρία → Μετά 4 ώρες θεωρία + 2 lab
↳ Μετά 2 ώρες θεωρία + 4 lab.

* Εξέταση: Εργαστηριακές Ασκήσεις (~ 3) μέγιστη
+ Τελική Εξέταση (ή Τελικό Project)

* Φοιητές: Applied Math, Math, CSB

Q → ποσοί;

- Προαπαιτούμενα: Πιθανότητες, Στατιστική.

Q: Τι έχετε κάνει;

⇒ Πιθανότητες: χώροι πιθανότητας, ανεξαρτησία, τύπος του Bayes, Διασπορά, Κατανομές (διωνομική, Bernoulli, Poisson) νόμος των μεγάλων αριθμών

⇒

⇒ Στατιστική : Παραμετρικά στατιστικά μοντέλα
εξαρτητικά δείγματα
Παραμετρικοί χώροι
Διαστήματα εμπιστοσύνης
Προβλήματα ελέγχου κανονικών πληθυσμών

- Q. . Από που μπορεί να προέρχονται τα φυσικά προβλήματα;

⇒ Φυσικές επιστήμες (Φυσική, Χημεία, Βιολογία, ...)

⇒ Ιατρική

⇒ Οικονομικές επιστήμες

⇒ Κοινωνικές

→ Παιδαγωγικά θέματα

⇒ - - - -

Προσοχή : Η εφαρμοσμένη στατιστική σημαίνει πάντα την ανάλυση δεδομένων

Applied statistics \equiv Στατιστική Ανάλυση Δεδομένων

Άλλο παράδειγμα : Κατανομή ύψους . Έστω δεδομένα ύψους N .
Βρείτε την κατανομή πιθανότητας.

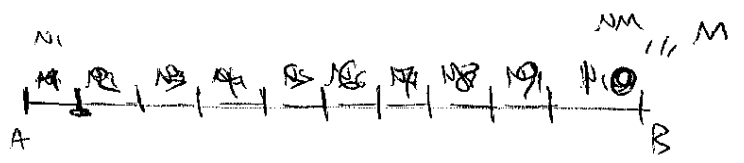
\hookrightarrow Μεταβλητή X : ύψος

\hookrightarrow Πιθανότητα $P(X)$

Q : Έστω δείγμα N . πως ρίχνουμε την πιθανότητα. $P(X \in \dots)$

$$P(X) = \frac{N(\dots)}{N}$$

- Ορίζουμε ομα A, B ,
διαστήμα h

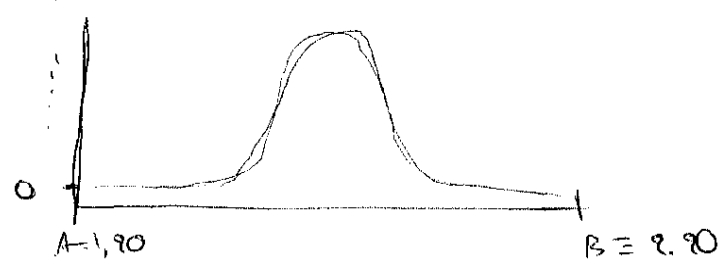


$$h = \frac{B-A}{M} \quad \text{ε.ρ.} \quad A = 1.20 \quad , \quad B = 2.00$$

$$P_1 = P(A+h) = \frac{N_1}{N}$$

$$\sum_{i=1}^M P_i = \frac{N}{N} = 1$$

- κατανομή πιθανότητας



$$P_1 = P_0 \pm \underline{\underline{\epsilon_1}}$$

- χρειάζεται : Αιτιολογία του προβλήματος

- Q : επίδραση N . πόσο είναι το σφάλμα ; $\epsilon_r = \epsilon_r(N)$
 $= ?$

π.κ. Αν $N=1000$. τι μπορούμε να πούμε για $\bar{N} = 10 \text{ cm} \pm 0.001$

Q : Στην προ-μαθητικότητα η κατανομή ύψους ^{συνήθως} δεν είναι κανονική! είναι περίπου κανονική.

Q : Είναι το πρόβλημα καλά ορισμένο;

ΟΧΙ! Γιατί? Δεν διευκρινίζει!

LECTURE 2

- Print → share syllabus

- Introduce for the rest

→ Ευζήτηση γλως.

Εισαγωγικά (Εισήγηση 1)

* Ορολογία : Δεδομένα (data)

Δειγματοληψία (sampling) : συλλογή δεδομένων

Περιγραφική Στατιστική (descriptive statistics) : Περιγραφή δεδομένων

Στατιστική Συμπερασματολογία (statistical inference) : αιτιολόγηση συμπερασμάτων

↳ Δεδομένα συλλέγονται από πληθυσμό (population)

Δείγμα (sample) : υποσύνολο του πληθυσμού

Τυχόν μεταβλητή (random variable) : οποιαδήποτε ^{καταγεγραμμένη} ποσότητα του οποίου η τιμή αλλάζει μέσα στο στοιχείο του πληθυσμού.

↳ Ποιοτική (qualitative) : είναι κατηγορίες

Ποσοτική (quantitative) : παίρνει αριθμ. τιμές

Συνεχής (continuous) : τιμές σε ένα συνεχές διάστημα

Διακριτή (discrete) : τιμές σε ένα σύνολο διακριτών τιμών

- Παραδείγματα :

Δείγμα : Πληθυσμός

Σ.Μ. : ύψος, ^{ποιοτή} φύλο, ^{ποσοτικές} βάρος, ηλικία

→ 18 ετών

→ " " " " X

Συνεχ. / Διακριτή Σ.Μ. ⇒

* Περιγραφή στατιστική

- Συχνότητες (frequency) : Έστω δείγμα x_1, x_2, \dots, x_n με πιθανές τιμές (διακριτές)
 F a_1, a_2, \dots, a_m

F_i : πόσες φορές εμφανίζεται η τιμή a_i στο δείγμα

- Εξέταση συχνότητας (relative frequency)

Ποσοστό (percent) $P_i = \frac{F_i}{n}$

- Για κατηγορίες που μπορούν να μπουν σε διάταξη ορίζεται και η

Αθροιστική συχνότητα $\tilde{F}_i = \sum_{j=1}^i F_j$ $a < i$

Αθροιστικό ποσοστό $\tilde{P}_i = \frac{\tilde{F}_i}{n}$ $a < i$

- Παρουσίαση : πίνακα συχνότητας (frequency table), βαρσογράμμο (bar chart), κυκλικό διάγραμμα (pie chart)

Παράδειγμα : Βαρος πλυσίματος

Εξέταση Πρωτοβάθμια / μήνα

1% 2% 3% 1% 4% 5% 6% 3% 8% 2% 10% 15%

$a_1 = 1\%$ $F_1 = \frac{1}{19} = 1/19$

$a_2 = 2\%$ $F_2 = \frac{2}{19} = 2/19$

$a_3 = 3\%$ $F_3 = 1/6$

$F_3 = 1/6 + 1/6 + 1/6 = 3/6 = 1/2$

; ποσοστό στο το κός είναι
 ως 3% είναι 1/2

Πίνακας συχνότητας

a	F	F_i
1%	1/19	1/19
2%	1/9	2/19
3%	1/6	3/19

↳ Ενδιαφέρον $\binom{n}{r} = \frac{n!}{(n-r)! r!}$

- Μέτρα Στατιστικών Δεδομένων

↳ Μέτρα Κέντρου (Statistics of Location)

⇒ Μέσος Όρος (Mean, Average) $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$

Γεωμετρικός $GM_x = \sqrt[n]{\prod_{i=1}^n x_i}$

Θ: Πότε χρησιμοποιείται ο γεωμετρικός μέσος όρος;

A: (ωικιπαιδία): όταν χρειάζεται ο υπολογισμός ^{ή των ποσοτήτων} "αποτελεσμάτων"

π.χ. Χρηματιστήριο με αύξηση 10% στο 1^ο χρόνο, 90% στο 2^ο και πτώση 15% στο 3^ο.

Τότε $GM_x = \frac{1,1 \cdot 0,9 \cdot 0,85}{3} = 1,0391$: μέση αύξηση καφέτας!

Αρμονικός (harmonic) $\frac{1}{H_x} = \frac{1}{n} \sum_{i=1}^n \frac{1}{x_i}$

⇒ $H_x = \frac{n}{\sum_{i=1}^n \frac{1}{x_i}}$

Ω: Πότε χρησιμοποιείται ο αρμονικός μέσος όρος;

A (ωικιπαιδία): π.χ. σε ένα ταξίδι έκανε στη μέση διαδρομή ταχύτητα 40 km/h και στην άλλη με 60 km/h → μέση ταχύτητα = H_x

$= \frac{2}{\frac{1}{40} + \frac{1}{60}} = 48 \text{ km/h}$

⇒ Δειγματική Διαμέσος (Median)

: κεντρική τιμή όταν διατάξουμε τα δεδομένα κατά αύξουσα σειρά

n : περιτός $\tilde{x} = x_{(n+1)/2}$

n : άρτιος $\tilde{x} = \frac{x_{(n/2)} + x_{(n/2+1)}}{2}$

⇒ Percentile (Δεκακταία Ποσοστά)

: Ποσοστό σημείων που βρίσκονται πριν από μια τιμή αναφοράς.

Διατάξτε κατά αύξουσα σειρά

η k percentile $\rightarrow \frac{(n+1) \cdot k}{100}$
 η τιμή της σειράς \rightarrow

↳ Μέτρα Μεταβλητότητας (Statistics of dispersion)

: Διαφορετικά δείγματα μπορεί να έχουν το ίδιο μέσο αλλά διαφορετική διασπορά

⇒

~~Εύρος~~ $(= x_{\max} - x_{\min})$
 Εύρος (Range)

⇒ Διακύμανση ή Διασπορά (Variance)

πληθυσμίο $\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$

Δείγμα $s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

⇒

(Q: Ποια η Διασπορά?)

⇒ Τυπική απόκλιση (standard deviation)

$$s = \sqrt{s^2} \quad , \quad \sigma = \sqrt{\sigma^2}$$

Ευρεθεί τιμή τυπικής απόκλισης $V = \frac{\sigma \times 100}{\bar{x}} = \frac{s \times 100}{\bar{x}}$

Συνολικά
 ⇒ Εφαρμογή (Standard Error)

$$S_{\bar{x}} = \frac{s}{\sqrt{n}}, \quad \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}}$$

⇒ Ενδο τεταρτομοριακό εύρος F (Interquartile Range)

$$F = Q_3 - Q_1$$

" " " "
 75-εταροστιακό σημείο 25-εταροστιακό σημείο
 σημείο

p-εταροστιακό σημείο είναι ποσο του ατόμου είναι μικρότερος από αυτή τη παρατήρηση.

⇒ Εύρος 5 αριθμών (5 number summary) : Διάμετρος, Q_1 , Q_3
 x_{min} , x_{max}

⇒ Στατιστική ποσότητα (Moment Statistic)

1η χαρακτηριστική ποσότητα $\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$

2η $\Rightarrow \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x}) = 0$

3η $\Rightarrow \equiv s^2$

3η, 4η, 5η ως προς να υπολογιστούν ασυμμετρίας (skewness) και κούρτωσης

→ Μισα ποσότητας → έχει ανοχή?

LECTURE 3 (22/02)

Είδη : Βασικές έννοιες
Ορολογία
Μερισμοί και στατιστικά

- Λίστα Γραμμάτων, Μικράς
ΕΜ 272

ΚΑΤΑΝΟΜΕΣ (Distributions)

* Διακριτές (Discrete)

↳ Διωνυμική (Binomial) : Έστω δείγμα με 2 δυνατά αποτελέσματα, τιμές A, B

Πιθανότητες $P(X=x) = \binom{n}{x} p^x q^{n-x}$

$P(A) = p$
 $P(B) = q$

Παράδειγμα : σύνολο φρούτων (μήλο (A), μπανάνα (B))

1 φρούτο : { A, B }
2 φρούτα : { AA, AB, BB }

Πιθανότητες : { p, q, q } $(p+q)^n$

$$P(X=x) = \frac{n!}{x!(n-x)!} p^x (1-p)^{n-x}$$

μέσος $\mu = np$ $\sigma = \sqrt{npq}$ $\sigma^2 = npq$

↳ Poisson : καλή προσέγγιση της διωνυμικής όταν το ένα γεγονός είναι σπάνιο (π.χ. $p < 0.1$) και το μέγεθος του δείγματος μεγάλο (n >> 1)

X : Διωνυμική με μέγεθος n και μίγμα p

$$P(X=x) = \frac{e^{-\lambda} \lambda^x}{x!}$$

$\lambda = np$
 $\lambda = \mu = \sigma^2$

εξήγησε το!

⇒ Παράδειγμα : (Birthday Paradox)

Έστω n άνθρωποι. Υπάρχουν $\binom{n}{2}$ ζευγάρια. Ορίζουμε ως επιτυχία αν υπάρχει ζευγάρι το οποίο έχει κοινή ημερομηνία γεννησης

πιθανότητα $1/365 = P(B|A)$

Αριθμός επιτυχιών $\mu = nP = \frac{n(n-1)}{2} \cdot \frac{1}{365} = \frac{n(n-1)}{730}$

πιθανότητα ότι δεν υπάρχουν 2 άτομα με ίδια ημερομηνία γεννησης

$P(X=0) = e^{-\lambda} \frac{\lambda^0}{0!} = e^{-\lambda} = \exp\left(-\frac{n(n-1)}{730}\right)$

Μή ώστε η πιθανότητα να είναι $< 0,5$: $\exp\left(-\frac{n(n-1)}{730}\right) \leq 0,5$

$\Rightarrow \exp\left\{\frac{n(n-1)}{730}\right\} \geq 2 \rightarrow n(n-1) \geq 730 \ln(2) \Rightarrow n \approx 28$

- Q: Ποια είναι η πιθανότητα αν $n=70$, $n=500$?
- Q: Ποια είναι η πιθανότητα κάποιος να έχει την ίδια ημερομηνία με τη δική μας?

* ΣΥΝΕΧΙΣ

↳ Ομοιομορφία (Uniform)

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{if } a < x < b \\ 0 & \text{otherwise} \end{cases}$$

$$\mu = \frac{a+b}{2} \quad \sigma^2 = \frac{(b-a)^2}{12}$$

↳ Κανονική (Normal)

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

: συμμετρική δόξα από κέντρο

$\mu \pm \sigma$	περιέχει	68.3%	σημ.
$\pm 2\sigma$	"	95.5%	"
$\pm 3\sigma$	"	99.7%	"

↳ Γάμμα

: χρόνος έως ότου η σχετικό πρόβλημα ολοκληρωθεί

$$f(x) = \frac{\lambda e^{-\lambda x} (\lambda x)^{n-1}}{\Gamma(n)} \quad \text{for } x \geq 0$$

$$\Gamma(n) = \int_0^{\infty} e^{-x} x^{n-1} dx = (n-1)! \quad \text{για ακέραιες τιμές του } n$$

↳ Έξθεση

: χρόνος έως ότου ένα γεγονός συμβεί

(συνέχεια ανάδοχο της Γάμμα με $n=1$)

$$f(x) = \lambda e^{-\lambda x} \quad \text{for } x \geq 0 \quad \mu = \frac{1}{\lambda}, \quad \sigma^2 = \frac{1}{\lambda^2}$$

$$F(x) = 1 - e^{-\lambda x}$$

$$P(X > t+s | X > t) = P(X > s) \quad \text{for } t \geq 0 : \text{Δεν έχει μνήμη!}$$

⇒ (Μι-σφιδαρ χ²) : τετραγωνο βαθμους στο n-δελσταο κίρο

Αν Y_i είναι κανονικές κατανομή

Σ_{i=1}ⁿ Y_i² ακολουθεί χ² κατανομή

↳ t η student κατανομή

$\bar{X} - \mu$ ένας δείκτης είναι κανονική κατανομή

$\frac{\bar{X} - \mu}{\sigma}$ >> >> μ=0, σ=1

αλλά $\frac{\bar{Y}_i - \mu}{\sum Y_i}$: πιο "πλάγια" από την κανονική κατανομή.
των
βγαίνει

Εκτίμηση - Διαστήματα εμπιστοσύνης

Η κατανομή μιας τ.μ. χαρακτηρίζεται από μ : μέση τιμή
 σ^2 : διασπορά

Εκτίμηση μιας παραμέτρου : σημειακή εκτίμηση (φραντ βελιτισμο.)

