



ΠΑΝΕΠΙΣΤΗΜΙΟ ΚΡΗΤΗΣ

ΣΧΟΛΗ ΘΕΤΙΚΩΝ ΚΑΙ ΤΕΧΝΟΛΟΓΙΚΩΝ
ΕΠΙΣΤΗΜΩΝ

ΤΜΗΜΑ ΜΑΘΗΜΑΤΙΚΩΝ

&

ΣΧΟΛΗ ΕΠΙΣΤΗΜΩΝ ΥΓΕΙΑΣ

ΤΜΗΜΑ ΙΑΤΡΙΚΗΣ

ΠΤΥΧΙΑΚΗ ΕΡΓΑΣΙΑ

«Έλεγχος συσχετίσεων και ανάλυση λογιστικής παλινδρόμησης μελώντας τη σχέση ενεργητικής ανοσοποίησης (εμβολιασμών) και νευροανάπτυξης σε παιδιά προσχολικής ηλικίας στα πλαίσια της μελέτης μητέρας παιδιού ΡΕΑ»

ΙΩΑΝΝΗΣ ΒΑΡΔΑΞΗΣ

ΑΜ 3595

Υπεύθυνος Καθηγητής/Ερευνητής Εργαστηρίου Υποδοχής: Λήδα Χατζή, Λέκτορας
Επιδημιολογίας Διατροφής Τμήματος Ιατρικής Πανεπιστημίου Κρήτης

Υπεύθυνο μέλος ΔΕΠ Τμήματος Μαθηματικών Κλωνιάς Βασίλειος, Αναπληρωτής
Καθηγητής στον Τομέα Εφαρμοσμένων Μαθηματικών και Στατιστικής Πανεπιστημίου
Κρήτης

ΗΡΑΚΛΕΙΟ, ΚΡΗΤΗ

ΑΠΡΙΛΙΟΣ 2012

Σε αυτό το σημείο θα ήθελα να ευχαριστήσω τον υπεύθυνο της εργασίας μου, την κυρία Λήδα Χατζή, λέκτορα επιδημιολογίας και διατροφής του τμήματος ιατρικής Πανεπιστημίου Κρήτης, που με δέχτηκε να κάνω την πτυχιακή μαζί της. Καθώς και την κυρία Ελένη Φθενού και τον κύριο Εμμανουήλ Μπαγκέρη, ερευνητές στο τμήμα επιδημιολογίας και διατροφής του τμήματος ιατρικής Πανεπιστημίου Κρήτης, για την βοήθεια που μου προσέφεραν.

Αφιερωμένο στην οικογένεια μου και στους φίλους μου, σε εκείνους που είναι εδώ και σε όσους έχουν φύγει.

Περιεχόμενα

Κεφάλαιο 1 Εισαγωγή.....	5
1.1 Στατιστική μεθοδολογία	5
1.2 Ανοσοποίηση.....	6
1.3 Νευροανάπτυξη	7
1.4 Νευροανάπτυξη και εμβολιασμός.....	8
Κεφάλαιο 2 Στόχος Διπλωματικής Εργασίας.....	11
Κεφάλαιο 3 Μεθοδολογία Διεξαγωγής Έρευνας	12
3.1 Θεμελιώδη Θεωρήματα Και Έννοιες	12
3.1.1 Θεμελιώδη Θεωρήματα:	12
3.1.2 Σημαντικές Έννοιες:	14
3.2 Παρουσίαση Στατιστικών Μεθόδων.....	17
3.2.1 Έλεγχος χ^2 (Pearson's chi square test):	17
3.2.2 Έλεγχος t (t-test):	22
3.2.3 P-value:	28
3.2.4 Γραμμική Παλινδρόμηση	29
3.2.5 Απλή Γραμμική Παλινδρόμηση:	30
3.2.6 Έλεγχοι Υποθέσεων Με Βάση Την Παλινδρόμηση:	32
3.2.7 Απλή Γραμμική Συσχέτιση (συντελεστής συσχέτισης):	33
3.2.8 Λογιστική παλινδρόμηση:	35
3.2.9 Έλεγχοι Υποθέσεων Με Βάση Την Λογιστική Παλινδρόμηση:	37
3.3 Επιλογή κατάλληλης στατιστικής μεθοδολογίας	38
Κεφάλαιο 4 Συζήτηση:	43
Πηγές:.....	44

Κεφάλαιο 1

Εισαγωγή

1.1 Στατιστική μεθοδολογία

Στην στατιστική, η επιλογή του κατάλληλου στατιστικού test για την ανάλυση των δεδομένων είναι πολύ σημαντική. Καλούμαστε να επιλέξουμε ανάμεσα σε δύο οικογένειες στατιστικών test, την *παραμετρική* και την *μη παραμετρική*. Στην παραμετρική οικογένεια ανήκουν εκείνα τα test τα οποία βασίζονται στην υπόθεση ότι τα δεδομένα προέρχονται από την κανονική κατανομή (Gaussian distribution). Παραμετρικά test είναι το chi-squared test, το t-test, η απλή γραμμική και η λογιστική παλινδρόμηση. Η ανάλυση με βάση αυτά τα test και η επιλογή του καταλληλότερου γίνεται αναλόγως τα δεδομένα και τον σκοπό της έρευνας. Στη μη παραμετρική οικογένεια ανήκουν εκείνα τα test τα οποία δεν βασίζονται σε κάποια συγκεκριμένη υπόθεση για την κατανομή των δεδομένων. Ο συντελεστής συσχέτισης του Spearman, το Mann-Whitney U-test και το Wilcoxon signed rank-test ανήκουν στη κατηγορία των μη παραμετρικών. Στην παρούσα εργασία δεν θα ασχοληθούμε με μη παραμετρικά test.