> Διαστήματος τιμής $(1-\alpha)\%$ confidence (interval) :

α : (σταθμική) σημαντικότητα της δοκιμής (significance level)

π.χ. Εκτίμηση μέσης τιμής μ - χρησιμοποιείται να γνωρίζουμε \bar{x} ,
 σ^2 , μορφή της κατανομής

κανονική Διαστήματα εμπιστοσύνης : $\bar{X} \pm (χρ. τιμή τ.μ.) \cdot \frac{\sigma}{\sqrt{n}}$
τιμής σφάλματος

$$\bar{X} \pm z_{1-\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

- Εμπριέχει την αληθινή τιμή της παραμέτρου με πιθανότητα $(1-\alpha)\%$

- Για τα δείγματα $(n < 30)$: δεν μπορεί να οριστεί παραμετρικά
 μη κανονικής κατανομής

χρησιμοποιείται μη-παραμετρική μέθοδος

- Διαφορά μέσων τιμών

για κανονικές $(n_1, n_2 > 30)$

$$(\bar{X}_1 - \bar{X}_2) \pm z_{1-\alpha/2} \cdot \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$$

LECTURE 4 (26/02/2010)

Είδαμε : κατανομές
 Διαστήματα εμπιστοσύνης

Confidence Intervals

Ορίζεται

→ Δίστα : em 272
 email : masorodomo@stat.uoi.gr
 subscribe : subscribe em 272
 - Lab : Τρίτη ή Πέμπτη
 15:00-18:00.
 Παρασκευή : 15:00-15:00 Η-205
 Δευτέρα : ΕΓΓΡΑΦΟ γυμνασίου
 ΕΡΓ : SIS : ΔΕΝ ΕΧΑΝ ΚΑΝΕΙ ΚΑΛΩΣΗΛΑ
 ΜΑΤΛΑΒ!

$$\bar{X} \pm (\text{κρίσιμη τιμή}) \times S_{\bar{X}}$$

$$\pm z_{1-\alpha/2} \frac{s}{\sqrt{n}}, \text{ κανονική, } n > 30$$

$$\pm t_{n-1, 1-\alpha/2} \frac{s}{\sqrt{n}}, \gg, n < 30$$

$z_{1-\alpha/2}$: κρίσιμη τιμή της τυπικής κανονικής κατανομής

$t_{n-1, 1-\alpha/2}$: $\gg \gg \gg t - \gg \gg$ με $n-1$ βαθμούς ελευθερίας.

: μη κανονική, $n < 30$: μη-παραμετρική

(δεν μπορεί να οριστεί με βάση κάποια γνωστή κατανομή)

→ μη παραμετρική μέθοδος.

$$L_1 = \bar{X} - \dots$$

$$L_2 = \bar{X} + \dots$$

$$P(L_1 \leq \mu \leq L_2) = 1 - \alpha$$

π.χ. $\alpha = 0,05$

$z_{1-\alpha/2} = 1,96$

$$P(\bar{X} - 1,96 \cdot S_{\bar{X}} \leq \mu \leq \bar{X} + 1,96 \cdot S_{\bar{X}}) = 0,95$$

Παράδειγμα : $n = 10$ 95% εμπιστοσύνη ($\alpha = 0.05$)

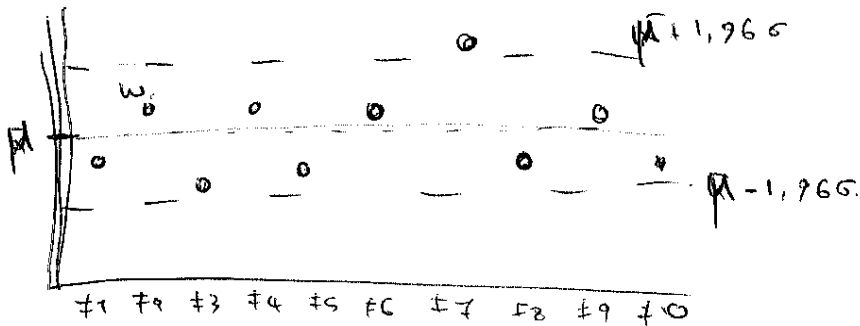
$$L_1 = \bar{x} - 2.3 s\bar{x}$$

$$L_2 = \bar{x} + 2.3 s\bar{x}$$

Για να ελαττώσουμε το εύρος του διαστήματος εμπιστοσύνης το τυπικό σφάλμα ως μέγιστη τιμή ($s\bar{x}$) πρέπει να ελαττωθεί.

Πως : $\downarrow \sigma$ ή $\uparrow n$

Διάστημα εμπιστοσύνης



Είναι ποσοστό των μετρήσεων
είναι μέσα στο $\mu \pm 1,96\sigma$ διάστημα.

$z_{1-\alpha/2} \sigma$: tolerance (ανεπιβιβάσιμος) ϵ

Ευκολότερο σφάλμα ως ποσοστό της μετρήσιμης τιμής

$$\epsilon = \frac{100 z_{1-\alpha/2} \sigma}{w.}$$

Προβλ. : μ : άγνωστος

6.2

Central Limit Theorem (Κεντρικό Ορίσιο Θεώρημα) CLT

Κατανομή μέσω των Τιμών

X_1, X_2, \dots, X_n : τυχαίο δείγμα ^{από} οποιαδήποτε κατανομή με μέση τιμή μ και διασπορά σ^2

n ~~μεγαλώνει~~ ^{μεγαλώνει}

$\frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$ προσεγγίζει την κανονική $N(0, 1)$ κατανομή

$$\Rightarrow \frac{n \bar{X} - n\mu}{\sqrt{n} \sigma} = \frac{\sum_{i=1}^n X_i - n\mu}{\sqrt{n} \sigma} \approx N(0, 1)$$

Q: Τι συμβαίνει όταν $n \rightarrow \infty$? Δείχνει function!

$$\bar{X} \approx N\left(\mu, \frac{\sigma^2}{n}\right)$$

- Νόμος των Μεγάλων Αριθμών (Law of Large Numbers) LLN

\hookrightarrow Ασυμπτωτικά αποτελέσματα προς μέσης τιμής

$$n \rightarrow \infty \quad \bar{X} - \mu \rightarrow 0$$

$$\frac{\text{CLT}}{n \rightarrow \infty} \quad \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \approx N(0, 1)$$

LECTURE 4

Διόστα: εμρτε
 μαζορδομο subscribe εμρτε
 - Διαστήματα εμπιστοσύνης
 - Central Limit Theorem

ΕΛΕΓΧΟΣ ΥΠΟΘΕΣΕΩΝ

π.χ.: Ερώτημα αν διαφέρουν οι μέσες τιμές των Z_1 , X_1 και X_2

ΕΛΕΓΧΟΣ ΕΣΤΑΤΙΣΤΙΚΗΣ ΥΠΟΘΕΣΗΣ

(α)
 ↳ ξεκινάμε με μια μηδενική υπόθεση (null hypothesis) H_0 :
 συνήδως δέλομε να απορρίτουμε να να δεχτούμε την εναλλακτική υπόθεση
 (alternative hypothesis) H_1

π.χ. $H_0: \mu_1 = \mu_2$ $H_1: \mu_1 \neq \mu_2$

(β)
 ↳ επιλέξουμε κατάλληλη στατιστική ελέγχου (test statistic) g που ακολουθεί γνωστή
 κατανομή (π.χ. $g \equiv z \sim N(0,1)$, $g \equiv t \sim t_{n_1+n_2-2} \equiv$ ταυτότητα)

Ορίζουμε ένα βυνοδο κρίσης της g που είναι σίδανο (σε στάδμη
 σημασιότητας α) να πάρει n g αν ισχύει n H_0 το οποίο σημαίνεται
 απορριπτική περιοχή (rejection region) R . Αυτό το βυνοδο αντιστοιχεί
 στις ούρες της κατανομής, π.χ. για $H_0: \mu_1 = \mu_2$ $g \equiv z$ $R = \{ |z| > z_{1-\alpha/2} \}$

(γ)
 ↳ Υπολοοίζουμε την τιμή της g στο δείγμα \tilde{g} (δειγματική στατιστική ελέγχου)
 και εζέτα ζουμε αν ανήκει στο R στε να απορρίτουμε την H_0 .

p -τιμή (p-value): χαμηλότερη στάδμη σημασιότητας οι για την
 οποία μπορούμε να απορρίτουμε την H_0 με βάση το
 δείγμα

π.χ. $g \equiv z$ $p = P(|z| > \tilde{z})$: \tilde{z} : τιμή της στατιστικής g από το δείγμα
 που είναι ίδιο με την κρίσιμη τιμή $z_{1-\alpha/2}$

↳ (δ) Υπάρχει πλήρης συμφωνία των αποτελεσμάτων από το διάστημα εμπιστοσύνης και του στατιστικού έλεγχου στην ίδια στάθμη σημαντικότητας α .

π.χ. αν απορρίψουμε την $H_0 : \mu = \mu_0$ για $\alpha = 0.05$ τότε το 95% διάστημα εμπιστοσύνης του μ δεν περιέχει την μ_0 . Το ίδιο ισχύει για την σύγκριση των \bar{x}_1 και \bar{x}_2 .

Ο έλεγχος μπορεί να είναι μονόπλευρος (one-sided test),

π.χ. $H_0 : \mu \leq \mu_0$ $H_1 : \mu > \mu_0$, τότε ορίζεται μονόπλευρα και η απορριπτική περιοχή, π.χ. $R = \{z > z_{1-\alpha}\}$.

- Type I / Type II errors

Αποφασί	True State	
	H_0 true	H_0 : false
Απορριψη H_0	Type I error α	correct $1-b$
$\mu_0 >> H_0$	correct $1-\alpha$	Type II error b

Type I error: όταν η μηδενική υπόθεση απορριπτεται ενώ είναι true (False Alarm Rate)

Type II error: >> >> >> δεν >> >> >> false (Miss Rate)

$P(\text{Type I error}) = \alpha \rightarrow$ απορριψη της H_0 ενώ ισχύει

$P(\text{Type II error}) = b \rightarrow$ αποδοχή της H_0 ενώ δεν ισχύει

45
- Δύναμη (power) πω. 1 - β

Ερώτηση: Υποθέτουμε ότι υπάρχει μια διαφορά. Θα είμαστε ικανοί να την καταλάβουμε (βρούμε);

Power είναι η πιθανότητα να ανιχνεύσουμε μια πραγματική διαφορά (σε δεδομένο α -διάστημα εμπιστοσύνης) δεδομένων των παραμέτρων του προβλήματος.

Η power εξαρτάται από:

1) Μέγεθος δείγματος n : $\uparrow n \rightarrow \uparrow \text{power} \equiv 1 - \beta$

2) Μέγεθος διαφοράς: $\mu_1 - \mu_2 \downarrow$, $\text{power} \downarrow$

3) α :

LECTURE 5 (01/03/10) \rightarrow Διασπορά

Είδαμε : Διαστήματα εμπιστοσύνης

Έλεγχος υποθέσεων : Βασικές Αρχές

Αποδοχή H_0 (No evidence) \rightarrow $F = \text{false}$

Εφαρμογή τύπου Z , Z^2

Αποδοχή H_0 \rightarrow $F = \text{true}$

Έλεγχος Διαφορών μεταξύ 2 Δειγμάτων

Εστω δύο δείγματα : x_1, x_2, \dots, x_{n_x} , y_1, y_2, \dots, y_{n_y}

t-test : Έλεγχος αν η μέση τιμή δύο συνόλων δεδομένων (δείγματα) είναι διαφορετική

$$t = \frac{\text{Διαφορά μέσων τιμών}}{\text{Διαφορά διασποράς μέσων τιμών}} = \frac{\bar{x} - \bar{y}}{SE(\bar{x} - \bar{y})}$$

$$SE(\bar{x} - \bar{y}) = \sqrt{\frac{S_x^2}{n_x} + \frac{S_y^2}{n_y}}$$

αν $n_x = n_y$

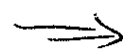
$$SE = \frac{S_x^2 + S_y^2}{2}$$

$$SE(\bar{x} - \bar{y}) = \sqrt{\frac{S_x^2}{n} + \frac{S_y^2}{n}}$$

- Βαθμοί ελευθερίας στο t-test :

$$df = n_x + n_y - 2$$

$$= n_x + n_y - 2$$



- Υπολογίζουμε p-value χρησιμοποιώντας την t-κατανομή

Η p-value είναι η πιθανότητα να απορρίψουμε τη μηδενική υπόθεση (H_0) ενώ αυτή είναι αληθινή (σφάλμα τύπου I)

→ Δεν είναι κατάλληλη μέθοδος για περισσότερα από 2 δείγματα.

Μοναχική Ικανότητα

Παράδειγμα: Μεντιούμ: Διατάγμα τραπεζοκαρτών!

Έλεγχος: Έχουμε 25 κάρτες τις οποίες θα εξετάσουμε και ρωτάμε το Μεντιούμ να μας πει τι "χρώμα" είναι

1 κάρτα

$p = 1/4$

H_0 : το πιάτο δεν είναι χρώμα φεταούκ!
 C : \neq οι κάρτες που θα ~~επιλέγεται~~ εστιάσει σωστά.

Q: Πότε απορρίπτουμε την H_0 ; $\rightarrow A$ $C = 25, 24, 23, 17?$

C: κρίση-Απορριπτική περιοχή \rightarrow συνδέεται με το σφάλμα τύπου I.

$C = 25$ πιθανότητα (τύπου I) $\equiv P(\text{reject } H_0 | H_0 \text{ is valid}) =$

$P(X \geq 25 | p = 1/4) = \left(\frac{1}{4}\right)^{25} \approx 10^{-15}$: πιθανότητα ταχείας πρόβλεψης 25 κάρτων!

$C = 10$. $P(X \geq 10 | p = 1/4) = \approx 0,07?$

Διαδικασία:

- a) H_0
- b) τ στοιχεία \rightarrow μέτρηση κάρτων, επιλέγει α.
- γ) υποδορίζουμε το q.

π.χ. $\alpha = 1\%$

$P(X \geq c | p = 1/4) \leq 0,01 \rightarrow C = 19$: το ελάχιστο είναι

Παράδειγμα R. Fisher (1935)

J. Neyman & Pearson
E. ...

Ανάπτυξη του
Έλεγχου Υπόθεσης 5.3

Μια χώρα υποστηρίζει ότι έχει τη δυνατότητα να προσδιορίσει το
γένη παραγωγής και το είδος τεσσάρων από τη χέρση

a) H_0 : η χώρα δεν έχει αυτή τη δυνατότητα

b) Test statistic : ≠ επιτυχιών σε 8 προστάσεις

c) Distribution-κατανομή : Διωνομική (Yes/No)

Κρίσιμη περιοχή βασισμένη σε κρίσιμο α : " 8 επιτυχίες

8 επιτυχίες > 98% εμπιστοσύνη

Αποτέλεσμα : H_0 απορρίπτεται

Παράδειγμα : Δύο γρήγορα καταγεγραμμένα ενός εξαρτήματος δίνουν την ακόλουθη

μέση διάρκεια ζωής (συνεχής ώρα λειτουργίας)

$\mu_A = 1100$ h $\mu_B = 1300$ h . κανονική κατανομή με γνωστό

απόκλιση $\sigma = 270$ h

Ενός νέου καταγεγραμμένου ισχυρίζεται ότι είναι γινόμενο παρόμοιο εξάρτημα
έχει $\mu = \mu_B$ και χαμηλότερο κόστος παραγωγής

Τέσσερις 12 εξαρτήματα ελέγχονται και δίνει $\bar{x} = 1260$ h

Ερώτηση : είναι $\mu = \mu_B$;

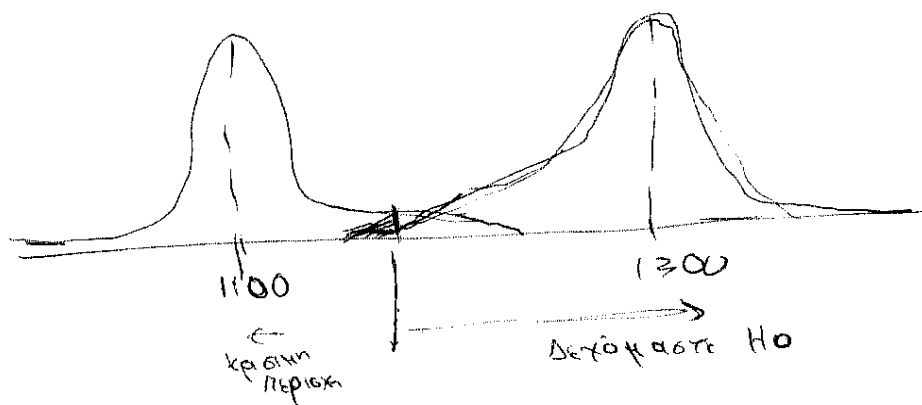
a) $H_0 : \mu = \mu_B = 1300$

$H_1 : \mu = \mu_A = 1100$

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} = \frac{270}{\sqrt{12}} = 77,94$$

b) Στατιστική Ελέγχου : Μέτρηση νέου εξαρτημάτων

Επίδοσθ ο κανονική κατανομή



$$p\text{-value} = 0,05$$

↳ 2 φύλλα έλεγχου αν η εξαγωγική ελίεχου ανήκει στην απορριπτική περιοχή.

1) υπολογίζουμε την απορριπτική περιοχή (για την κατανομή B)

$$\mu. \kappa. \quad \alpha = 0,05 \quad (5\%)$$

$$P(Z \leq 0,05) = -1,64 \Rightarrow \bar{X}_\alpha = \mu_B - 1,64 \cdot \sigma_{\bar{X}} \\ = 1300 - 1,64 \times 77,94 = 1179,2$$

→ Άρα έφωσον $\bar{X} = 1960 > \bar{X}_\alpha$: Δεχόμαστε την H_0

2) υπολογίζουμε την πιθανότητα να έχουμε απόκλιση από τη θεωρητική τιμή της H_0

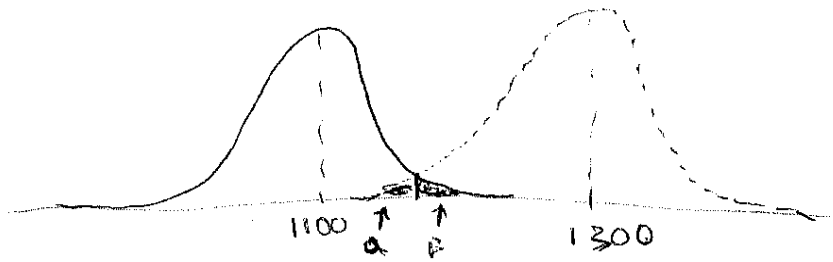
$$P = P(T \geq t_n(X)) \quad \text{ή} \quad p = P(T \in t_n(X))$$

Αυτή χρησιμοποιείται από τα statistic software.

$$P = P\left(Z \leq \frac{\bar{X} - \mu_B}{\sigma_{\bar{X}}}\right) = P\left(Z \leq \frac{1960 - 1300}{77,94}\right) = 0,304$$

Έφωσον $P > 0,05$ δεν έχουμε καμία απόδειξη απόρριψης της H_0 .

Επείδη τωρον ΕΙ (αποδοχή της Η₀ ενώ είναι False)

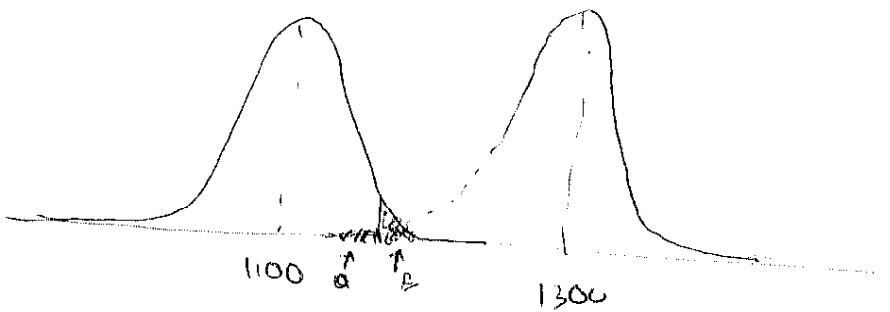


$\bar{X}_0 = 1200$: Όριο απόφασης

$$\alpha = \beta = P \left(z \leq \frac{(1200 - 1300)}{77,94} \right) = N_{0,1}(-1,283) = 0,10$$

$$\beta = P \left(z \leq \frac{(1200 - 1100)}{77,94} \right) = \gg = 0,10$$

Όταν ελαττώνουμε το α $\rightarrow \uparrow \beta$



π.χ $\alpha = 0,05$

$$\beta = P \left(z \geq \frac{(X_0 - \mu_a)}{\sigma_X} \right) = P \left(z \geq \frac{1172,8 - 1100}{77,94} \right) = 0,177!$$

Επιλέγουμε τις Η₁

$\alpha = 0,05$

Η₁: μ_A

$z = (\mu_A - \bar{X}_{0,05}) / \sigma_X$

β

$1 - \beta$

1100	0,93	0,18	0,82
1172,2	0,00	0,50	0,50
1200	-0,36	0,64	0,36
1300	-1,64	0,95	0,05

- Για να αυξηθούν τα $(1-\beta)$ μπορούμε να αυξηθούν τα n

$\rightarrow \downarrow \sigma_{\bar{x}}$

π.χ. $n = 24, \alpha = 0,05$

$$\bar{X}_A = \mu_B - 1,64 \times \sigma_{\bar{x}} = 1300 - 1,64 \cdot 55,11 = \underline{\underline{1209,6}}$$

μ_A	$z = (\mu_A - \bar{X}_{0,05}) / \sigma_{\bar{x}}$	β	$1-\beta$
1100	1,99	0,09	0,91
1150	1,08	0,14	0,86
1200	0,17	0,43	0,57
1300	-1,64	0,95	0,05

LECTURE 6 (08/03/10)

Είδη: Έλεγχος Στατιστικών Υποθέσεων

Βασικά Παραδείγματα

Διαστήματα Εμπιστοσύνης ε.γ. $P(\mu \in [95, 105]) = \delta = 1 - \alpha$

Εργαστήριο - Lab

Παρασκευή 13:00-15:00 Η205

Πέμπτη

15:00-17:00 Η205

Τετάρτη 13:00-16:00

Η203 ε.γ. = 2 x 1.5 per group

Common Test Statistics

e.g. One-sample z-test

$$z_{\text{calc}} = z_0 = \frac{\bar{X} - \mu_0}{\left(\frac{\sigma}{\sqrt{n}}\right)}$$

$$\text{or } \mu_0 = \bar{X} - (z_{1-\alpha/2} \times \frac{\sigma}{\sqrt{n}})$$

$$= \bar{X} - z_{1-\alpha/2} \left(\frac{\sigma}{\sqrt{n}}\right)$$

$$H_0: \mu = \mu_0$$

$$\text{or } \bar{X} = \mu_0 \pm z_{1-\alpha/2} \left(\frac{\sigma}{\sqrt{n}}\right)$$

$$P = P(T \geq t_{n,\alpha}) \text{ or } P(T \leq t_{n,\alpha})$$

$$P = P(|Z| \leq \frac{\bar{X} - \mu_0}{\left(\frac{\sigma}{\sqrt{n}}\right)}) = P(|Z| \leq z_0)$$

IF $P(|Z| \leq z_0) \geq \alpha$ accept H_0
 $< \alpha$ reject H_0

↳ Test statistics → Παραμετρικά
μη παραμετρικά.

↳ περίληψη διαφορετικών test.

Εταριστικός έλεγχος 2 ευνόμων δεδομένων (δείγματων)

π.χ. A, B μ_A, μ_B sets

$H_0 : \mu_A = \mu_B$ (ή $\mu_A - \mu_B = 0$)

$H_1 : \mu_A \neq \mu_B$: two-sided test

ή $H_0 : \mu_A \leq \mu_B$, $H_1 : \mu_A > \mu_B$

ή $H_0 : \mu_A \geq \mu_B$, $H_1 : \mu_A < \mu_B$

} One-sided test.

π.χ. sets $X = \{x_1, x_2, \dots, x_n\}$ $Y = \{y_1, y_2, \dots, y_n\}$

N_{μ_A, σ_A}

N_{μ_B, σ_B}

παραμέτροι: μ_A, μ_B (or $\mu_A + \underline{\Delta}$) , σ_A, σ_B

Αν ορίσουμε $\underline{\Delta} = (x_1 - y_1), (x_2 - y_2), \dots, (x_n - y_n)$, σ_{Δ}

Έχουμε ως παραμέτρους μόνο Δ, σ_{Δ}

Βασική ιδέα : ανεξάρτητα δείγματα ("independent samples")

↳ τα δεδομένα έχουν επιλεγεί με τέτοιο τρόπο ώστε όλοι οι παράγοντες που ^{τα} επηρεάζουν, εκτός αυτών που μας ενδιαφέρουν να λαμβάνονται υπόψη με τον ίδιο τρόπο

Παράδειγμα: Ανεξάρτητα Δείγματα.

(α) ⇒ Δύο, Placebo : Δράση ενός φαρμάκου. Δύο δείγματα,

στο ένα χορηγείται φάρμακο στο άλλο placebo. Κατόπιν μετρείται η επίδραση στην ασθένεια.

- Αν υπάρχουν και άλλοι παράμετροι που επηρεάζουν, όπως φύλο, ηλικία τότε για να είναι τα 2 δείγματα ανεξάρτητα πρέπει να επεκταθούν ώστε να είναι πραγματικά τυχαία όσον αφορά φύλο και ηλικία.

(β) ⇒ Πρόγραμμα TV : Δείγμα → Άνδρες, Γυναίκες

- Πρέπει τα δείγματα να είναι τυχαία όσο αφορά άλλες ιδιότητες όπως : ηλικία, μορφωτικό επίπεδο, εισόδημα, κλπ.

Ανάλογα κατά ζευγάρια (paired samples)

(α) : εύνολο ασθενών που τη μια φορά λαμβάνουν placebo και την άλλη φάρμακο.

(β) ζευγάρια "συνζυγών". Ο/Η σύζυγος μπορεί (συνήθως) να έχουν παρόμοια δεδομένα όσο αφορά την ηλικία, παιδιά, συνήθειες. Προσοχή : κάλλιστα αυτές οι παραδοχές μπορεί να είναι λάθος

Παράδειγμα : (two-tailed test) : Αξιολόγηση Διαστάσεων. 64

Εε ένα σχολείο οι μαθητές κατανέμονται ελάχιστα σε 2
καθηγητές μαθηματικών. κ. Α κ. Β

Μετα την κατανομή κ. Α : 30 μαθητές
κ. Β : 25 >>

Ετο τέλος του χρόνου όλοι οι μαθητές πέρνουν το ίδιο
διαγώνισμα. οι μαθητές του κ. Α έχουν μέσο όρο βαθμολογίας 78
και τυπική απόκλιση 10

οι >> >> κ. Β έχουν μέσο όρο 85
και τυπική απόκλιση 15

Έλεγχος της υπόθεσης ότι οι κ. Α και κ. Β είναι το ίδιο
καλοί καθηγητές. Επίπεδο σημαντικότητας $\alpha = 0,1$

α) Μηδενική υπόθεση

$$H_0 : \mu_1 - \mu_2 = 0$$

$$H_1 : \mu_1 - \mu_2 \neq 0$$

Two-sided test : Η H_0 θα απορριφθεί αν $\mu_1 - \mu_2 > \epsilon$
ή $\mu_1 - \mu_2 < -\epsilon$

β) Στατιστική Έλεγχου : οι βαθμοί

t-κατανομή (διαφορετικές αποκλίσεις - διασπορές)

$$SE = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}} = \sqrt{\frac{10^2}{30} + \frac{15^2}{25}} = \sqrt{3,33 + 9} = 3,51$$

$$s^2 = \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\frac{\left(\frac{s_1^2}{n_1}\right)^2}{n_1-1} + \frac{\left(\frac{s_2^2}{n_2}\right)^2}{n_2-1}} = \frac{\left(\frac{10^2}{30} + \frac{15^2}{25}\right)^2}{\frac{\left(\frac{10^2}{30}\right)^2}{29} + \frac{\left(\frac{15^2}{25}\right)^2}{24}}$$

$$\Rightarrow s^2 = \frac{(3,33 + 9)^2}{\frac{3,33^2}{29} + \frac{9^2}{24}} = \frac{152,03}{0,382 + 3,375}$$

$$= \frac{152,03}{3,757} = 40,47$$

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - \delta_0}{SE} = \frac{(78 - 85) - 0}{3,51} = \frac{-7}{3,51} = -1,99$$

\Rightarrow P-value : πιθανότητα ένα t-score με 40

πλάγια να
καταρτηθεί για την
της στατιστικής εξέχου
16η με την t-εξέχου

βαθμούς ελευθερίας (degrees of freedom) να

είναι $\leq -1,99$ ή $> 1,99$

t - κατανομή

$$P(t < -1,99) = 0,027$$

$$P(t > +1,99) = 0,027$$

$$\Rightarrow \text{p-value} = 0,054$$

Αποτέλεσμα Εφόσον p-value $< \alpha = 0,10$ δεν μπορούμε να δεχτούμε την H_0 !

Common Test Statistics

In the table below, the symbols used are defined at the bottom of the table. Many other tests can be found in the literature (other articles).

Name	Formula	Assumptions or notes
One-sample <u>z</u> -test	$z = \frac{\bar{x} - \mu_0}{(\sigma/\sqrt{n})}$	(Normal population or $n > 30$) and σ known. (z is the distance from the mean in relation to the standard deviation of the mean). For non-normal distributions it is possible to calculate a minimum proportion of a population that falls within k standard deviations for any k (see: <i>Chebyshev's inequality</i>).
Two-sample <u>z</u> -test	$z = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$	Normal population and independent observations and σ_1 and σ_2 are known
One-sample <u>t</u> -test	$t = \frac{\bar{x} - \mu_0}{(s/\sqrt{n})},$ $df = n - 1$	(Normal population or $n > 30$) and σ unknown. For non-normal populations, n should be large enough to ensure both that dist. of mean is close to normal and that s is a good estimate of σ .
Paired t-test	$t = \frac{\bar{d} - d_0}{(s_d/\sqrt{n})},$ $df = n - 1$	(Normal population of differences or $n > 30$) and σ unknown
Two-sample pooled t-test, equal variances*	$t = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}},$ $s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2},$ $df = n_1 + n_2 - 2_{[S]}$	(Normal populations or $n_1 + n_2 > 40$) and independent observations and $\sigma_1 = \sigma_2$ and σ_1 and σ_2 unknown
Two-sample unpooled t-test, unequal variances*	$t = \frac{(\bar{x}_1 - \bar{x}_2) - d_0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}},$	(Normal populations or $n_1 + n_2 > 40$) and independent observations and $\sigma_1 \neq \sigma_2$ and σ_1 and σ_2 unknown

	$df = \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\frac{\left(\frac{s_1^2}{n_1}\right)^2}{n_1-1} + \frac{\left(\frac{s_2^2}{n_2}\right)^2}{n_2-1} \quad [6]$	
One-proportion z-test	$z = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$	$n \cdot p_0 > 10$ and $n(1 - p_0) > 10$ and it is a SRS (Simple Random Sample).
Two-proportion z-test, pooled	$z = \frac{(\hat{p}_1 - \hat{p}_2) - d_p}{\sqrt{\hat{p}(1 - \hat{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$ $\hat{p} = \frac{x_1 + x_2}{n_1 + n_2}$	$n_1 p_1 > 5$ and $n_1(1 - p_1) > 5$ and $n_2 p_2 > 5$ and $n_2(1 - p_2) > 5$ and independent observations
Two-proportion z-test, unpooled	$z = \frac{(\hat{p}_1 - \hat{p}_2) - d_p}{\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}}$	$n_1 p_1 > 5$ and $n_1(1 - p_1) > 5$ and $n_2 p_2 > 5$ and $n_2(1 - p_2) > 5$ and independent observations
One-sample chi-square test	$\chi^2 = \frac{(n - 1)s^2}{\sigma_0^2}$	One of the following <ul style="list-style-type: none"> • All expected counts are at least 5 • All expected counts are > 1 and no more than 20% of expected counts are less than 5
*Two-sample F test for equality of variances	$F = \frac{s_1^2}{s_2^2}$	Arrange so $s_1^2 \geq s_2^2$ and reject H_0 for $F > F(\alpha/2, n_1 - 1, n_2 - 1)$ [7]

Definitions of symbols:

- α , the probability of Type I error (rejecting a null hypothesis when it is in fact true)
- n = sample size
- n_1 = sample 1 size
- n_2 = sample 2 size
- \bar{x} = sample mean
- μ_0 = hypothesized population mean
- μ_1 = population 1 mean
- μ_2 = population 2 mean
- σ = population standard deviation
- σ^2 = population variance

- s = sample standard deviation
- s^2 = sample variance
- s_1 = sample 1 standard deviation
- s_2 = sample 2 standard deviation

- t = t statistic
- df = degrees of freedom
- \bar{d} = sample mean of differences
- d_0 = hypothesized population mean difference
- s_d = standard deviation of differences

- $\hat{p} = x/n$ = sample proportion, unless specified otherwise
- p_0 = hypothesized population proportion
- p_1 = proportion 1
- p_2 = proportion 2
- d_p = hypothesized difference in proportion
- $\min\{n_1, n_2\}$ = minimum of n_1 and n_2
- $x_1 = n_1 p_1$
- $x_2 = n_2 p_2$
- χ^2 = Chi-squared statistic
- F = F statistic

In general, the subscript 0 indicates a value taken from the null hypothesis, H_0 , which should be used as much as possible in constructing its test statistic.