Η **Μελέτη PEA** (Μελέτη PEA, <http://rhea.med.uoc.gr>) η οποία πραγματοποιήθηκε στο Ηράκλειο Κρήτης εξέτασε ένα δείγμα εγκύων γυναικών εντός ενός έτους, που ζουν στο νομό Ηρακλείου, ξεκινώντας από το Φεβρουάριο του 2007. Οι γυναίκες που επιλέχθηκαν στη μελέτη έπρεπε να κατανοούν καλά την ελληνική γλώσσα και να είναι άνω των 17 ετών. Η πρώτη επαφή έγινε πριν από τις πρώτες 15 εβδομάδες της κύησης, κατά τη στιγμή του πρώτου σημαντικού υπερηχογραφικού ελέγχου. Οι συμμετέχουσες κλήθηκαν να δώσουν δείγματα αίματος και ούρων και να συμμετάσχουν σε μια συνέντευξη πρόσωπο με πρόσωπο. Στη συνέχεια επικοινωνήσαμε με τις γυναίκες στον έκτο μήνα της εγκυμοσύνης, στη γέννηση και μετά τη γέννηση, στους 6 και 18 μήνες. Οι προσωπικές συνεντεύξεις,

μαζί με ερωτηματολόγια αυτοαξιολόγησης και ιατρικά αρχεία, χρησιμοποιήθηκαν για την απόκτηση πληροφοριών από διάφορους παράγοντες, συμπεριλαμβανομένων των κοινωνικο-δημογραφικών χαρακτηριστικών, τον τρόπο ζωής και τις διατροφικές συνήθειες, τη μητρική υγεία, κλπ. Λεπτομερή χαρακτηριστικά του πληθυσμού έχουν περιγραφεί. Κατά τη διάρκεια της μελέτης, προσεγγίστηκαν 1765 γυναίκες, 1610 (91%) συμφώνησαν να συμμετάσχουν, και 1388 (86%) παρακολούθηθηκαν έως την γέννα. Η μελέτη εγκρίθηκε από την Επιτροπή Ηθικής του Πανεπιστημιακού Νοσοκομείου Ηρακλείου. Μετά την πλήρη περιγραφή όλων των διαδικασιών, δόθηκε γραπτή συγκατάθεση από όλους τους συμμετέχοντες. [1]

1.2 Ανοσοποίηση

Οι εμβολιασμοί είναι μία από τις πιο σημαντικές ανακαλύψεις για την δημόσια υγεία στην ιστορία του ανθρώπου. Κάθε χώρα εφαρμόζει μια πολιτική εμβολιασμών, που ορίζεται με ένα χρονοδιάγραμμα εμβολιασμών του πληθυσμού της σύμφωνα με τις τρέχουσες επιδημιολογικές συνθήκες και τις διεθνείς οδηγίες. Η πολιτική εμβολιασμών της χώρας ενδέχεται να αλλάξει τροποποιώντας τη χρονολογική σειρά ή την ηλικία έναρξης ή και όλο το πρόγραμμα των εμβολιασμών. Το Πρόγραμμα Εμβολιασμών για την Ελλάδα, προτείνεται από την Εθνική Επιτροπή Εμβολιασμών στο Υπουργείο Υγείας και Κοινωνικής Αλληλεγγύης, και το Υπουργείο το εγκρίνει ή το τροποποιεί. Το Πρόγραμμα Εμβολιασμών για την Ελλάδα φαίνεται στον πίνακα 1.

Ωστόσο, τα προγράμματα εμβολιασμού αντιμετωπίζουν πολλές προκλήσεις παγκοσμίως. Παρότι τα εμβόλια έχουν μειώσει αποτελεσματικά τους κινδύνους των νοσηρών και θνησιγόνων ασθενειών οι οποίες ήταν ευρέως διαδεδομένες στο παρελθόν, οι σημερινές πολιτικές εμβολιασμού γίνονται όλο και πιο αμφιλεγόμενες λόγω αμφιβολιών για την ασφάλεια των εμβολίων. Τα εμβόλια, όπως και άλλα φαρμακευτικά προϊόντα, δεν είναι εντελώς ακίνδυνα. Ενώ οι περισσότερες γνωστές δυσμενείς επιπτώσεις είναι μικρές, ορισμένα εμβόλια έχουν συσχετισθεί με πολύ

σπάνιες αλλά σοβαρές παρενέργειες, όπως για παράδειγμα ασθένειες που σχετίζονται με νευροαναπτυξιακές λειτουργίες (π.χ. αυτισμός) [2, 3].

Εικ. 1. Χρονοδιάγραμμα εμβολιασμών για παιδιά και εφήβους

Ηλικία Εμβόλιο	Γέννηση	1 μην.	2 μην.	4 μην.	6 μην.	12 μην.	15 μην.	18 μην.	24 μην.	4-6 ετ.	11-12 ετ.	13-18 ετ.
Ηπατίτιδας Β (Hep B) ¹	Hep B ^{1a1a}	Hep B ^{1β} (1-2 δόσεις)			Hep B							
			Hep B ^{1γ}	Hep B	Hep B				Hep B (όλες οι δόσεις)			
Διφθερίτιδας, Τετάνου, Κοκκύτη (DTaP) ²			DTaP	DTaP	DTaP		DTaP			DTaP	Tdap ^{2a,2β}	
Πολιομυελίτιδας (IPV) ³			IPV	IPV	IPV					IPV		
Αιμόφιλου τύπου Β ⁴			Hib	Hib	Hib	Hib						
Μηνιγγιτιδόκοκκου C (MCC) ⁵			MCC	MCC	MCC							
Πνευμονιόκοκκου (PCV) ⁶			PCV	PCV	PCV	PCV			PCV (PPV) ^{6a}			
Ιλαράς, Παρωτίτιδας, Ερυθράς (MMR) ⁷						MMR				MMR		
Ανεμευλογιάς (VAR) ⁸						Var				Var		
Ιός ανθρώπινων θηλωμάτων (HPV) ⁹												HPV κορίτσια 12-15 ετ. 3 δόσεις
Ηπατίτιδας Α (Hep A) ¹⁰					Hep A (2 δόσεις)							
Φυματίωσης (BCG) ¹¹						Mantoux				Mantoux ^{11a} BCG	Mantoux ^{11b}	
Γρίπης (INFL) ¹²										INFL (ετησίως)		

--- Τα εμβόλια κάτω από τη διακεκομμένη γραμμή συνιστώνται για επιλεκτικό εμβολιασμό (βλέπε επεξηγήσεις της εκ.1)

Εύρος ηλικιών διενέργειας του εμβολιασμού. Στην παρένθεση αναγράφονται οι δόσεις του εμβολίου που γίνονται σ' αυτό το εύρος ηλικιών, όταν είναι περισσότερες από μία. Το εύρος ηλικιών διενέργειας του εμβολιασμού δίνει τη δυνατότητα να χρησιμοποιούνται μονοδύναμα ή πολυδύναμα (συνδυασμένα) εμβόλια ή/και συνδυασμός μονοδύναμων-συνδυασμένων

Εύρος ηλικιών διενέργειας του εμβολιασμού όταν αυτός δεν έχει προηγηθεί κατά το συνιστώμενο σχήμα ως προς την ηλικία και τις δόσεις (βλέπε επεξηγήσεις πινάκων 1 και 2)