LECTURE 7 (15/03/10)

Είδαμε: Έλεγχος Στατιστικών Υποθέσεων

Παραδείγματα

Lab : Παραγωγή 13:00 - 15:00 ΗΡΟΣ
 Τετάρτη 13:00 - 15:00 ΗΡΟΣ → Free OK
 ή 13:00 - 16:00 ΗΡΟΣ

CH.3

Ανάλυση Διακύμανσης (ή Διασποράς) - Analysis of Variance.

Είδαμε ότι ο έλεγχος 2 στατιστικών δειγμάτων γίνεται με τους συνδυαστικούς έλεγχους.

Τι γίνεται σε περίπτωση που έχουμε έλεγχο > 2 δειγμάτων;

π.χ k ανεξάρτητων πληθυσμών

μ_j
 σ_j $j = 1, 2, \dots, k.$

→ τυχαία δείγμα

\bar{x}_j, \bar{s}_j

Πληθυσμός	1	2	...	k
Μέση τιμή	μ_1	μ_2	...	μ_k
Τυπική Αποκλίση	σ_1	σ_2	...	σ_k
Δείγματα				
Μέση τιμή	\bar{x}_1	\bar{x}_2	...	\bar{x}_k
Τυπική Αποκλίση	\bar{s}_1	\bar{s}_2	...	\bar{s}_k
Μέγεθος Δείγματος	n_1	n_2	...	n_k

π.χ. < διαφορετικά πειράματα (ή δειγμάτεις) ενός προβλήματος

πλήρως τυχαιοποιημένα σχέδια (completely randomized designs)

Έστω $H_0 : \mu_1 = \mu_2 = \dots = \mu_k$

H_1 : δύο τουλάχιστον από τους μ_i διαφέρουν μεταξύ τους.

Μικτή λύση : πολλά πολλά t-tests.

Πρόβλημα :

Α) Χρειαζόμαστε $\binom{k}{q}$ t-tests. π.χ. $k=5 \quad \binom{5}{q} = \frac{5!}{q! 3!} = \frac{1 \cdot 2 \cdot 3 \cdot 4 \cdot 5}{1 \cdot 2 \cdot 1 \cdot 2 \cdot 3} = 10$

$\binom{M}{q}$ η απεικόνιση ανά $q \Rightarrow \equiv n C_q = \frac{n!}{q! (n-q)!}$ π.χ. $\binom{10}{q} = \frac{10!}{q! 8!} = 45$

Άρα πολλά t-tests !

Β) μπορεί να οδηγήσουμε σε λάθος συμπεράσματα

Η πιθανότητα σφάλματος τύπου I (H_0 rejected / H_0 valid) αυξάνεται !

Γιατί :

π.χ. $\alpha = 0,05$ για 2 εξερχούς $H_0^1 : \mu_1 = \mu_2$
 $H_0^2 : \mu_1 = \mu_2$

$P(H_0^1, H_0^2) = P(H_0^1) \cdot P(H_0^2) = (1-\alpha) \cdot (1-\alpha) = 0,95^2 = 0,9025$

Άρα συνολικό σφάλμα τύπου I = $1 - 0,9025 = 0,0975$.

Γενικά $\gg \gg \gg \Gamma = 1 - (1-\alpha)^m$ m : # εξερχούς

Ανάλυση Διακύμανσης ως προς έναν Παράγοντα

↳ υποθέτουμε ότι υπάρχουν ένας ή περισσότεροι παράγοντες ως προς τον οποίο οι πληθυσμοί διαφοροποιούνται μεταξύ τους.

Βασική παραδοχή (όπως και στα t-test) : $\sigma_1^2 = \sigma_2^2 = \dots = \sigma_k^2 = \sigma^2$

Ίδια διακύμανση

- Εκτίμηση κοινής διακύμανσης σ^2 μέσω της

$$S_w^2 = \frac{(n_1 - 1)S_1^2 + (n_2 - 1)S_2^2 + \dots + (n_k - 1)S_k^2}{n_1 + n_2 + \dots + n_k - k} = \text{διασπορά στο εσωτερικό των ομάδων}$$

(within groups variability)

Γενικά η συνολική μεταλλητότητα k πληθυσμών ονομάζεται σε

↳ 1) διασπορά ανά διακύμανση σύμφωνα με την μέση τιμή του πληθυσμού

↳ 2) >>> συνολική μέση τιμή

$$S_B^2 = \frac{n_1(\bar{x}_1 - \bar{x})^2 + n_2(\bar{x}_2 - \bar{x})^2 + \dots + n_k(\bar{x}_k - \bar{x})^2}{k-1}$$

(between groups variability)

$$\text{όπου } \bar{x} = \frac{n_1\bar{x}_1 + n_2\bar{x}_2 + \dots + n_k\bar{x}_k}{n_1 + n_2 + \dots + n_k} = \frac{n_1\bar{x}_1 + n_2\bar{x}_2 + \dots + n_k\bar{x}_k}{N}$$

Εάν ισχύει η H_0 : $\mu_1 = \mu_2 = \dots = \mu_k$ ^{πρόκειται} $S_B^2 = S_w^2$ ή

$$F = \frac{S_B^2}{S_w^2} = 1. \text{ Αν υπάρχουν διαφορές } F \neq 1 \text{ (}> 1)$$

Παρένθεση: Κατανομή F (ομογενείς κατανομών)

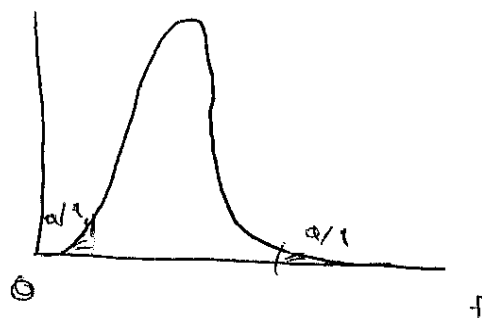
Αφορά έλεγχους διακυμάνσεων 2 κανονικά κατανομμένων πληθυσμών

$$F = \frac{s_1^2 / \sigma_1^2}{s_2^2 / \sigma_2^2} \quad : \text{ομολοποιεί την κατανομή}$$

F με $n_1 - 1, n_2 - 1$ βαθμούς ελευθερίας

π.χ $H_0 : \sigma_1^2 = \sigma_2^2$

$H_1 : \sigma_1^2 \neq \sigma_2^2$



Αν η τιμή της $F = \frac{s_1^2}{s_2^2}$ που υπολογίζεται από τα διαθέσιμα δεδομένα είναι $> F_{n_1-1, n_2-1, \alpha/2}$ που ορίζει μια περιοχή εμπέδου $\alpha/2$ στο δεξί άκρο της αντίστοιχης κατανομής τότε η H_0 απορρίπτεται! (πρόσκληση: $\sigma_1^2 > \sigma_2^2$)

Παράδειγμα : Δύο αναλυτικά σκευάσματα χορηγούνται σε ασθενείς

$n_1 = 13, n_2 = 13 \rightarrow$ Διακυμάνσεις $s_1^2 = 64, s_2^2 = 16$

$\alpha = 0,05$ Να ελεγχθεί αν $\sigma_1^2 > \sigma_2^2$

$H_0 : \sigma_1^2 = \sigma_2^2$

$H_1 : \sigma_1^2 > \sigma_2^2$

$\rightarrow F_{12, 12, 0,05} = 2,69 \equiv F_{crit.}$

$F = \frac{s_1^2}{s_2^2} = \frac{64}{16} = 4 > F_{crit.} \rightarrow$

\rightarrow απορρίπτουμε $H_0 \rightarrow$ άρα δεχόμαστε ότι $\sigma_1^2 > \sigma_2^2$.

75
Αρα, εφόσον έχουμε διατυπώσεις χρησιμοποιούμε
την κατανομή F . \Rightarrow

$\Rightarrow F = \frac{S_B^2}{S_W^2}$ ακολουθεί κατανομή F με $k-1, N-k$ β.ε.

$$F_{k-1, N-k, \alpha} \equiv F_{crit}$$

Διαδικασία : Υποθέτουμε F . Αν $F > F_{crit}$. $H_0 \Rightarrow$ reject

Αν $F \leq F_{crit}$. H_0 : accept, δηλαδή \rightarrow

$$\rightarrow \mu = \mu_1 = \mu_2 = \dots = \mu_k$$

$$\sigma = \sigma_1 = \sigma_2 = \dots = \sigma_k.$$

ΑΝΟΝΑ : Συμβολισμός

Έχουμε k πληθυσμούς $\mu_j, \sigma_j \quad j=1, 2, \dots, k$

και ε.δ. $\bar{x}_j, s_j, \quad x_{ij} : \begin{matrix} i=1, 2, \dots, n_j \\ j=1, 2, \dots, k \end{matrix}$

Συνολική διακύμανση των πληθυσμών:

$$S_T^2 = \frac{\sum_{j=1}^k \sum_{i=1}^{n_j} (x_{ij} - \bar{x})^2}{N-1} \quad (1)$$

Οι δύο συνιστώσες της συνολικής διασποράς είναι:

$$- S_W^2 = \frac{(n_1-1)s_1^2 + (n_2-1)s_2^2 + \dots + (n_k-1)s_k^2}{N-k} = \frac{\sum_{j=1}^k (n_j-1)s_j^2}{N-k}$$

$$\left[S_T^2 - \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 \right] = \frac{\sum_{j=1}^k \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_j)^2}{N-k} \quad (2)$$

$$- S_B^2 = \frac{n_1(\bar{x}_1 - \bar{x})^2 + n_2(\bar{x}_2 - \bar{x})^2 + \dots + n_k(\bar{x}_k - \bar{x})^2}{k-1} = \frac{\sum_{j=1}^k n_j(\bar{x}_j - \bar{x})^2}{k-1} \quad (3)$$

Οι αριθμητές στις παραπάνω σχέσεις έχουν για

συγκεκριμένη ονομασία.



7.7
⇒ (α) τῆς S_T^2 : συνολικό άδροισμα τετραγώνων (Total Sum of Squares) TSS

$$TSS = \sum_{j=1}^k \sum_{i=1}^{n_j} (x_{ij} - \bar{x})^2$$

⇒ (β) τῆς S_W^2 : άδροισμα τετραγώνων στο εσωτερικό των ομάδων (Sum of Squares within-groups) SSW

$$SSW = \sum_{j=1}^k \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_j)^2$$

⇒ (γ) τῆς S_B^2 : άδροισμα τετραγώνων μεταξύ των ομάδων (Sum of Squares between-groups) SSB

$$SSB = \sum_{j=1}^k n_j (\bar{x}_j - \bar{x})^2$$

$$TSS = SSB + SSW \quad \rightarrow \text{Απόδειξη 7.7}$$

Αντίστοιχα :

$$S_B^2 = \frac{SSB}{k-1} : \text{ μέσο τετραγώνο μεταξύ των ομάδων}$$

$$S_W^2 = \frac{SSW}{N-k} : \text{ μέσο τετραγώνο στο εσωτερικό των ομάδων}$$

LECTURE 8

ANOVA Πινάκας

Είδημε :

Πηγή Διασποράς	Άθροισμα Τετραγώνων	Βαθμοί Ελευθερίας	Μέσο Τετράγωνο	F
Μεταξύ των Ομάδων	SSB	k - 1	$S_B^q = \frac{SSB}{k-1}$	$\frac{S_B^q}{S_W^q}$
Ενδο-ομάδων	SSW	N - k	$S_W^q = \frac{SSW}{N-k}$	
Εύννοο	SST	N - 1		

Οι υπολογισμοί είναι αρκετά πολύπλοκοι και χρονοβόροι →
 → χρήση SOFTWARE.

Προβλεψη - Αποκρίνεται : Έλεγχος προϋποθέσεων της μεθόδου
 Ορδη ανάγνωση των αποτελεσμάτων

Οότε
 ⇒ να έχουμε άσφαλτη συμπεράσματα.

$$SSB = \sum_{j=1}^k n_j (\bar{x}_j - \bar{x})^2$$

$$SSW = \sum_{j=1}^k \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_j)^2$$

$$SST = \sum_{j=1}^k \sum_{i=1}^{n_j} (x_{ij} - \bar{x})^2$$

Παράδειγμα: Επιδημιολογική μελέτη

Μετρήθηκε ο αριθμός των λευκών αιμοσφαιρίων σε 3 ομάδες υγιών ατόμων με διαφορετικές συνήθειες καπνίσματος.

Όλοι ήταν άνδρες ηλικίας 25-34 ετών.

1^η ομάδα: 49 μη καπνιστές $n_1 = 49$

$$\bar{X}_1 = 8071,4, \quad S_1 = 1771,6$$

2^η ομάδα: 96 καπνιστές

$$n_2 = 96$$

με < 20 τσιγάρα/ημέρα

$$\bar{X}_2 = 8680,8, \quad S_2 = 1750,2$$

3^η ομάδα: 39 καπνιστές

$$n_3 = 39$$

> 20 τσιγάρα/ημέρα

$$\bar{X}_3 = 10809,4, \quad S_3 = 1795,1$$

Ερώτηση: αν οι πληθυσμιακές μέσες τιμές των λευκών αιμοσφαιρίων των τριών ομάδων διαφέρουν σημαντικά μεταξύ τους.

Αρχικά:

$$\begin{aligned} \bar{X} &= \frac{n_1 \bar{X}_1 + n_2 \bar{X}_2 + n_3 \bar{X}_3}{n_1 + n_2 + n_3} = \frac{49(8071,4) + 96(8680,8) + 39(10809,4)}{(49 + 96 + 39)} = 9106. \end{aligned}$$

Κατόπιν υπολογίζουμε τις διαφόρες συνιστώσες της διασποράς

$$\begin{aligned} -SSB &= n_1 (\bar{X}_1 - \bar{X})^2 + n_2 (\bar{X}_2 - \bar{X})^2 + n_3 (\bar{X}_3 - \bar{X})^2 = \\ &= 49(8071,4 - 9106)^2 + 96(8680,8 - 9106)^2 + 39(10809,4 - 9106)^2 \\ &= 14950314. \end{aligned}$$

$$- SSW = (n_1 - 1)S_1^2 + (n_2 - 1)S_2^2 + (n_3 - 1)S_3^2 =$$

$$= (49 - 1)(1771,6^2) + (96 - 1)(1750,2^2) + (39 - 1)(1795,4^2) =$$

$$= 305188596$$

και

$$S_B^2 = \frac{SSB}{k-1} = \frac{149503114}{3-1} = 71951557$$

$$S_W^2 = \frac{SSW}{N-k} = \frac{305188596}{100-3} = 3146273, \text{ ~~ok~~ } \underline{\underline{ok}}$$

Υποθέσεις:

$$H_0 : \mu_1 = \mu_2 = \mu_3$$

H_A : Δύο τουλάχιστον από τους μ_1, μ_2, μ_3 διαφέρουν μεταξύ τους
 $\mu_1 \neq \mu_2$ ή $\mu_1 \neq \mu_3$ ή $\mu_2 \neq \mu_3$

- Επίπεδο σημαντικότητας ελέγχου $\alpha = 0,05$

= Τιμή κριτηρίου ελέγχου $F = \frac{S_B^2}{S_W^2} = \frac{71951557}{3146273} = 22,88$

Κρίσιμη τιμή της κατανομής $F_{\alpha, k, N-k} = F_{0,05, 2, 97} < 3,11$

- Πρόβωση $F > F_{\alpha, k, N-k} : H_0 : \text{γερπρεδ}$

↳ Ο μέγος αριθμός λευκών αιμοσφαιρίων διαφοροποιείται σε δύο τουλάχιστον από τους τρεις πληθυσμούς.