Πίνακας 1: Χρονοδιάγραμμα Εμβολιασμών για παιδιά και εφήβους. [4]

1.3 Νευροανάπτυξη

Τα πρώτα χρόνια της ζωής ενός παιδιού είναι θεμελιώδη για την ανάπτυξη του. Είναι γνωστό πως ο παιδικός εγκέφαλος αναπτύσσεται στο 75% μέχρι την ηλικία των τριών ετών.[5]. Τυχόν ελλείψεις βασικών συστατικών από τη διατροφή των παιδιών μέχρι αυτή την ηλικία επιβραδύνουν τη σωστή ανάπτυξη του εγκεφάλου

και οδηγούν σε διαταραχές του κεντρικού νευρικού συστήματος [5, 6]. Πιο αναλυτικά οι διαταραχές αυτές αναφέρονται σε διάφορα προβλήματα όπως: ανεπαρκής νοητική εξέλιξη που αντανακλάται σε μειωμένο δείκτη νοημοσύνης (IQ), δυσλειτουργία των νεύρων που επηρεάζουν την όραση, ελλιπής διαδικασία εκμάθησης, αδυναμία στη μνήμη του παιδιού, πρόωμη άνοια καθώς και πιο ακραίες ασθένειες[6-8]. Οι διατροφικές συνήθειες στον πρώτο χρόνο ζωής ενός παιδιού, αλλά και στην ενήλικη ζωή του έχουν συσχετιστεί με σοβαρά νοητικά προβλήματα [5]. Συστατικά που η απουσία τους από την παιδική διατροφή οδηγεί σε διαταραχές στην νευροανάπτυξη είναι ο σίδηρος, τα λιπαρά οξέα, το ιώδιο, οι βιταμίνες και ψευδάργυρο [9-11].

Επίσης, ένας σημαντικός παράγοντας που επηρεάζει την νευροανάπτυξη των παιδιών είναι το περιβάλλον. Πολλές είναι οι τοξίνες, τα βαρέα μέταλλα και τα σύνθετα μίγματα ρύπων που χρησιμοποιούνται ευρέως από τη δεκαετία του 1930, [12, 13]. Οι περισσότερες, βιομηχανικές χώρες έχουν πλέον απαγορεύσει ή περιορίσει την παραγωγή των ρύπων [14]. Οι κυριότερες πηγές των για τον περιβαλλοντικών ρύπων είναι τα πλαστικά, οι ορμονικοί διαταράκτες, τα βαρέα μέταλλα, η κατανάλωση ψαριού και αλιευτικών προϊόντων και τα ζωικά λίπη [15]. Τα νεογνά είναι εκτεθειμένα σε τοξίνες μέσω του πλακούντα και του θηλασμού. Πειραματικές μελέτες σε ζώα έχουν δείξει ότι πολλά από αυτά κατατάσσονται στις νευροτοξίνες (όπως PCBs, EtHg) ,παρόλα αυτά οι νευρολογικές επιπτώσεις αυτών των ουσιών στα παιδιά δεν είναι σαφείς [13].

1.4 Νευροανάπτυξη και εμβολιασμός

Πολλές μελέτες ερευνούν τη σχέση μεταξύ των εμβολιασμών και των νευρολογικών διαταραχών. Πρόσφατες μελέτες προτείνουν, τη συσχέτιση ανάπτυξης νευρολογικών διαταραχών (όπως το Σύνδρομο Guillain-Barre, εγκεφαλίτιδα, οξεία διάχυτη εγκεφαλομυελίτιδα) σε παιδιά μετά τον εμβολιασμό τους στα εμβόλια του

κίτρινου πυρετού, του κοκκύτη και του εμβολίου της παρωτίτιδας-ερυθράς-ιλαράς [16-20]

Εκτιμάται ότι το 3% των νευροαναπτυξιακών αναπηριών (NND) είναι άμεσα συνδεδεμένες με νευροτοξικές ουσίες και ότι το 25% αυτών των αναπηριών μπορεί να προκύψουν από ατομικές γενετικές προδιαθέσεις [21]. Ο υδράργυρος Hg στην μορφή του αιθυλ-υδράργυρος EtHg (αιθυλικός υδράργυρος), φαίνεται να είναι η πρώτη νευροτοξική ουσία με την οποία έρχεται σε επαφή ένα βρέφος κατά τη διάρκεια του εμβολιασμού. Ο αιθυλ-υδράργυρος είναι ο μεταβολίτης της θιμεροσάλης η οποία χρησιμοποιείται ευρέως για τη διατήρηση των ανοσοσφαιρινών και των εμβολίων που δίνονται σε εγκύους μητέρες, νεογνά, βρέφη και μικρά παιδιά. [22] Κατά τη διάρκεια του εμβολιασμού μπορεί να προκληθούν παρενέργειες στον εμβολιασμένο πληθυσμό, οι οποίες παρακολουθούνται για να εξακριβωθεί η ασφάλεια του εμβολίου. Νευρολογικά σύνδρομα και παθήσεις μπορεί να εμφανιστούν ως αποτέλεσμα των αντιγόνων του εμβολίου [17, 18, 23]. Σε κάποιες περιπτώσεις, αυτές οι παρενέργειες μπορεί να είναι παροδικές ή μπορεί να εμφανιστούν μετά από χρόνια και μόνο κάτω από συνθήκες που μέχρι σήμερα δεν έχουν χαρακτηριστεί σαφώς.

Σε πολλές χώρες, οι νευρολογικές και άλλες παρενέργειες που προκαλούνται μετά τον εμβολιασμό, παρακολουθούνται στενά από ένα “ενιαίο σύστημα αναφοράς παρενεργειών των εμβολιασμών” (Vaccine Adverse Event Reporting System) όπως το VAERS στις ΗΠΑ. Σοβαρά γεγονότα, όπως νευρολογικές παθήσεις, είναι σπάνια αλλά έχουν παρατηρηθεί σε ενήλικες από την VAERS [20]. Στα παιδιά οι πιο κοινές παρενέργειες των εμβολίων είναι: ήπιος πυρετός, ρίγη και δυσφορία. Παρενέργειες οι οποίες απαιτούν ιατρική φροντίδα είναι: υποτονικά-επεισόδια, πυρετός και σπασμοί, τα οποία θεραπεύονται χωρίς καμία συνέπεια στο 98,4% των περιπτώσεων [25]. Καμία από τις παραπάνω παρενέργειες δεν αποδίδεται στην τοξικότητα του υδραργύρου. Ωστόσο, πολλές είναι οι αναφορές που προτείνουν το εμβόλιο MMR (ιλαράς, παρωτίτιδας και ερυθράς) ως πιθανό παράγοντα κινδύνου για την έκφραση νευροαναπτυξιακών ασθενειών, όπως ο αυτισμός στα παιδιά [26].

Επιδημιολογικές μελέτες χρησιμοποιούν στατιστικά test για την διεξαγωγή ερευνών. Η επιλογή των κατάλληλων στατιστικών test, παραμετρικών η μη, εξαρτάται από τα εκάστοτε δεδομένα της έρευνας. Στον παρακάτω πίνακα αναφέρονται μερικές επιδημιολογικές μελέτες καθώς και τα στατιστικά test που χρησιμοποιήθηκαν.

Πίνακας 2. Στατιστικά test που χρησιμοποιήθηκαν σε επιδημιολογικές μελέτες.

Όνομα Έρευνας	Ερευνητές	Χρήση στατιστικής μεθόδου	Στατιστική Μέθοδος
Frequency of Human Papillomavirus Infection, Coinfection, and Association with Different Risk Factors in Colombia	Camargo, M. Soto-De Leon, S. C. Sanchez, R. Perez-Prados, A. Patarroyo, M. E. Patarroyo, M. A.	Για την αξιολόγηση της σχέσης μεταξύ των παραγόντων κινδύνου και μόλυνσης.	Έλεγχος χ^2 . Λογιστική παλινδρόμηση.
Association between cervical dysplasia and human papillomavirus in HIV seropositive women from Johannesburg South Africa	Firnhaber, C. Van Le, H. Pettifor, A. Schulze, D. Michelow, P. Sanne, I. M. Lewis, D. A. Williamson, A. L. Allan, B. Williams, S. Rinas, A. Levin, S. Smith, J. S.	Για την αξιολόγηση των διαφορών του επιπολασμού του HPV. Για τον καθορισμό των συντελεστών που σχετίζονται με βλάβες του τραχήλου της μήτρας.	Έλεγχος χ^2 . Λογιστική παλινδρόμηση.
Does influenza vaccination improve pediatric asthma outcomes?	Ong, B. A. Forester, J. Fallot, A.	Για τον προσδιορισμό της σχέσης μεταξύ του εμβολιασμού κατά της γρίπης, επιλεγμένων δημογραφικών δεδομένων και της επιδείνωσης του άσθματος.	Λογιστική παλινδρόμηση(SPSS).
Hepatitis B vaccine and the risk of CNS inflammatory demyelination in childhood	Mikaeloff, Y. Caridade, G. Suissa, S. Tardieu, M.	Για τον έλεγχο της σχέσης των odds(ORs)της φλεγμονώδους απομυελίνωσης (CNS) με τον εμβολιασμό HB.	Λογιστική παλινδρόμηση.

[27] [28] [29] [30]

Κεφάλαιο 2

Στόχος Διπλωματικής Εργασίας

Στόχος της παρούσας διπλωματικής εργασίας είναι η διερεύνηση της κατάλληλης στατιστικής μεθοδολογίας για τον έλεγχο της σχέσης ενεργητικής ανοσοποίησης και νευροανάπτυξης σε παιδιά προσχολικής ηλικίας στα πλαίσια της προοπτικής μελέτης μητέρας – παιδιού Κρήτης, Μελέτη Ρέα.

Κεφάλαιο 3

Μεθοδολογία Διεξαγωγής Έρευνας

3.1 Θεμελιώδη Θεωρήματα Και Έννοιες

3.1.1 Θεμελιώδη Θεωρήματα:

Ο ρόλος των θεωρημάτων: σύγκλιση κατά πιθανότητα, σύγκλιση κατά κατανομή, κεντρικό οριακό Θεώρημα, θεώρημα Slutsky και ANMA, είναι πολύ σημαντικός στις στατιστικές μελέτες και γενικά στην θεωρία πιθανοτήτων. Ο λόγος που τα αναφέρουμε σε αυτό το υποκεφάλαιο είναι, κυρίως, διότι χρησιμοποιούνται στις αποδείξεις των στατιστικών test που θα αναλύσουμε στο κεφάλαιο 3.2.

Τα εν λόγω θεωρήματα είναι τα εξής:

1. Σύγκλιση κατά πιθανότητα:

Η ακολουθία των τυχαίων μεταβλητών $X_n, n=1,2,3\dots$ συγκλίνει κατά πιθανότητα στη τυχαία μεταβλητή X όταν $n \rightarrow \infty$, αν για κάθε $\varepsilon > 0$ ισχύει η σχέση:

$$\lim_{n \rightarrow \infty} P(|X_n - X| < \varepsilon) = 1 \text{ ή ισοδύναμα:}$$

$$\lim_{n \rightarrow \infty} P(|X_n - X| \geq \varepsilon) = 0$$

[31]

2. Σύγκλιση κατά κατανομή:

Η ακολουθία των τυχαίων μεταβλητών $X_n, n=1,2,3\dots$ συγκλίνει κατά κατανομή στη τυχαία μεταβλητή X αν η ακολουθία των συναρτήσεων κατανομής $F_n(x) \equiv F_{X_n}(x)$ συγκλίνει στη συνάρτηση κατανομής $F(x) \equiv F_X(x)$ δηλαδή:

$$3. \lim_{n \rightarrow \infty} F_n(x) = F(x)$$

4. για όλα τα σημεία $x \in \mathbb{R}$ στα οποία η $F(x)$ είναι συνεχής. Η $F(x)$ καλείται οριακή συνάρτηση κατανομής.

[31]

5. Κεντρικό Οριακό Θεώρημα(ΚΟΘ):

Έστω X_1, X_2, \dots ανεξάρτητες, ομοκατανεμημένες τυχαίες μεταβλητές με μέσο μ και πεπερασμένη μη μηδενική διασπορά σ^2 . Θέτουμε $S_n = X_1 + X_2 + \dots + X_n$. Τότε:

$$\lim_{n \rightarrow \infty} P\left(\frac{S_n - n\mu}{\sigma\sqrt{n}} \leq x\right) = \Phi(x), \quad -\infty \leq x \leq \infty.$$

6. Θεώρημα Slutsky:

Έστω ακολουθίες A_n, B_n, X_n τέτοιες ώστε:

$A_n \xrightarrow[n \rightarrow \infty]{} a$ και $B_n \xrightarrow[n \rightarrow \infty]{} b$ κατά πιθανότητα και $X_n \xrightarrow[n \rightarrow \infty]{} X$ κατά κατανομή.