ΑΝΟΒΑ : πολλαπλές συγκρίσεις

Είδαμε ότι η $H_0 : \mu_1 = \mu_2 = \dots = \mu_k$ έχει ως εναλλακτική υπόθεση

$H_A : \geq$ τουλάχιστον δύο μ είναι διαφορετικές μεταξύ τους. \Rightarrow

\Rightarrow Δεν μας βοηθά στον προσδιορισμό των μέσων τιμών που διαφέρουν μεταξύ τους. \Rightarrow

\Rightarrow Απαιτούνται έλεγκοί πολλαπλών συγκρίσεων.

(Multiple Comparison Tests)

Πολλές διαφορετικά Test

Α) Έλεγχος του Fisher (Least Significant Difference LSD) (1935)

: Διευρημένη έκδοση του t-test.

Πρώτα κάνουμε το τυπικό F-test. Όταν έχει απορριφθεί η H_0 για επίπεδο σημαντικότητας α , η ελάχιστη σημαντική διαφορά δύο οποιονδήποτε μέσων τιμών μ_i και μ_j είναι:

$$LSD = t_{n-k, \alpha/2} \sqrt{s_w^2 \left(\frac{1}{n_i} + \frac{1}{n_j} \right)} \quad \text{Αν } |\bar{x}_i - \bar{x}_j| \geq LSD$$

η $100(1-\alpha)\%$ Δ.Ε.

$\Rightarrow \mu_i, \mu_j$: διαφέρουν σημαντικά μεταξύ τους

$$(\bar{y}_i - \bar{y}_j) > \pm t_{n-k, \alpha/2} \sqrt{s_w^2 \left(\frac{1}{n_i} + \frac{1}{n_j} \right)}$$

↳ Αρκετά "ελαστικός" έλεγχος. Δεν διασφαλίζει πάντα ότι το συνολικό σφάλμα τύπου I είναι 0!

B) Έλεγχος του Tukey

: Χρησιμοποιεί την τυποποιημένη κατά Student κατανομή εύρους (studentized range distribution) και υπολογίζεται

στο :

$$\frac{\bar{X}_{max} - \bar{X}_{min}}{\sqrt{\frac{S_w^2}{n}}}$$

$\bar{X}_{max} - \bar{X}_{min}$ = εύρος των δειγματικών μέσων τιμών κ ομάδων.

⇒ Για $n_1 = n_2 = n$

Έλεγχος με την κρίσιμη τιμή

$$W = q_{\alpha, k, n-k} \sqrt{\frac{S_w^2}{n}}$$

σημ. ως studentized range distribution

AN $|\bar{x}_i - \bar{x}_j| \geq q_{\alpha, k, n-k} \sqrt{\frac{S_w^2}{n}} \Rightarrow$

⇒ οι μ_i, μ_j διαφέρουν σημαντικά μεταξύ τους.

AN $n_1 \neq n_2$ $W_{i,j} = q_{\alpha, k, n-k} \sqrt{\frac{S_w^2}{2} \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$

LECTURE 9

ΕΙΣ/ΟΕ)

(Πως μένει το πρόβλημα +10%)

Είδαμε:

= ANOVA με γ μεταβλητή

κ : ανεξάρτητοι πληθυσμοί

$$\left\{ \begin{array}{l} \bar{x}_j, \mu_j \\ s_j, \sigma_j \end{array} \right\} j = 1, 2, \dots, \kappa$$

↳ μινιμουμ ANOVA

Πως Διασφραξ	Άθροισμα Τετραγώνων	Βαθμιαί (λευδέρια)	Μεσο Τετραγώνω	F
Μεταξύ των Ομάδων	$SSB = \sum_{j=1}^{\kappa} n_j (\bar{x}_j - \bar{x})^2$	$\kappa - 1$	$S_B^2 = \frac{SSB}{\kappa - 1}$	$\frac{S_B^2}{S_W^2}$
Εσωτερικά των Ομάδων	$SSW = \sum_{j=1}^{\kappa} \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_j)^2$	$N - \kappa$	$S_W^2 = \frac{SSW}{N - \kappa}$	
	$SS_T = \sum_{j=1}^{\kappa} \sum_{i=1}^{n_j} (x_{ij} - \bar{x})^2$	$N - 1$		

- $H_0 : \mu_1 = \mu_2 = \dots = \mu_{\kappa}$
- $H_A : \neq H_0$

: Δείνουμε λοιπόν στον προσδιορισμό των μέσων τιμών που διαφέρουν μεταξύ τους. Αποκρίνεται =>

=> Πολλαπλά Test

Διαφορετικές Μέθοδοι:

A) Test Fisher (Ελάχιστη σημαντική διαφορά) LSD

$$= t_{n-\kappa, \alpha/2} \sqrt{S_W^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$$

B) Test Tukey

$|\bar{x}_1 - \bar{x}_2| \geq LSD \Rightarrow$
 μ_1, μ_2 διαφέρουν σημαντικά μεταξύ τους
 Δεν μας δείχνει όμως ότι το συνολικό σύνολο τύπου F είναι α!

↳ Studentized Range Distribution

$$\frac{\bar{x}_{max} - \bar{x}_{min}}{\sqrt{S_W^2/n}} \quad \text{Έλεγχος με } \omega = q_{\alpha, \kappa, N-\kappa} \sqrt{\frac{S_W^2}{n}} \quad |\bar{x}_i - \bar{x}_j| \geq \omega$$

Γ ⇒ Έλεγχος του Bonferroni (Bonferroni Correction)

: Απλοί t-test με τροποποιημένο επίπεδο σημαντικότητας α .

Είδαμε ότι σε k πηλυσμούς $\binom{k}{2} = \frac{k(k-1)}{2}$ διαφορετικά

t-test οδηγούν σε αύξηση του συνολικού σφάλματος

τύπου I. \Rightarrow

\Rightarrow Τροποποιημένο $\alpha^* \rightarrow \alpha^* = \frac{\alpha}{\binom{k}{2}}$

α^* : επίπεδο σημαντικότητας σε κάθε επιμέρους έλεγχο ώστε το συνολικό σφάλμα τύπου I να είναι α .

π.χ. $k = 3$ $\binom{3}{2} = 3$ $\alpha = 0,05$ $\alpha^* = \frac{0,05}{3} = 0,016$

\Rightarrow Άρα. $H_0: \mu_1 = \mu_2$

\rightarrow υπολογίζουμε την $t_{12} = \frac{(\bar{x}_1 - \bar{x}_2)}{\sqrt{s_w^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}}$

\rightarrow κρίσιμη τιμή της και τα νομής. $t_{n-2, \alpha^*/2}$

Δ.Ε. $(\bar{x}_1 - \bar{x}_2) \pm t_{n-2, \alpha^*/2} \sqrt{s_w^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$

Γενικά :

- * Για μικρό αριθμό πληθυσμών (π.χ. $k < 10$) ο έλεγχος Bonferroni, είναι πολύ χρήσιμος.
- * Για μεγαλύτερο αριθμό πληθυσμών το $\alpha^* = \frac{\alpha}{\binom{k}{2}}$ γίνεται πολύ μικρό και επομένως το σφάλμα τύπου I κάθε επιμέρους ελέγχου γίνεται πολύ μικρό και το σφάλμα τύπου II πολύ μεγάλο! Τότε \Rightarrow χρησιμοποιείται περισσότερο το test Tukey.

* Προσοχή: Υπάρχουν ακόμη πολλά διαφορετικά εναλλακτικά test πολλαπλών συγκρίσεων.

π.χ. Duncan, Scheffé, Newman-Keuls, Sidak, Dunnett κ.α.

Βασικό κριτήριο επιλογής: το είδος των συγκρίσεων που θα γίνουν και το πείραμα που ακολουθεί η μελέτη.

ANOVA - Γραμμικό Μοντέλο

Ομάδες (groups)

	1	2	3	...	k	
	X_{11}	X_{12}	X_{13}	...	X_{1k}	
	X_{21}	X_{22}	X_{23}	...	X_{2k}	
	\vdots	\vdots	\vdots	\vdots	\vdots	
	$X_{n_1 1}$	$X_{n_2 2}$	$X_{n_3 3}$...	$X_{n_k k}$	
Αθροισμα	$T_{.1}$	$T_{.2}$	$T_{.3}$...	$T_{.k}$	$T_{..}$
Μέση τιμή	$\bar{X}_{.1}$	$\bar{X}_{.2}$	$\bar{X}_{.3}$...	$\bar{X}_{.k}$	$\bar{X}_{..}$

Οπου

$$\bar{X}_{.j} = \frac{T_{.j}}{n_j}$$

$$T_{..} = \sum_{j=1}^k T_{.j} = \sum_{j=1}^k \sum_{i=1}^{n_j} X_{ij}$$

$$\bar{X}_{..} = \frac{T_{..}}{N}, \quad N = \sum_{j=1}^k n_j$$

$$X_{ij} = \mu_j + \epsilon_{ij}, \quad \epsilon_{ij} = X_{ij} - \mu_j = \text{εντρον (σφάλμα)}$$

$$\mu = \frac{\sum_{j=1}^k \mu_j}{k} = \text{συνολική μέση τιμή όλων των παρατηρήσεων}$$

$$\mu_j = \mu + \tau_j, \quad \tau_j = \mu_j - \mu = \text{group effect}$$

$$X_{ij} = \mu + \tau_j + \epsilon_{ij}$$

Οι υποθέσεις της ANOVA (ανάλυσης διακύμανσης) επαναδιατυπώνονται ως εξής:

α) k δεδομένα: τυχαία και ανεξάρτητα δείγματα απίθους αντίστοιχου πληθυσμού με $\sigma_1^2 = \sigma_2^2 = \dots = \sigma_k^2 = \sigma^2$

β) τ_j : άγνωστες με $\sum_{j=1}^k \tau_j = 0$

$$\left(\tau_j = \mu_j - \mu \Rightarrow \sum_{j=1}^k \tau_j = \sum_{j=1}^k (\mu_j - \mu) = \sum_{j=1}^k \mu_j - \frac{\sum_{j=1}^k \mu_j}{k} = 0 \right)$$

γ) $\varepsilon_{ij} = x_{ij} - \mu_j \Rightarrow$

ε_{ij} : ανεξάρτητα, κανονική κατανομή $\left(\sum_{j=1}^k \varepsilon_{ij} = \sum_{j=1}^k (x_{ij} - \mu_j) \right)$
 με μέση τιμή 0 και διακύμανση ίση με τη διακύμανση των x_{ij} , εφόσον x_{ij} και ε_{ij} διαφέρουν κατά μια σταθερά

\Rightarrow Οι υποθέσεις της ανάλυσης $H_0: \mu_1 = \mu_2 = \dots = \mu_k$

H_A : Δύο τουλάχιστον οπότες μ_j διαφέρουν μεταξύ τους

δίνονται:

$$H_0: \tau_j = 0, \quad j=1, 2, \dots, k$$

$$H_A: \text{ένα τουλάχιστον } \tau_j \neq 0$$

$$x_{ij} = \mu + \tau_j + \varepsilon_{ij} \Rightarrow x_{ij} - \mu = \underbrace{\tau_j}_{\mu_j - \mu} + \underbrace{\varepsilon_{ij}}_{x_{ij} - \mu_j}$$

$$\Rightarrow x_{ij} - \mu = (\mu_j - \mu) + (x_{ij} - \mu_j), \quad \begin{matrix} i = 1, 2, \dots, n_j \\ j = 1, 2, \dots, k \end{matrix} \quad [1]$$

Προσχή: ο πληθυσμιακός μέσος μ εκτιμάται από το δειγματικό \bar{x} , η διαφορά $x_{1j} - \mu$ εκτιμάται από την ποσότητα $x_{1j} - \bar{x}_{..}$. Επίσης η διαφορά $(\mu_j - \mu)$ εκτιμάται από την $(\bar{x}_{.j} - \bar{x}_{..})$ και η $(x_{1j} - \mu_j)$ από την $(x_{1j} - \bar{x}_{.j})$

Άρα η (1) γίνεται

$$x_{1j} - \bar{x}_{..} = (\bar{x}_{.j} - \bar{x}_{..}) + (x_{1j} - \bar{x}_{.j}) \quad (2)$$

Παίρνοντας τα αθροίσματα των τετραγώνων έχουμε

NOTE: TO BE
PROVEN BY
STUDENTS

[see by Γωστό!]

$$\begin{aligned} \sum_{j=1}^k \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_{..})^2 &= \sum_{j=1}^k \sum_{i=1}^{n_j} \left[(\bar{x}_{.j} - \bar{x}_{..}) + (x_{ij} - \bar{x}_{.j}) \right]^2 \\ &= \sum_{j=1}^k \sum_{i=1}^{n_j} (\bar{x}_{.j} - \bar{x}_{..})^2 + \sum_{j=1}^k \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_{.j})^2 + \underbrace{2 \sum_{j=1}^k \left[(\bar{x}_{.j} - \bar{x}_{..}) \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_{.j}) \right]}_{(A)} \end{aligned} \quad (3)$$

Αλλά από τον τύπο $\sum_{i=1}^{n_j} (x_{ij} - \bar{x}_{.j}) = 0 \Rightarrow (A) = 0$

για $j = 1, 2, \dots, k$.

και $\sum_{j=1}^k \sum_{i=1}^{n_j} (\bar{x}_{.j} - \bar{x}_{..})^2 = \sum_{j=1}^k n_j (\bar{x}_{.j} - \bar{x}_{..})^2$

Άρα (3) $\Rightarrow \sum_{j=1}^k \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_{..})^2 = \sum_{j=1}^k n_j (\bar{x}_{.j} - \bar{x}_{..})^2 + \sum_{j=1}^k \sum_{i=1}^{n_j} (x_{ij} - \bar{x}_{.j})^2$

(4)

L

Τελικά η (4) είναι η γνωστή μας ιδιότητα.

$$SS_T = SS_B + SS_W$$

Note : $SS_W \rightarrow$ συνδέεται με την $(X_{ij} - \bar{X}_{.j})$ ή την $(X_{ij} - \mu_j)$

δηλαδή \rightarrow ειναι \rightarrow άρα αποτελεί ένα μέτρο της διασποράς των παρατηρήσεων.
(error sum of squares)

και $SS_B \rightarrow$ (error mean square)

Υπολογιστικά :

$$SS_T = \sum_{j=1}^k \sum_{i=1}^{n_j} X_{ij}^2 - \frac{T_{.j}^2}{n_j}$$

$$SS_B = \sum_{j=1}^k \frac{T_{.j}^2}{n_j} - \frac{T_{..}^2}{n}$$

$$SS_W = \sum_{j=1}^k \sum_{i=1}^{n_j} X_{ij}^2 - \sum_{j=1}^k \frac{T_{.j}^2}{n_j}$$

LECTURE 10 (19/04/10)

Ασκήσεις τ : Απορίες? - Τι είναι PDF, CDF

Παρατηρήσεις: Αναφορά!

Είδαμε: ANOVA με τ μεταλλή - Γραμμικά Υπόδειγμα

- Γιατί κάνουμε ANOVA?

- Τι είναι \Rightarrow

* Στην ANOVA με τ μεταλλή σκοπός είναι η σύγκριση μέσων τιμών μιας τυχαίας μεταλλής σε k πληθυσμιακές ομάδες

Παραδοχή: Όλοι οι άλλοι παράγοντες οι οποίοι μπορούν να επηρεάσουν στη μεταλλότητα των τιμών στο εσωτερικό των ομάδων, είναι ελεγχόμενοι.

Όταν η παραπάνω παραδοχή δεν ισχύει μπορούμε να ελέγξουμε την επίδραση ενός δεύτερου παράγοντα, με τη σκέδαση τυχαιοποιημένων μπλοκ (randomized block design)

(Fisher 1925)

Γενικεύσεις

- Block: οι παρατηρήσεις χωρίζονται σε n block με τέτοιο τρόπο ώστε μέσα σε κάθε block οι μονάδες να είναι όσο το δυνατόν περισσότερο ομοιογενείς.

- Ο αριθμός των παραγόντων που ελέγχονται > 1

Παράδειγμα: Σε μια ιατρική μελέτη μελετάται η σχέση της πυκνότητας χοληστερόλης στο αίμα σε γυναίκες αναφορικά με την ηλικία και τις καπνιστικές συνήθειες

Παράγοντας B Ηλικία (π.κ. $\Rightarrow n < 50$ ετών)

Παράγοντας A Κάπνισμα

μη καπνιστριες
 καπνιστριες (< 10 τσιγάρα/μέρα)
 καπνιστριες (> 10 \rightarrow)

π.κ. $n = 100$: n_1 (ηλικία < 40) = 50
 n_2 50
 n_3 50

$n_1 = 50$
 $n_2 = 50$
 $n_3 = 50$

Παράγοντας A

Παράγοντας B

< 40 $40 \leq 60$ > 60

(μη καπνιστριες)

καπνιστριες < 10 τσιγ/μέρα

καπνιστριες > 10 τσιγ/μέρα

ANOVA με 2 Μεταβλητές Παράγοντες (Factors) (Generalized Randomized Block Design)

Παράγοντας A	Παράγοντας B Group			Σύνολο	Μέση Τιμή	
	1	2	...			k
1	X_{111}	X_{121}	...	X_{1k1}	$T_{1..}$	$\bar{X}_{1..}$
⋮	X_{112}	X_{122}	...	X_{1k2}		
	X_{11n}	X_{12n}	...	X_{1kn}		
2	X_{211}	X_{221}	...	X_{2k1}	$T_{2..}$	$\bar{X}_{2..}$
⋮	X_{212}	X_{222}	...	X_{2k2}		
	X_{21n}	X_{22n}	...	X_{2kn}		
⋮	⋮	⋮	⋮	⋮	⋮	⋮
m	X_{m11}	X_{m21}	...	X_{mk1}	$T_{m..}$	$\bar{X}_{m..}$
⋮	X_{m12}	X_{m22}	...	X_{mk2}		
	X_{m1n}	X_{m2n}	...	X_{mkn}		
Σύνολο	$T_{.1.}$	$T_{.2.}$...	$T_{.k.}$	$T_{...}$	$\bar{X}_{...}$
Μέση Τιμή	$\bar{X}_{.1.}$	$\bar{X}_{.2.}$...	$\bar{X}_{.k.}$		$\bar{X}_{...}$

→ 2 μεταβλητές με αριθμό κατηγοριών m και k αντίστοιχα.

• Two-way Factorial Design.

* ΥΠΟΘΕΣΗ ΖΩΟΜΕ

$$T_{ij} = \sum_{b=1}^n X_{ijb}$$

$$\bar{X}_{ij} = \frac{T_{ij}}{n}$$

$$i = 1, 2, \dots, m$$

$$j = 1, 2, \dots, k$$

$$X_{ijb} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \epsilon_{ijb} \quad (1)$$

μ = συνολική μέση τιμή (grand mean)

α_i = επίδραση του επιπέδου i του παράγοντα A στη τιμή X_{ijb}

$$\left(\bar{X}_{i..} - \bar{X}_{...} \right) \quad (2a)$$

β_j >> >> >> >> >> >> >> X_{ijb}

$$\left(\bar{X}_{.j.} - \bar{X}_{...} \right) \quad (2b)$$

$(\alpha\beta)_{ij}$: αλληλεπίδραση των επιπέδων i, j >> >> X_{ijb}

$$\begin{aligned} & \left[\bar{X}_{ijs.} - \bar{X}_{i..} - \bar{X}_{.j.} + \bar{X}_{...} \right] = \left[\bar{X}_{ijs.} - (\bar{X}_{i..} - \bar{X}_{...}) \right. \\ & \quad \left. - (\bar{X}_{.j.} - \bar{X}_{...}) - \bar{X}_{...} \right] \quad (2c) \end{aligned}$$

Είδη: Σφάλμα κατά την εκτίμηση των X_{ijb} , δηλαδή η επίδραση όλων των άλλων εξωγενών παραγόντων οι οποίοι δεν λαμβάνονται υπόψη.

$$\epsilon_{ijb} = X_{ijb} - \mu = (X_{ijb} - \bar{X}_{ij}) \quad (2d)$$

Εξουμε (1) \rightarrow q_a, q_b, q_c, q_d

$$(X_{121} - \mu) = (\bar{X}_{1..} - \bar{X}_{000}) + (\bar{X}_{.2.} - \bar{X}_{000})$$

$$+ (\bar{X}_{12.} - \bar{X}_{1..} - \bar{X}_{.2.} + \bar{X}_{000}) + (X_{121} - \bar{X}_{12.}) \Rightarrow$$

$$\Rightarrow \sum_{i=1}^m \sum_{j=1}^k \sum_{l=1}^n (X_{121} - \bar{X}_{000})^2 = \sum_{i=1}^m \sum_{j=1}^k \sum_{l=1}^n (\bar{X}_{1..} - \bar{X}_{000})^2 + \sum_{i=1}^m \sum_{j=1}^k \sum_{l=1}^n (\bar{X}_{.2.} - \bar{X}_{000})^2$$

$$+ \sum_{i=1}^m \sum_{j=1}^k \sum_{l=1}^n (\bar{X}_{12.} - \bar{X}_{1..} - \bar{X}_{.2.} + \bar{X}_{000})^2 + \sum_{i=1}^m \sum_{j=1}^k \sum_{l=1}^n (X_{121} - \bar{X}_{12.})^2$$

$$= \underbrace{m k \sum_{i=1}^m (\bar{X}_{1..} - \bar{X}_{000})^2}_{SSA} + \underbrace{m m \sum_{j=1}^k (\bar{X}_{.2.} - \bar{X}_{000})^2}_{SSB}$$

$$+ \underbrace{m \sum_{i=1}^m \sum_{j=1}^k (\bar{X}_{12.} - \bar{X}_{1..} - \bar{X}_{.2.} + \bar{X}_{000})^2}_{SSAB} + \underbrace{\sum_{i=1}^m \sum_{j=1}^k \sum_{l=1}^n (X_{121} - \bar{X}_{12.})^2}_{SSE} \quad (3)$$

$TSS =$
Παραλλαγές A
Παραλλαγές B
Αλληλεπίδραση AB
Εσφαλμα

$$TSS = SSA + SSB + SSAB + SSE$$

βαθμοί ελευθ. $(mk(n-1))$ $(m-1)$ $(k-1)$ $(m-1)(k-1)$ $mk(n-1)$

Μέση τετραγωνικά

$$\Rightarrow MSA = \frac{SSA}{(m-1)}, \quad MSB = \frac{SSB}{k-1}, \quad MSAB = \frac{SSAB}{(m-1)(k-1)}, \quad MSE = \frac{SSE}{mk(n-1)}$$

Έλεγχος Υποθέσεων.

α. Επίδραση
Παράγοντας A

$$H_0 : \alpha_i = 0, \quad i = 1, 2, \dots, m$$

$$H_A : \text{τουλάχιστον ένα } \alpha_i \neq 0$$

↳ ελέγχεται η επίδραση του A επί των τιμών χ_{ijk}

$$F = \frac{MSA}{MSE}$$

β. Επίδραση
Παράγοντας B

$$H_0 : \beta_\lambda = 0, \quad \lambda = 1, 2, \dots, k$$

$$H_A : \text{τουλάχιστον ένα } \beta_\lambda \neq 0$$

$$F = \frac{MSB}{MSE}$$

γ. Έλεγχος επίδρασεων αλληλεπίδρασεων A, B

$$H_0 : (\alpha\beta)_{i\lambda} = 0, \quad i = 1, 2, \dots, m, \quad \lambda = 1, 2, \dots, k$$

$$H_A : \text{τουλάχιστον ένα } (\alpha\beta)_{i\lambda} \neq 0$$

$$F = \frac{MSAB}{MSE}$$

Γίνοντας ANOVA (Τυχαίοποιημένο σχέδιο, 2 Παράγοντες)

Γενική Διαστρέβλωση	Άσφαση Τετραγωνάκι	Βαθμοί Ελευθερίας	Μέσο Τετραγωνάκι	F
Παράγοντας A	SSA	$m-1$	$MSA = \frac{SSA}{(m-1)}$	$\frac{MSA}{MSE}$
Παράγοντας B	SSB	$k-1$	$MSB = \frac{SSB}{(k-1)}$	$\frac{MSB}{MSE}$
Αλληλεπίδραση AB	SSAB	$(m-1)(k-1)$	$MSAB = \frac{SSAB}{(m-1)(k-1)}$	$\frac{MSAB}{MSE}$
Εφαδμα	SSE	$mk(m-1)$	$MSE = \frac{SSE}{mk(m-1)}$	
Εύνολο	TSS	$mkk-1$		

LECTURE 11 (26/02/10)

ANOVA - 2 παράγοντες

Παράδειγμα :

Εε μια ιατρική μελέτη εξέτασμε τη σχέση χολινστερόλης στο αίμα σε γυναίκες αναφορικά με την ηλικία και τις καπνιστικές συνήθειες

36 γυναίκες

= Παράγοντας A : καπνισμα { μη καπνιστριες
καπνιστριες I (< 10 τσιγάρα/μέρα)
II (> 10 »)

- Παράγοντας B : ηλικία { ηλικια I (< 40 ετών)
II (40 < ⁴⁰⁻⁶⁰ και < 60 ετών)
III (> 60 ετών)

Παράγοντας A	Παράγοντας B			Αδροσημ	Μέση τιμή
	ηλικία I	ηλικία II	ηλικία III		
μη καπνιστριες	202	210	240		
	226	201	230		
	200	190	290	2520	210
	180	175	210		
Αδροσημ	214	206	200		
Μέση τιμή	203,5	201,5	225		
καπνιστριες I	210	222	245		
	207	239	247	2610	217,5
	190	212	232		
	180	207	221		
Αδροσημ	792	773	945		
Μέση τιμή	198	218,25	236,25		
καπνιστριες II	209	215	260		
	227	203	252	2720	226,67
	210	235	235		
	200	230	250		
Αδροσημ	240	883	997		
Μέση τιμή	210	220,75	249,25		
Αδροσημ	2446	2562	2842	7850	

Υπολογίζουμε όλες τις ποσότητες που χρειάζεστε:

$$TSS = \sum_{i=1}^m \sum_{j=1}^k \sum_{b=1}^n x_{ijb}^2 - \frac{\left(\sum_{i=1}^m \sum_{j=1}^k \sum_{b=1}^n x_{ijb} \right)^2}{m \cdot k \cdot n}$$

$$\left. \begin{array}{l} \text{Έχουμε: } m = 3 \\ \quad \quad \quad k = 3 \\ \quad \quad \quad n = 4 \end{array} \right\} C$$

$$C = \frac{7850^2}{36} = 1711736,111$$

$$\begin{aligned} TSS &= (208^2 + 296^2 + 200^2 + \dots + 250^2) - 1711736,111 = \\ &= \underline{528366 + 572174 + 621010} \\ &\quad 1721550 - 1711736,111 = 9813,889 \end{aligned}$$

$$SSA = \frac{\sum_{i=1}^m T_{i.}^2}{k \cdot n} - C = \frac{2590^2 + 2610^2 + 2790^2}{12} - C = 1672,222$$

$$SSB = \frac{\sum_{j=1}^k T_{.j}^2}{m \cdot n} - C = \frac{2440^2 + 2562^2 + 2842^2}{12} - C = 6907,55$$

$$SSAB = \frac{\sum_{i=1}^m \sum_{j=1}^k T_{ij.}^2}{n} - \frac{\sum_{i=1}^m T_{i.}^2}{k \cdot n} - \frac{\sum_{j=1}^k T_{.j}^2}{m \cdot n} + C$$

$$= \underline{814^2 + 806^2 + 900^2 + 792^2 + 873^2 + 945^2 + 840^2 + 883^2 + 997^2}$$

$$- \frac{2590^2 + 2610^2 + 2790^2}{12} - \frac{2440^2 + 2562^2 + 2842^2}{12} + 1711736,111$$

$$SSA = 1790.987 - 1713,408,333 - 1718643,667 + (711730,111) \\ = 671,111$$

$$SSC = \sum_{i=1}^m \sum_{j=1}^x \sum_{k=1}^n \sum_{l=1}^q x_{ijkl}^2 - \frac{\sum_{i=1}^m \sum_{j=1}^x T_{ij}^2}{n} \\ = 1795550 - 1790.987 = 4563$$

Βασικοί ελευθερίες :

$$\hookrightarrow \text{Συνολικό άθροισμα} \quad m \cdot x \cdot n - 1 = 3 \cdot 3 \cdot 4 - 1 = 35$$

$$\text{Παράγοντα A} \quad m - 1 = 3 - 1 = 2$$

$$\hookrightarrow \text{B} \quad x - 1 = 3 - 1 = 2$$

$$\hookrightarrow \text{Αλληλεπίδραση} \quad (m-1)(x-1) = 4$$

$$\hookrightarrow \text{Εφαρμογών} \quad m \cdot x \cdot (n-1) = 3 \cdot 3 \cdot 3 = 9$$

Μέσα τετραγωνά :

$$MSA = \frac{SSA}{m-1} = 836,111$$

$$MSB = \frac{SSB}{x-1} = 3453,775$$

$$MSAB = \frac{SSAB}{(m-1)(x-1)} = 167,78$$

$$MSE = \frac{SSC}{m \cdot x \cdot (n-1)} = 507$$

Έλεγχος 1

Note : $X_{126} = \mu + \alpha_1 + \beta_2 + (\alpha\beta)_{12} + \epsilon_{126}$

Επίπεδο σημαντικότητας $\alpha = 0,10$

A) Έλεγχος παράγοντα A (κατηγορίας)

$$H_0 : \mu_A^1 = \mu_A^2 = \mu_A^3$$

$$\text{or} : \alpha_1 = \alpha_2 = \alpha_3 = 0$$

$$H_A : \alpha_1 \neq 0 \text{ ή } \alpha_2 \neq 0 \text{ ή } \alpha_3 \neq 0$$

$$F = \frac{MSA}{MSE} = \frac{836,111}{507} = 1,649$$

$$F_{2,9,\alpha} = F_{2,9,0,05} = 4,26 \quad (\alpha = 0,05 = 5\%)$$

$F < F_{2,9,0,1} : H_0 : \text{Accepted}$

B) Έλεγχος παράγοντα B (ηλικία)

$$H_0 : \beta_1 = \beta_2 = \beta_3 = 0$$

$$H_A : \beta_1 \neq 0 \text{ ή } \beta_2 \neq 0 \text{ ή } \beta_3 \neq 0$$

$$F = \frac{MSB}{MSE} = \frac{3453,475}{507} = 6,81$$

$$F_{2,9,\alpha} : F > F_{2,9,\alpha} : H_0 : \text{Απορρίπτεται!}$$

C) Έλεγχος αλληλεπιδράσεων A, B

$$H_0 : (\alpha\beta)_{11} = (\alpha\beta)_{12} = (\alpha\beta)_{13} = \dots = (\alpha\beta)_{33} = 0$$

$$H_A : \text{τουλάχιστον ένα } (\alpha\beta)_{ij} \neq 0$$

$$F = \frac{MSAB}{MSE} = 0,331$$

$$F_{4,9,0,05} = F_{4,9,0,05} = 3,63$$

$$F < F_{4,9,0,1}$$

Note



LECTURE 12 (03/05/10)

Είδοςμε : ΑΝΟΝΥΑ - 2 παράγοντες

- ΕΚΕΙΣΕΣ ΕΠΡΟΙΟΜΑΤΩΝ → ΠΙΝΑΚΑ (j)

$$= \sum_{i=1}^m \sum_{a=1}^k \sum_{b=1}^n (X_{iab} - \bar{X}_{i..})^2$$

$$SST = SSA + SSB + SSAB + SSE$$

$$\sum_{i=1}^m \sum_{a=1}^k \sum_{b=1}^n (X_{iab} - \bar{X}_{i..})^2 \quad \sum_{i=1}^m \sum_{a=1}^k (X_{i..} - \bar{X}_{i..})^2 \quad \sum_{i=1}^m \sum_{a=1}^k \sum_{b=1}^n (X_{iab} - \bar{X}_{i..} - \bar{X}_{.a.} + \bar{X}_{...})^2$$

= Παράδειγμα : Ιατρική μελέτη

Πίνακα (j)

Εφαρμογή ΜΑΤΡΩΝ (π.ρ. ορθογώνιος πίνακας, $m \neq n$) ?