Τότε $A_n X_n + B_n \xrightarrow[n \rightarrow \infty]{} aX + b$ κατά κατανομή.

Γενικευμένη μορφή Slutsky:

Έστω $X_n, n=1,2,\dots$ μια ακολουθία τυχαίων μεταβλητών η οποία συγκλίνει στοχαστικά στην τυχαία μεταβλητή X . Αν $g(x), x \in \mathbb{R}$ είναι μια συνεχής συνάρτηση, τότε η ακολουθία των τυχαίων μεταβλητών $g(X_n), n=1,2,\dots$ συγκλίνει στοχαστικά (κατά πιθανότητα) στην τυχαία μεταβλητή $g(x)$.

[31]

7. ΑΝΜΑ(Ασθενής νόμος μεγάλων αριθμών):

Έστω $X_n, n=1,2,\dots$ μια ακολουθία ανεξάρτητων τυχαίων μεταβλητών με μέση τιμή $E(X_n) = \mu_n$ και διασπορά $V(X_n) = \sigma_n^2 < \infty, n=1,2,\dots$.

Αν $\lim_{k \rightarrow \infty} \frac{1}{k^2} \sum_{n=1}^k \sigma_n^2 = 0$, τότε η ακολουθία των διαφορών $\underline{X}_k - \underline{\mu}_k, k=1,2,\dots$ όπου

$\underline{X}_k = \frac{1}{k} \sum_{n=1}^k X_n, \underline{\mu}_k = \frac{1}{k} \sum_{n=1}^k \mu_n$ συγκλίνει στοχαστικά στο μηδέν.

Αν θέσουμε $S_n = X_1 + X_2 + \dots + X_n$, τότε αντίστοιχα:

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{S_n}{n} - \mu\right| \geq \varepsilon\right) = 0.$$

[31]

3.1.2 Σημαντικές Έννοιες:

Σε αυτό το υποκεφάλαιο θα αναφέρουμε κάποιες σημαντικές έννοιες της στατιστικής, όπως είναι το τυχαίο δείγμα, η εκτιμήτρια, η στατιστική συνάρτηση και κυρίως οι έλεγχοι υποθέσεων. Οι έννοιες του τυχαίου δείγματος, της εκτιμήτριας και της στατιστικής συνάρτησης αναφέρονται με λίγα λόγια, σε αντίθεση με αυτήν του ελέγχου υποθέσεων, διότι η πλήρης ανάπτυξη τους θα ξέφευγε από τα πλαίσια αυτής της εργασίας. Παρόλα αυτά, δίνονται κατάλληλες πηγές τις οποίες μπορεί να συμβουλευτεί ο αναγνώστης, για την πλήρη κατανόηση των εννοιών αυτών.

3.1.2.1 Τυχαίο δείγμα:

Τυχαίο δείγμα $\underline{X}_1, \dots, \underline{X}_n$, είναι το δείγμα που παίρνουμε από έναν, υπό μελέτη, πληθυσμό, του οποίου τα μέλη έχουν την ίδια πιθανότητα να εκλεγούν. Δηλαδή, οι τυχαίες μεταβλητές \underline{X}_i είναι ανεξάρτητες και ισόνομες, δηλαδή ακολουθούν την ίδια κατανομή, αυτή του τυχαίου δείγματος.[32]

3.1.2.2 Εκτιμήτρια:

Το τυχαίο δείγμα που έχουμε επιλέξει ακολουθεί κάποια κατανομή (όπως Κανονική, Poisson κλπ.). Έστω λοιπόν $f(x; \underline{\theta})$ η συνάρτηση πιθανότητας της τυχαίας μεταβλητής x , του δείγματος μας. Το $\underline{\theta}$ είναι η παράμετρος της κατανομής, όπως π.χ. το p της Διωνυμικής, ή το λ της Poisson, η τιμή της οποίας μας είναι άγνωστη. Γενικότερα το $\underline{\theta} = (\theta_1, \dots, \theta_m)$ μπορεί να είναι μια άγνωστη m -διάστατη παράμετρος, όπως η διδιάστατη $\underline{\theta} = (\mu, \sigma^2)$ της Κανονικής. Το πρόβλημα της εκτιμητικής είναι πώς με βάση ένα τυχαίο δείγμα $\underline{X}_1, \dots, \underline{X}_n$ είναι δυνατόν να προσδιορίσουμε όσο το δυνατόν καλύτερα την παράμετρο $\underline{\theta}$. Ο προσδιορισμός της καλείται *εκτιμήτρια*. [33]

(Για το πώς υπολογίζεται μια εκτιμήτρια μπορείτε επίσης να δείτε: [34])

3.1.2.3 Στατιστική συνάρτηση:

Μια στατιστική συνάρτηση του τυχαίου δείγματος $\underline{X}_1, \dots, \underline{X}_n$ είναι οποιαδήποτε (μετρήσιμη) συνάρτηση του παραπάνω τυχαίου δείγματος, η οποία δεν εξαρτάται (άμεσα ή έμμεσα) από άγνωστες ποσότητες. Κάθε μια από τις τυχαίες μεταβλητές $\underline{X}_1, \dots, \underline{X}_n$ αποτελεί μια στατιστική συνάρτηση. Αυτές οι στατιστικές συναρτήσεις δίνουν διάφορες πληροφορίες για την παράμετρο $\underline{\theta}$. [34]

3.1.2.4 Έλεγχοι υποθέσεων:

Το κομμάτι των ελέγχων, θα μπορούσαμε να πούμε ότι είναι το πιο σημαντικό αυτού του εισαγωγικού κεφαλαίου, διότι πάνω στην θεωρία των ελέγχων υποθέσεων βασίζονται όλα τα στατιστικά test. Για αυτόν τον λόγο είναι πολύ σημαντικό από τον αναγνώστη να το κατανοήσει πλήρως.

Έστω λοιπόν $\underline{x} = x_1, \dots, x_n$ ανεξάρτητες και ισόνομες(α.ι.) παρατηρήσεις από μία κατανομή με πυκνότητα $f(\underline{x}|\theta)$ ή συνάρτηση μάζας πιθανότητας $p(\underline{x}|\theta)$, $\theta \in \Theta \subseteq \mathbb{R}^k$. Το πρόβλημα του ελέγχου μίας υπόθεσης H , σχετικά με την άγνωστη παράμετρο θ (π.χ. ότι η αληθής τιμή της θ ανήκει στο υποσύνολο Θ_0 του παραμετρικού χώρου Θ την οποία υπόθεση συμβολίζουμε ως $H: \theta \in \Theta_0$), ρωτάει αν είναι αληθής η υπόθεση $H: \theta \in \Theta_0$ έναντι της μόνης εναλλακτικής δυνατότητας $K: \theta \in \Theta_1$, με $\Theta_1 \subseteq \Theta$ και $\Theta_0 \cap \Theta_1 = \emptyset$.