= Ασκήσεις 7

Νοίε : καλύτερα χαρακτηριστικά

κέρυθση:

$$\rho_2 = \frac{m_4}{m_2^2} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^4}{\left[\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \right]^2} = \frac{E(X-\mu)^4}{\sigma^4}$$

$$\text{or } = \frac{m_4}{m_2^2} - 3 = \dots - 3$$

ΜΑΤΡΩΝ :

k : κέρυθση (χ) : χ : vector → + number

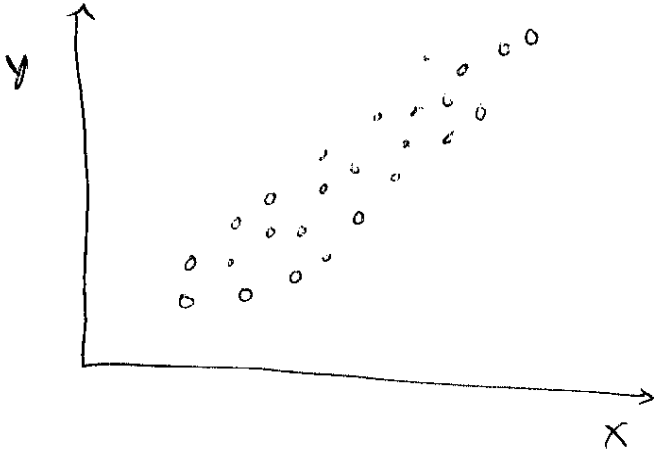
Θεωρητική σημ: Uniform $k = -\gamma/2$
 mod $\chi = 0$

χ : matrix → κέρυθση for each column of X

ΕΥΧΕΤΙΣΗ

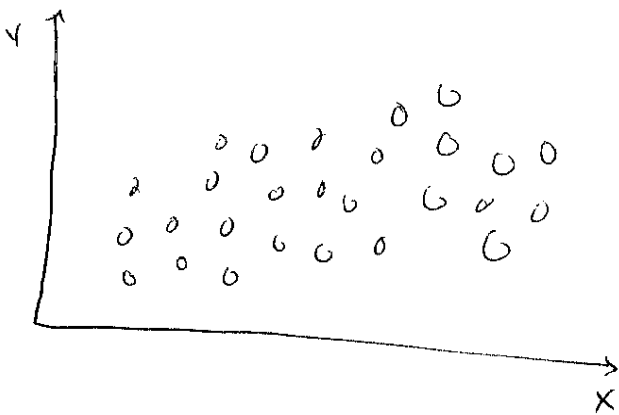
- ΕΧΕΘΕΙΣ ΜΕΤΑΞΙ ΤΥΧΑΙΑΣ ΜΕΤΑΒΛΗΤΩΝ.

π.χ. γραμμική εξάρτηση μιας τ.μ. από μια ή περισσότερες μεταβλητές. X, Y



Διαγράμματα Διασποράς (Scatter Plots)

← Ισχυρή γραμμική σχέση μεταξύ δύο συνεχών τ.μ.



← Αδύναμη γραμμική σχέση.

* ΕΥΧΕΤΙΣΤΗΣ ΕΥΣΧΕΤΙΣΗ *

Δύο τ.μ. X, Y με διασπορές $\sigma_x^2 = \text{Var}[X], \sigma_y^2 = \text{Var}[Y]$

και συνδιασπορά $\sigma_{xy} = \text{Cov}[X, Y] = E[(X - \mu_x)(Y - \mu_y)]$

Note: Συνεχές χώρο.

$$\mu_n = \int_{-\infty}^{\infty} (x - c)^n F(x) dx$$

E : expected value.

$$\mu_n = E(X^n) = \int_{-\infty}^{\infty} x^n dF(x), \quad (E(X) = \int_{-\infty}^{\infty} x \cdot f(x) dx)$$

F : cumulative distribution function x

$$f(x) = \frac{dF(x)}{dx} = \int_{-\infty}^{\infty} x \cdot f(x) dx$$

- Correlation Coefficient.

$$\rho = \frac{\sigma_{XY}}{\sigma_X \sigma_Y} \quad : \text{ ανεξάρτητος μονάδων}$$

ΤΙΜΕΣ ΣΤΟ ΔΙΑΣΤΗΜΑ [-1, 1]

$\rho \rightarrow 1$: ισχυρή θετική συσχέτιση

$\rho \rightarrow -1$: αρνητική \rightarrow

$\rho \rightarrow 0$: γραμμική ανεξαρτησία των X, Y

Ποιοτικός έλεγχος : scatter diagram

* Εκτίμηση του συντελεστή συσχέτισης
(συντελεστής συσχέτισης Pearson)

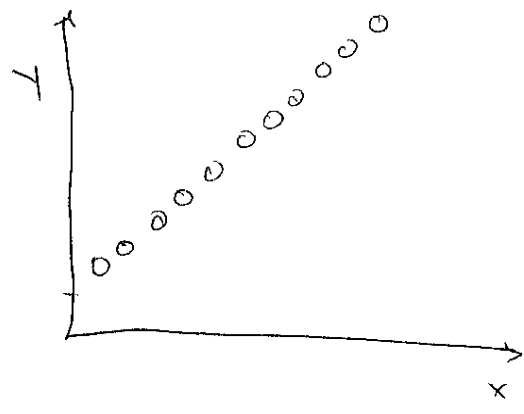
Έστω n ζευγάρια παρατηρήσεων των X, Y

$$\{ (x_1, y_1), (x_2, y_2), \dots, (x_n, y_n) \}$$

$$r = \frac{s_{XY}}{s_X s_Y}, \quad s_{XY} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \frac{1}{n-1} \left(\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y} \right)$$

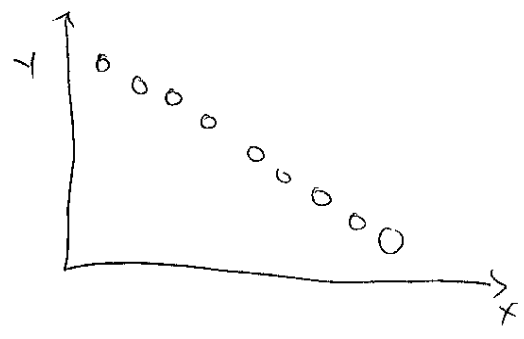
$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}}$$

Παράδειγμα : Scatter plots

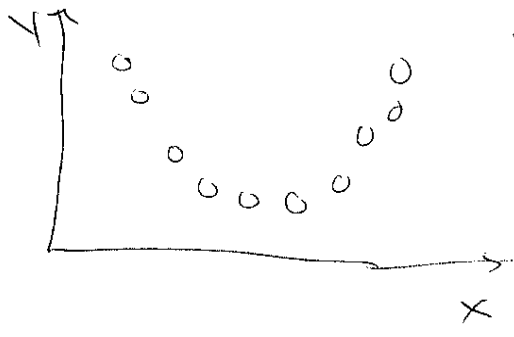
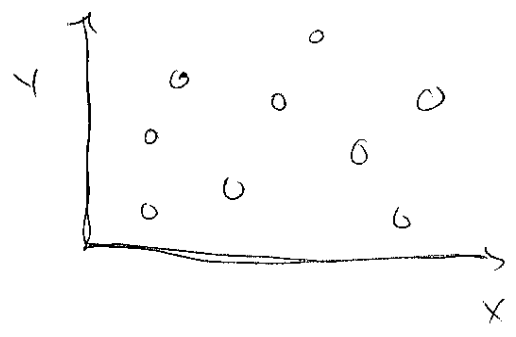


: Πλήρως θετική συσχέτιση $r=1$

$$Y = Y_0 + \alpha X$$



: Πλήρως αρνητική συσχέτιση $r=-1$



: Αυστηρά ίδια γραμμική σχέση $r=0$

$$Y = 0$$

Έλεγχος υποθέσεων

18.5
: παρόμοιος με τους έλεγχους

υποθέσεων που έχουμε αναφέρει

π.χ.

$H_0: \rho = 0$ (ελέγχει αν τα X, Y είναι γραμμικά ανεξάρτητα)

Το τυπικό σφάλμα του στατιστικού συντελεστή r εκτιμάται από την ποσότητα

$$\hat{\sigma}_e(r) = \sqrt{\frac{1-r^2}{n-2}}$$

Έλεγχος με το κριτήριο: (για $H_0: \rho = \rho_0$)

$$t = \frac{r - \rho_0}{\sqrt{\frac{(1-r^2)}{(n-2)}}} = r \sqrt{\frac{n-2}{1-r^2}} \quad \left(r = \frac{\sum xy}{\sum x^2 \sum y^2} \right)$$

Όταν τα ζεύγη (x_i, y_i) λαμβάνονται τυχαία από του αντίστοιχα πληθυσμό που X, Y (και η X, Y ακολουθούν κανονική κατανομή) η παραπάνω ποσότητα ακολουθεί κατανομή t με $n-2$ βαθμούς ελευθερίας

Προβλεπή:

Για $H_0: \rho = \rho_0 \neq 0$ ακολουθείται διαφορετική διαδικασία έλεγχου

Μετασχηματίζοντας την τιμή του δειγματικού συντελεστή συσχέτισης ως

$$z_r = \frac{1}{2} \ln \left(\frac{1+r}{1-r} \right) \text{ αποδεικνύεται}$$

(Fisher (1921) Μετρώ, 1, 3-21)
Μετασχηματισμός Z του Fisher

ότι η δειγματική κατανομή της ποσότητας z_r είναι κανονική με μέση τιμή

$$z_p = \frac{1}{2} \ln \left(\frac{1+p}{1-p} \right) \text{ και τυπική απόκλιση } \frac{1}{\sqrt{n-3}}$$

Άρα αν

$$H_0: \rho = \rho_0 \neq 0$$

$$H_A: \rho \neq \rho_0$$

χρησιμοποιείται η $z = \frac{z_r - z_p}{\frac{1}{\sqrt{n-3}}}$ η οποία ακολουθεί τυπική κανονική κατανομή.

LECTURE 13ΑΝΑΛΥΣΗ ΠΑΛΙΝΔΡΟΜΗΤΕΣ (Regression Analysis)

: Περιγράφει τη μεταβολτικότητα μιας τ.μ Y από τιμές άλλων τ.μ. X_1, X_2, \dots

ΕΤΟΧΟΣ: Εύρεση μοντέλου που περιγράφει την εξάρτηση της Y (εξαρτημένη μεταβλητή) από μια ή άλλες μεταβλητές X (ανεξάρτητη μεταβλητή)

Απλή Γραμμική Παλινδρόμηση (Simple Linear Regression)

: Η ανεξάρτητη τ.μ. είναι x και η εξάρτηση θεωρείται γραμμική.

Υποθέσεις:

1. Η μέση τιμή της τ.μ. Y για κάθε x του X , $E\{Y|X=x\}$ είναι γραμμική συνάρτηση της x

$$\mu_{Y|X} \equiv E\{Y|X=x\} = \alpha + \beta x$$

2. Η διασπορά της Y είναι σταθερή $\text{Var}\{Y|X=x\} = \sigma^2_{Y|X} = \sigma^2$

3. Η κατανομή της Y ως προς τη X είναι κανονική

$$Y|X=x \sim N(\alpha + \beta x, \sigma^2)$$

Επιπλέον
παραμετρικό
έξοχο και
επιλογή
των α, β

Γενικά : $Y_i = \alpha + \beta X_i + \epsilon_i$

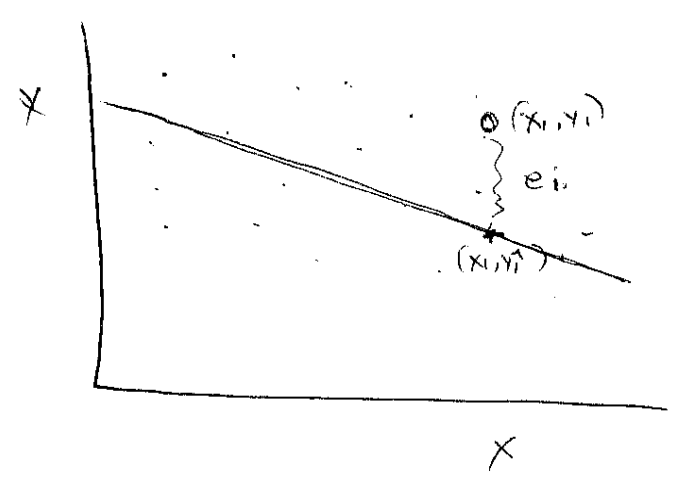
ϵ_i : σφάλμα παλινδρόμησης (regression error)

Var [ϵ_i] = $\sigma_{\epsilon}^2 = \sigma^2$ $\epsilon_i = Y_i - \mu_{Y|X_i}$

Εκτίμηση Παρομιέτρων Απλής Γραμμικής Παλινδρόμησης

Έστω ζεύγη $\{(x_1, y_1), (x_2, y_2) \dots (x_n, y_n)\}$

Μέθοδο Ελάχιστων Τετραγώνων (Method of Least Squares)



Ελαχιστοποίηση των

$$\sum_{i=1}^n \epsilon_i^2 = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta} X_i$$

$$\hat{\beta} = \frac{S_{XY}}{S_X^2} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x} \quad (\text{δίασ } \bar{y} = \hat{\alpha} + \hat{\beta} \bar{x})$$

→ Πληθυσμιακή συνάρτηση της παλινδρόμησης $\mu_{Y|X} = \alpha + \beta X$

Επαγωγική Μέθοδος Εγρήμωσις:

13.3

Τυπικές αποκλίσεις κατανομών

$$SE(\hat{\beta}) = \frac{SYX}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

Δειγματικές αποκλίσεις:

$$SE(\hat{\alpha}) = SYX \sqrt{\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

$$SYX = \sqrt{\frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2}}$$

Τυπική απόκλιση της Παλινδρόμησης
(standard deviation from regression)

$\sum_{i=1}^n (y_i - \hat{y}_i)^2$: άθροισμα των τετραγώνων των σφαλμάτων

$e_i \equiv y_i - \hat{y}_i$: υπολοίπο (residual) - σφάλμα ελαχίστων τετραγώνων

οι $SE^2 \equiv SE^2 = \frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2$: $n-2$ σκέτη από τους βαθμούς ελευθέριας n του δείγματος αφαιρούμε δύο στις τις δύο παραμέτρους που έχουν ήδη εγρήμωσει!

↳ Βασισμένοι στην κρίση της ευθείας $\hat{\beta}$ μπορούμε να ελέγξουμε την τιμή της κρίσης β της πληθυσμιακής ευθείας

δηλαδή να ελέγξουμε την υπόθεση

$$H_0 : \beta = \beta_0$$

$$H_A : \beta \neq \beta_0$$

Έλεγχος με τη βοήθεια της μοσότητας

$$t = \frac{\hat{\beta} - \beta_0}{SE(\hat{\beta})} \quad : \text{κατανομή } t \text{ με } n-2 \text{ βαθμ.}$$

$$\text{ή } t = \frac{\hat{\alpha} - \alpha_0}{SE(\hat{\alpha})}$$

π.χ. $\beta_0 = 0$ $μ_{Y|X} = \alpha + \beta X = \alpha$. (η γ δεν εξαρτάται γραμμικά από την X)

Διαστήματα εμπιστοσύνης για τις εκτιμώμενες παραμέτρους:

(α, β)

$$\hat{\alpha} \pm t_{n-2, 1-\alpha/2} \times SE(\hat{\alpha})$$

$$\hat{\beta} \pm t_{n-2, 1-\alpha/2} \times SE(\hat{\beta})$$

} $\alpha\%$ διάστημα εμπιστοσύνης (ΔΕ) των συντελεστών α, β της ευθείας της παλινδρόμησης.

Εκτίμηση - πρόβλεψη εξαρτημένης μεταβλητής:

Εμπειρική εκτίμηση \hat{y}_0 της μέγιστης τιμής της Y για κάθε τιμή x_0 της X . Το $(1-\alpha)\%$ Δ.Ε είναι

$$\hat{y} \pm t_{n-2, \alpha/2} SE(\hat{y}) \equiv (\alpha + \beta x_0) \pm t_{n-2, \alpha/2} SY/X \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

τυπικό σφάλμα της \hat{y}

$$SE(\hat{y}) = SY/X \sqrt{\frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}$$

Όρια της προβλεψής της Y για μια τιμή X_0 ;
(limits of prediction)

↳ \tilde{Y} : προβλεπόμενη τιμή νέου ατόμου-δεδομένου που προστίθεται στον πληθυσμό.