Μπορούμε να διατυπώσουμε το πρόβλημα ως εξής:

$$H: \theta \in \Theta_0 \text{ vs } K: \theta \in \Theta_1$$

Ζητούμε την κατασκευή μιας στατιστικής συνάρτησης απόφασης $d: \mathcal{X} \rightarrow [0,1]$, (όπου \mathcal{X} είναι ο χώρος στον οποίο παίρνει τιμές το δείγμα $\underline{x} = x_1, \dots, x_n$) την οποία καλούμε *ελεγχοσυνάρτηση* ή απλώς *έλεγχος* της υπόθεσης H , έναντι της εναλλακτικής δυνατότητας K . Αν $d(\underline{x})=0$ τότε δεχόμαστε την H ως αληθή, αν $d(\underline{x})=1$ τότε απορρίπτουμε την υπόθεση H , έναντι της εναλλακτικής δυνατότητας K .

Συγκεκριμένα, ορίζουμε στον χώρο \mathcal{X} ένα υποσύνολο του $C \in \mathcal{X}$ (την κρίσιμη περιοχή) και αν $\underline{x} \in C$ τότε απορρίπτουμε την H . Δηλαδή: $d(\underline{x}) = 1(\underline{x} \in C)$, (όπου για το κριτήριο $\underline{x} \in C$ χρησιμοποιούμε την κατάλληλη εκτιμήτρια).

Υπάρχουν δύο τρόποι να κάνουμε σφάλμα όσον αφορά την αποδοχή ή την απόρριψη της H :

- 1) Να απορρίψουμε την H ενώ αυτή είναι αληθής
- 2) Να δεχθούμε την H ενώ αυτή δεν είναι αληθής

Αυτές οι δυνατότητες λάθους είναι γνωστές ως τύπου (1) και (2) και οι αντίστοιχες πιθανότητες συμβολίζονται με:

$$\alpha_1(\theta) := P_\theta(d(\underline{x})=1) = P_\theta(\underline{x} \in \mathcal{X}) = E_\theta [d(\underline{x})], \theta \in \Theta_0$$

$$\alpha_2(\theta) := P_\theta(d(\underline{x})=0) = P_\theta(\underline{x} \notin \mathcal{X}) = E_\theta [1-d(\underline{x})], \theta \in \Theta_1$$

Χειριζόμαστε και τις δυο αυτές συναρτήσεις μέσω της συνάρτησης ισχύος $B_d(\theta)$ του ελέγχου d , η τιμή της οποίας καλείται ισχύς του ελέγχου d στη θέση $\theta \in \Theta$:

$$B_d(\theta) := E_\theta[d(\underline{x})] = \begin{cases} \alpha_1(\theta) & \text{αν } \theta \in \Theta_0 \\ 1 - \alpha_2(\theta) & \text{αν } \theta \in \Theta_1 \end{cases}, \theta \in \Theta$$

Η συνάρτηση ισχύος ενός "καλού" ελέγχου θα πρέπει να παίρνει μικρές τιμές πάνω στο Θ_0 (κοντά στο 0) και μεγάλες τιμές πάνω στο Θ_1 (κοντά στο 1). Ένας τέτοιος "καλός" έλεγχος d θα πρέπει να κρατά χαμηλό το μέτρο του σφάλματος τύπου (1) και επίσης χαμηλό το μέτρο του σφάλματος τύπου (2), ή ισοδύναμα μεγάλη ισχύ πάνω στο Θ_1 . [33]

Τα παραπάνω είναι τα βασικά πράγματα που πρέπει να κατανοήσει ο αναγνώστης για να μπορέσει να προχωρήσει στα επόμενα κεφάλαια

3.2 Παρουσίαση Στατιστικών Μεθόδων

Σε αυτό το κεφάλαιο θα ασχοληθούμε με τους ελέγχους της μορφής $H_0: \theta = \theta_0$ vs $H_1: \theta \neq \theta_0$ ενός τυχαίου δείγματος $\underline{x} = x_1, \dots, x_n$, με παραμετρικές μεθόδους, δηλαδή το τυχαίο δείγμα μας πρέπει να ακολουθεί την Κανονική ή την Πολυωνυμική κατανομή. Παρακάτω αναλύονται οι μέθοδοι-έλεγχοι χ^2 και t καθώς και οι έννοιες της απλής γραμμικής παλινδρόμησης και του συντελεστή συσχέτισης, οι οποίες έχουν άμεση σχέση με τις δύο προηγούμενες μεθόδους-ελέγχους.

3.2.1 Έλεγχος χ^2 (Pearson's chi square test):

Για αυτόν τον έλεγχο, υποθέτουμε ότι τα δεδομένα μας προέρχονται από την Πολυωνυμική κατανομή, άρα ασχολούμαστε με διακριτά δεδομένα. Πολλά προβλήματα μπορούν να μπου στο ακόλουθο πλαίσιο που αφορά την Πολυωνυμική κατανομή:
 Έστω $(N_1, N_2, \dots, N_k) \sim M_{k-1}(p_1, \dots, p_k, n)$ (Πολυωνυμική κατανομή) με $\underline{p} := (p_1, \dots, p_k) \in S := \{\underline{p} \in [0, 1]^k : \sum_{i=1}^k p_i = 1\}$ και $\sum_{i=1}^k N_i = n \in \mathbb{N}$ (φυσικοί) με $N_i \in \mathbb{N}_0 := \{0, 1, 2, \dots\}, i=1, 2, \dots, k$

Δηλαδή:

$$P_{\underline{p}, n}(N_1=n_1, N_2=n_2, \dots, N_k=n_k) = \frac{n!}{n_1! \dots n_k!} p_1^{n_1} \dots p_k^{n_k} 1(\underline{n} \in S')$$

$$\text{όπου } S' := \{\underline{n} \in \mathbb{N}_0^k : \sum_{i=1}^k n_i = n\}$$

Εδώ $\theta = \underline{p} \in \Theta = S$ με $\dim \Theta = k-1$ το n είναι γνωστό. Μας ενδιαφέρει ο έλεγχος της

$$H_0 : \theta \in \Theta_0 \text{ vs } H_1 : \theta \notin \Theta_0,$$

δηλαδή $H_1 : \theta \in \Theta_1$, όπου $\Theta_1 = \Theta \setminus \Theta_0$.

Ο Pearson, ήδη από τις αρχές του προηγούμενου αιώνα (1900), πρότεινε για το σκοπό αυτό τη χρήση μιας στατιστικής συνάρτησης η οποία, υπό την $H_0 : \theta \in \Theta_0$, ακολουθεί (προσεγγιστικά) κατανομή χ^2 , ενώ όταν δεν ισχύει η H_0 (δηλαδή ισχύει η

H₁) το θ λαμβάνει «μεγάλες» τιμές. Η στατιστική συνάρτηση είναι η ακόλουθη:

$$\sum_{i=1}^k \frac{(N_i - n\pi_i)^2}{n\pi_i} = \sum_{i=1}^k \frac{(\text{παρατηρούμενα δεδομένα} - \text{αναμενόμενα δεδομένα})^2}{\text{αναμενόμενα δεδομένα}}$$

Όπου το συγκεκριμένο άθροισμα ακολουθεί την κατανομή X_{k-1}^2 .

Μια προσεγγιστική απόδειξη:

Έστω δείγμα $\underline{N}=(N_1, N_2, \dots, N_k) \sim M_{k-1}(\underline{p}, n)$. Μας ενδιαφέρει ο έλεγχος της $H_0 : (p_1, \dots, p_k) = (\pi_1, \dots, \pi_k) = \text{δεδομένα}$. Η πιθανοφάνεια είναι η ακόλουθη:

$$L(p_1, \dots, p_k | \underline{N}) = \frac{n!}{N_1! \dots N_k!} p_1^{N_1} \dots p_k^{N_k} \quad p \in S \Rightarrow$$

$$\text{Log}(L(p | \underline{N})) = \text{Log}\left(\frac{n!}{N_1! \dots N_k!}\right) + \sum_{i=1}^{k-1} N_i \text{Log} p_i + N_k \text{Log}(1 - \sum_{i=1}^{k-1} p_i)$$

$$\text{Άρα } \frac{\partial \text{Log}(L(p | \underline{N}))}{\partial p_j} = N_j \frac{1}{p_j} + N_k \frac{-1}{1 - \sum_{i=1}^{k-1} p_i} = \frac{N_j}{p_j} - \frac{N_k}{p_k} = 0, \quad 1 \leq j \leq k-1 \Leftrightarrow \hat{p}_j = \frac{\hat{p}_k}{N_k} N_j,$$

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

όπου \hat{p}_j είναι οι εκτιμήτριες \Rightarrow

$$1 - \hat{p}_k = \sum_{j=1}^{k-1} \hat{p}_j = \frac{\hat{p}_k}{N_k} \sum_{j=1}^{k-1} N_j = \frac{\hat{p}_k}{N_k} (n - N_k) \Rightarrow 1 = \hat{p}_k \left(1 + \frac{n - N_k}{N_k}\right) \Rightarrow 1 = \hat{p}_k \frac{n}{N_k} \Rightarrow \hat{p}_k = \frac{N_k}{n}$$

$$\text{Άρα } \hat{p}_j = \frac{\hat{p}_k}{N_k} N_j = \frac{N_j}{n}, \quad j=1, 2, \dots, k-1$$

$$\text{Δηλαδή } \hat{p}_i = \frac{N_i}{n}, \quad i=1, 2, \dots, k$$

Άρα, εφόσον και $\dim \Theta_0 = 0 < \dim \Theta = k-1$, χρησιμοποιούμε το:

$$\lambda(\underline{N}) = \frac{L(\hat{p}_1, \dots, \hat{p}_k | \underline{N})}{L(\pi_1, \dots, \pi_k | \underline{N})} = \left(\frac{N_1}{n\pi_1}\right)^{N_1} \dots \left(\frac{N_k}{n\pi_k}\right)^{N_k} \Rightarrow 2 \text{Log} \lambda(\underline{N}) = \sum_{i=1}^k N_i \text{Log} \left(\frac{N_i}{n\pi_i}\right) =$$

$$2 \sum_{i=1}^k N_i \text{Log} \left(1 + \frac{N_i - n\pi_i}{n\pi_i}\right) \approx \sum_{i=1}^k \frac{(N_i - n\pi_i)^2}{n\pi_i} =: X^2 \text{ υπό την } H_0.$$

Εφόσον $E(N_i) = np_i = n\pi_i$ και $\hat{p}_i = \frac{N_i}{n} \xrightarrow{P} E\left(\frac{N_i}{n}\right) = p_i = \pi_i$ (υπό την H_0) καθώς $n \rightarrow \infty$, από τον ΑΝΜΑ, έχουμε ότι: 0 καθώς $n \rightarrow \infty$.

Άρα για μεγάλα n (από το ανάπτυγμα Taylor $\log(1+x) \approx x - \frac{x^2}{2}$) έχουμε:

$$\begin{aligned} 2\text{Log}\lambda(\underline{N}) &\approx 2 \sum_{i=1}^k N_i \left(\frac{N_i - n\pi_i}{n\pi_i}\right) - \sum_{i=1}^k N_i \left(\frac{N_i - n\pi_i}{n\pi_i}\right)^2 = \\ &= 2 \sum_{i=1}^k \frac{(N_i - n\pi_i)^2}{n\pi_i} + 2 \sum_{i=1}^k (N_i - n\pi_i) - \sum_{i=1}^k \frac{(N_i - n\pi_i)^3}{(n\pi_i)^2} - \sum_{i=1}^k \frac{(N_i - n\pi_i)^2}{n\pi_i} = \\ &= 2x^2 + 2(\sum_{i=1}^k N_i - n \sum_{i=1}^k p_i) - n \sum_{i=1}^k \frac{(\hat{p}_i - \pi_i)^3}{(\pi_i)^2} - x^2 = \end{aligned}$$

$$x^2 + 2(n - n \cdot 1) - \frac{1}{\sqrt{n}} \sum_{i=1}^k \frac{(\sqrt{n}(\hat{p}_i - \pi_i))^3}{(\pi_i)^2} \approx X^2 \text{ υπό την } H_0, \text{ από το ΚΟΘ και το}$$

θεώρημα του Slutsky.

$$\text{Δηλαδή, } 2\text{Log}\lambda(\underline{N}) \approx X^2 := \sum_{i=1}^k \frac{(o_i - e_i)^2}{e_i} \sim X_{k-1}^2 \text{ υπό την } H_0 .$$