Πρέπει να συνολοδογήσουμε μια επιπλέον μεταβλητότητα λόγω της διασποράς των τιμών της Y γύρω από τη μέση τιμή του υπο-πληθυσμού στον οποίο ανήκουν.

$$\Rightarrow se(\tilde{Y}) = \sqrt{S^2_{Y|X} se(\tilde{Y})^2} = S_{Y|X} \sqrt{1 + \frac{1}{n} + \frac{(X_0 - \bar{X})^2}{\sum_{i=1}^n (X_i - \bar{X})^2}}$$

Αξιολόγηση (υ)είας Μαθηματικής - Ευντελέστης Προσδιορισμού

$$\Rightarrow R^2 = r^2 \quad (\text{coefficient of determination}) \quad \left[r = \frac{S_{XY}}{S_X S_Y} \right]$$

$$r \in [-1, 1] \Rightarrow R^2 \in [0, 1]$$

$$R^2 = \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = \frac{\text{Άθροισμα τετραγώνων απημ γραμμική σταθμδρόμηση}}{\text{Ευνοητικό Άθροισμα τετραγώνων}}$$

$$r^2 = \frac{\left[\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) \right]^2}{\left[\sum_{i=1}^n (X_i - \bar{X})^2 \right] \left[\sum_{i=1}^n (Y_i - \bar{Y})^2 \right]}$$

Ⓞ : Απόδειξη ότι $R^2 = r^2$; $\left(R^2 \Rightarrow \hat{Y}_i = \hat{\alpha} + \hat{\beta}X_i - \hat{\alpha} - \hat{\beta}\bar{X} \right)$

Επίσης ισχύει

$$\sum_{i=1}^n (y_i - \bar{y})^2 = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2 + \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Έχουμε

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

Διορθωμένος συντελεστής προσδιορισμού (adjusted coefficient of determination)

$$\bar{R}^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2 / (n-2)}{\sum_{i=1}^n (y_i - \bar{y})^2 / (n-1)} = 1 - \frac{S^2_{Y|X}}{S_{Y^2}}$$

Έλεγχος υπόθεσης:

$$H_0: \beta = 0$$

$$H_A: \beta \neq 0$$

Με τη βοήθεια του λόγου $F = \frac{MSR}{MSE}$

$$MSR: \text{Μέσο τετράγωνο της παλινδρόμησης} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{k}$$

όπου $k = (p - 1)$
 είναι ελευθέρια βαθμύτητα παλινδρόμησης

$$MSE: \text{Μέσο τετράγωνο των σφαλμάτων} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n - 2}$$

F κατανομή με $k, n-2$ dof.

Μετασχηματισμοί των Μεταληπιών

Εσχύει η σχέση τ.μ. X, Y δεν είναι γραμμική.

↳ Με τη βοήθεια κατάλληλου μετασχηματισμού μπορεί να μετατραπεί σε γραμμική.

π.χ. $Y_t = Y_0 e^{rt} \rightarrow \tau.μ. Y, t \rightarrow$

$$\Rightarrow \ln Y_t = \ln(Y_0 e^{rt}) = \ln Y_0 + rt$$

$$Y' = \alpha' + \beta t$$

π.χ. $Y = a + bx^2 \rightarrow x' = \sqrt{x} \Rightarrow Y = a + bX'$
 $\sqrt{Y} = \sqrt{a + bx^2} \Rightarrow Y' = a' + b'X'$

- Γενικά χρειαζονται Μη Γραμμικά Μοντέλα (Non linear Models)

- Τα παραπάνω ονομάζονται και Γενικευμένα Γραμμικά Μοντέλα (Generalized linear models)

π.χ. $\log Y = b_1 + b_2 W$

$$\text{α. } \log\left(\frac{P}{1-P}\right) = b_1 + b_2 W \Rightarrow P = (1-P) e^{b_1 + b_2 W} = e^{b_1 + b_2 W} - P e^{b_1 + b_2 W}$$

$$\Rightarrow P (1 + e^{b_1 + b_2 W}) = e^{b_1 + b_2 W} \Rightarrow P = \frac{e^{b_1 + b_2 W}}{1 + e^{b_1 + b_2 W}} = \frac{1}{1 + e^{-b_1 - b_2 W}}$$

APPLIED STATISTICS

LECTURE 14

Πολλαπλή Γραμμική Παλινδρόμηση (Multiple Linear Regression)

Η ζ.μ. Y εξαρτάται γραμμικά από τις X_1, X_2, \dots, X_k

$$Y_i = \alpha_i + \beta_1 X_{1i} + \beta_2 X_{2i} + \dots + \beta_k X_{ki} + \epsilon_i$$

(β₀)

$i = 1, 2, \dots, n$

Υποθέσεις = όπως και στην απλή γραμμική παλινδρόμηση

π.χ. : X_1 : ηλικία

X_2 : ύψος

X_3 : βάρος

Y : ενεργειακή πρόσληψη ενηλίκων. $= \alpha + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3$

Ξεχωρή πίνακες

$$\underline{Y} = \underline{X} \underline{B} + \underline{\epsilon}$$

\underline{X} : dimension $n \times (k+1)$

\underline{Y} : $n \times 1$
 \underline{B} : $(k+1) \times 1$
 $\underline{\epsilon}$: $n \times 1$

$$\hookrightarrow E[\underline{\epsilon}] = 0$$

$$V(\underline{\epsilon}) = \sigma^2 \underline{I}$$

Λύση ελαχίστων τετραγώνων

$$\hat{\underline{B}} = (\underline{X}^T \underline{X})^{-1} \underline{X}^T \underline{Y}$$

(minimizes error)
 $\underline{\epsilon}^T \underline{\epsilon}$

$$\epsilon_i = Y_i - \hat{Y}_i \quad i = 1, \dots, n \quad : \text{ψηφίομα (residuals)}$$

Εξίσωση Παλινδρόμησης

$$\hat{Y} = \hat{\alpha} + \hat{\beta}_1 X_1 + \hat{\beta}_2 X_2 + \dots + \hat{\beta}_k X_k$$

→ Για περισσότερες από 2 ανεξ. μ. απαιτείται software.

Έλεγχος Υποθέσεων:

$$H_0: \beta_i = \beta_{i0}, \quad i = 1, 2, \dots, k$$

$H_A: \beta_i \neq \beta_{i0}$ (υποθέτουμε ότι οι τιμές των υποθέσεων ανεξ. μ. ~~μ~~ $\neq \chi^2$ παραμένουν σταθερές)

Κρίσιμη ποσότητα:

$$t = \frac{\hat{\beta}_i - \beta_{i0}}{se(\hat{\beta}_i)}$$

$$t = \frac{\hat{\alpha} - \alpha_0}{se(\hat{\alpha})}$$

Χρήση κατηγορικών μεταλλήτων

Κωδικοποίηση κατηγορικών μεταλλήτων (ψευδομεταλλήτων)

π.χ. Δίτιμες μεταλλήτες $\rightarrow (0, 1)$

Παράδειγμα: ενεργειακή πρόβλεψη.

X_0 : ώλο $\left(\begin{array}{l} 0 \text{ γυναίκες} \\ 1 \text{ άνδρες} \end{array} \right)$

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \beta_4 X_4$$

$\rightarrow \dots \in \beta_i \quad i=0, 1, 2, 3, 4.$

Εισαγωγή Όρων Αλληλεπίδρασης - Μη Γραμμικοί Όροι

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots \quad (\text{Γραμμικοί Όροι})$$

$$+ \beta_{12} X_1 X_2 + \beta_{13} X_1 X_3 + \dots \quad (\text{Όροι Αλληλεπίδρασης})$$

$$+ \beta_{111} X_1^3 + \beta_{222} X_2^3 + \dots \quad (\text{Όροι 2ης τάξης})$$

$$+ \beta_{1111} X_1^4 + \beta_{2222} X_2^4 + \dots \quad (\text{3ης τάξης})$$

⋮

⋮

APPLIED STATISTICS

LECTURE 14-15 (17/05/10)

ΛΟΓΑΡΙΘΜΙΚΗ ΠΑΛΙΝΔΡΟΜΗΣΗ (LOGISTIC REGRESSION)

Είδαμε την γραμμική παλινδρόμηση

$$\mu_{y/x_1, x_2, \dots, x_k} = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k$$

εγκρίνεται με τη βοήθεια γραμμικού υποδείγματος

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_k x_k + \varepsilon$$

Διζήτη εξορμημένη μεταλητική (π.χ. $y = 0 \text{ ή } 1$)
αποτυχία επιτυχία

$\Rightarrow \mu_y = P = P(Y=1)$: πιθανότητα P της επιτυχίας

$P = \alpha + \beta x$: Δεν είναι κατάλληλο $\rightarrow \alpha + \beta x \in \mathbb{R}$
 $P \in [0, 1]$

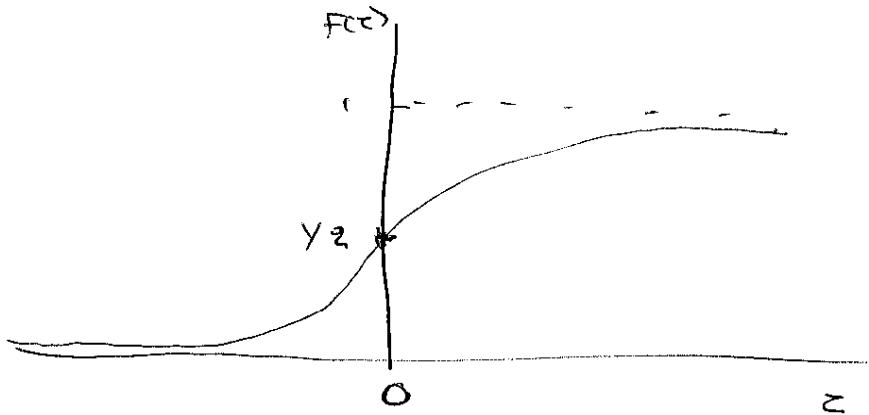
Λύση
 $\Rightarrow \frac{P}{1-P} \xrightarrow{\text{λογάριθμος}} \ln \left[\frac{P}{1-P} \right] = \alpha + \beta x$
 $\hookrightarrow \in [0, +\infty]$ $\hookrightarrow \in [-\infty, +\infty]$

$$\text{logit}(P) \equiv \ln \left[\frac{P}{1-P} \right] = \alpha + \beta x$$

$$\Rightarrow \frac{p}{1-p} = e^{\alpha + \beta x} = e^z \quad \Rightarrow \quad p = e^z - p e^z \Rightarrow$$

$$\Rightarrow p = \frac{e^z}{1 + e^z} = \frac{1}{1 + e^{-z}} = f(z)$$

: εντήμενη πιθανότητα επιτυχίας p της δίκης S τ.μ. Y .



$$f(z) \in [0, 1]$$

$$z \in [-\infty, +\infty]$$

Γενική μορφή

$$\ln \left\{ \frac{p}{1-p} \right\} = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k$$

$$\Rightarrow \frac{p}{1-p} = e^{\beta_0 + \beta_1 x_1 + \dots + \beta_k x_k}$$

□

- κατηγορία: Generalized Linear Models

Log-linear models

Κάθε ένας από τους β_i εκφράζει τη μεταβολή του $\ln \left\{ \frac{p}{1-p} \right\}$ (σχετικής πιθανότητας) για μία μονάδα αύξησης της αντίστοιχης τ.μ. x_i .

π.χ. Δίτιμη τ.μ. φύλο

$$\ln\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 X = \beta_0 + \beta_1 (\text{φύλο})$$

$X = 0$: γυναίκες
 $X = 1$: άνδρες

$$\frac{P_A}{1-P_A} = e^{\beta_0 + \beta_1}$$

$$\frac{P_T}{1-P_T} = e^{\beta_0}$$

$$\psi = \frac{P_A/(1-P_A)}{P_T/(1-P_T)} = e^{\beta_1}$$

π
 λόγος σχετικών πιθανοτήτων
 τις εμπνοχία >

π.χ. Πιθανότητα τροχαίου

Accident Analysis and Prevention
 37, 815-825 (2000)
 Chlidoutelis et al.

198 οδμήδου >

$X_1 = 0$: γυναίκες
 $X_1 = 1$: άνδρες >

$$\ln\left(\frac{P}{1-P}\right) = 0,488 - 0,669 X_1$$

$$\rightarrow \psi = e^{-0,669} = 0,51$$

⇒ Σχετική πιθανότητα τροχαίου = 0,51

άνδρες είναι 0,51 φορές μικρότερη απέναντι στις γυναίκες

Επιμέτρηση συντελεστών με τη χρήση μιας

Επιάρθνης Πιθανοφάνειας (Likelihood Function)

$$L(\beta_0, \beta_1, \dots, \beta_k) = \prod_{i=1}^n p_i^{y_i} (1-p_i)^{1-y_i}$$

y_i είναι η τιμή (0, ή 1) της Y στα n παρατηρήσιμα $i=1, 2, \dots, n$

p_i : πιθανότητα επιτυχίας της Y >> >>

$$p_i = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik})}}, \quad i = 1, 2, \dots, n$$

L : πιθανότητα να προκύψουν οι τιμές $y_i, i=1, 2, \dots, n$

(Υπενθυμίζουμε $P(Y=y) = p^y (1-p)^{1-y}$ διωνομική μεταλητική)
πιθανότητα να προκύψει μια απίτητη τιμή)

= Αποδεικνύεται ότι ένας επιτιμώμενος συντελεστής β_1 (στοιχείο με
μεγάλο δείγμα) ακολουθεί κ.κ με μέση τιμή β_1 και τυπικό
σφάλμα $se(\hat{\beta}_1)$ $\frac{\hat{\beta}_1 - \beta_1}{se(\hat{\beta}_1)} : \text{επιμή κανονική κατανομή (ε.κ.κ)}$

Υπόθεση υποθέσεων:

$$H_0 : \beta_1 = 0 \quad (\text{ή } e^{\beta_1} = 1)$$

$$H_A : \beta_1 \neq 0$$

με τη βοήθεια του κριτηρίου

$$z = \frac{B_1^{\wedge}}{\text{se}(B_1^{\wedge})}$$

: κριτήριο του Wald

⇒ φρόσον ισχύει η H_0 το z ακολουθεί $z \sim \chi_1$

Μεγάλα δείγματα

- ΜΕΓΑΛΑ ΔΕΙΓΜΑΤΑ: Όταν η B_1^{\wedge} είναι πολύ μεγάλη $\text{se}(B_1^{\wedge})$ είναι αναλογικά πολύ μικρότερη →

→ Likelihood ratio statistic (λόγος μέγιστων τιμών συνάρτησης πιθανοφάνειας)

$$-2 \ln \left(\frac{L_1^{\wedge}}{L_0^{\wedge}} \right) = -2 \ln L_1^{\wedge} - (-2 \ln L_0^{\wedge})$$

L_1^{\wedge} : ΜΟΛΧ της L όταν ο B_1 δεν περιλαμβάνεται στο υπόδειγμα (δηλ η αντίστοιχη X_i)

L_0^{\wedge} : $\gg \gg L \gg \gg 0$ B_1 περιλαμβάνεται $\gg \gg$
($\gg \gg$)

=

Μεγιστοποίηση της $L \rightarrow$ χη1 ανεξάρτητες $B_0, B_1 \dots B_k$.

Μινιμίζουσα τις μερικές παραγόμενες ως προς $B_0, B_1 \dots B_k$ της $\ln L$

$L, \ln L$: μεγιστοποιούνται για τις ίδιες τιμές των παραγόμενων →

και εξίσωση ή χη1 άγνωστων → (μερική αλγεβράς)

Διάστημα εμπιστοσύνης

100 (1-α)% Δ.Ε. για τον β₁

$$\beta_1 \pm z_{\alpha/2} \text{se}(\hat{\beta}_1)$$

$$\alpha \left[\hat{\beta}_1 \pm z_{\alpha/2} \text{se}(\hat{\beta}_1) \right]$$

π.χ. πιθανότητα τροχαίου ατυχήματος

$$\ln \left[\frac{p}{1-p} \right] = 0,488 - 0,662x_1$$

$$\text{se}(\hat{\beta}_1) = 0,309$$

$$H_0: \beta_1 = 0.$$

$$\frac{\hat{\beta}_1}{\text{se}(\hat{\beta}_1)} = -2,192 \rightarrow$$

< κρίσιμη τιμή ε. κ. κ.

που αντιστοιχεί σε α = 0,05

η H₀: απορρίπτεται

→ Άρα H₁: β₁ ≠ 0.

$$95\% \text{ Δ.Ε. } \left(e^{\left[\hat{\beta}_1 - z_{\alpha/2} \text{se}(\hat{\beta}_1) \right]}, e^{\left[\hat{\beta}_1 + z_{\alpha/2} \text{se}(\hat{\beta}_1) \right]} \right) = (0,285, 0,933)$$