Όπου $o_i \equiv N_i$ οι παρατηρήσεις και $e_i \equiv E_H(N_i) = n\pi_i$ οι μέσες τιμές. Το ότι για μεγάλα n , η κατανομή της X^2 είναι, υπό την H_0 , η X_{k-1}^2 . Το οποίο έπεται από την παρακάτω παρατήρηση:

Κάτω από πολύ γενικές συνθήκες: υπό την $H_0 : \theta \in \Theta_0$, έχουμε:

$$2\text{Log}\lambda(\underline{X}_n) \xrightarrow[n \rightarrow \infty]{} X^2, \text{ όπου } k = \dim\Theta - \dim\Theta_0$$

■

(Η πλήρης απόδειξη αποκλίνει πολύ από το περιεχόμενο της εργασίας διότι απαιτεί γνώσεις γραμμικής άλγεβρας και άλλων θεωρημάτων. Για την πλήρη απόδειξη βλέπε: C. R. Rao : *Linear Statistical Inference and its Applications*(σ.391-393))

Για τον έλεγχο της $H_0 : (p_1, \dots, p_k) = (\pi_1, \dots, \pi_k) = \text{δεδομένα}$, ο προσεγγιστικός ΕΠΜΠ(έλεγχος πηλίκου μεγίστων πιθανοφανειών), μεγέθους α , είναι ο ακόλουθος X^2 -έλεγχος(Pearson's X^2):

$$d(\underline{N})=1(X^2 > X_{k-1}^2(1 - \alpha))$$

Άρα, με βάση την παραπάνω στατιστική συνάρτηση X^2 μπορούμε να κατασκευάσουμε έναν έλεγχο για την υπόθεση $H_0 : \theta \in \Theta_0$. Συγκεκριμένα θα απορρίπτουμε την H_0 (σε ε.σ. α περίπου) όταν, με βάση τις παρατηρήσεις π_1, \dots, π_k

$$T(X) > c = X_{k-1}^2(\alpha) : \text{άνω } \alpha\text{-σημείο της } X_{k-1}^2$$

με αντίστοιχο (προσεγγιστικό) p-value:

$$p\text{-value} = P(T(X) > T(x)) = 1 - F_{X_{k-1}^2}(T(x))$$

Παρατηρήσεις:

1. Η προσέγγιση του $2\text{Logl}(\underline{N}) \approx X^2$ είναι πολύ καλή συνήθως και μάλιστα η κατανομή του X^2 για πεπερασμένο n , προσεγγίζει την X_{k-1}^2 καλύτερα από ότι η αντίστοιχη κατανομή του $2\text{Logl}(\underline{N})$. Επίσης ο υπολογισμός του X^2 είναι απλούστερος του υπολογισμού του $2\text{Logl}(\underline{N})$.
2. (έλεγχος X^2 όταν υπάρχουν άγνωστες παράμετροι). Παραπάνω προφανώς θεωρήσαμε ότι τα p_i είναι γνωστά (καθορίζονται πλήρως από την υπόθεση H_0). Υπάρχουν όμως περιπτώσεις όπου τα p_i δεν είναι απολύτως γνωστά, αλλά εξαρτώνται από κάποιες άγνωστες παραμέτρους, δηλαδή $p_i = p_i(\theta)$ με $\theta = (\theta_1, \theta_2, \dots, \theta_r)$ άγνωστο. Η περίπτωση αυτή εμφανίζεται π.χ. κατά τον έλεγχο καλής προσαρμογής δεδομένων σε μία γνωστή κατανομή (π.χ. κανονική) με άγνωστες όμως παραμέτρους (π.χ. μ, σ , δηλ. $p_i = p_i(\mu, \sigma)$) ή π.χ. κατά τον έλεγχο ανεξαρτησίας σε πίνακες συνάφειας (χρησιμοποιώντας το X^2 τεστ). Στην περίπτωση αυτή χρησιμοποιούμε την τροποποιημένη στατιστική συνάρτηση:

$$T'(X) = \sum_{i=1}^k \frac{(N_i - np_i(\hat{\theta}))^2}{np_i(\hat{\theta})}$$

Όπου $\hat{\theta}$ είναι η εκτίμηση του θ από τα δεδομένα. Τώρα, υπό την H_0 ,

αποδεικνύεται ότι η T' ακολουθεί ασυμπτωτικά την X^2 κατανομή με $k - r - 1$ βαθμούς ελευθερίας, όπου r είναι το πλήθος των παραμέτρων που χρειάστηκε να εκτιμηθούν από τα δεδομένα (αρκεί να χρησιμοποιηθούν οι εκτιμήτριες μέγιστης πιθανοφάνειας των παραμέτρων από τα ομαδοποιημένα στις k κλάσεις δεδομένα). Επομένως τώρα, απορρίπτουμε την H_0 σε ε.σ. α (περίπου) όταν :

$$T(X) > X_{k-1-r}^2(\alpha)$$

με αντίστοιχο (προσεγγιστικό) p-value:

$$\text{p-value} \approx P(T'(X) \geq T'(x)) = 1 - F_{X_{k-1-r}^2}(T'(x)).$$

3. Ο έλεγχος X^2 τις περισσότερες φορές δεν είναι ο καλύτερος έλεγχος καλής προσαρμογής για συνεχή δεδομένα διότι προϋποθέτει ομαδοποίηση των δεδομένων (διαμερίζουμε το πεδίο τιμών των παρατηρήσεων σε k σύνολα A_1, A_2, \dots, A_k) με συνέπεια την απώλεια πληροφορίας (επίσης η διαμέριση είναι τις περισσότερες φορές αυθαίρετη). [Σε αυτήν την περίπτωση (δεδομένα από συνεχή κατανομή) συνήθως προτιμάται ο έλεγχος Kolmogorov-Smirnov (K-S)]. Ο έλεγχος X^2 προτιμάται όταν έχουμε διακριτά δεδομένα που παίρνουν τιμές σε ένα πεπερασμένο σύνολο.
4. Το X^2 test, ελέγχει αν υπάρχει σημαντική σχέση μεταξύ των μεταβλητών, αλλά δεν καθορίζει το πόσο σημαντική είναι. Το V (ή αλλιώς Phi) του Cramer είναι ένα test το οποίο μας παρέχει αυτές τις πρόσθετες πληροφορίες. Πρώτα υπολογίζουμε το X^2 και μετά το V του Cramer από τον ακόλουθο τύπο:

$$V = \sqrt{\frac{x^2}{n(k-1)}}$$

Όπου το x^2 είναι το X^2 και k είναι ο αριθμός των γραμμών ή των στηλών του πίνακα.

Το V του Cramer κυμαίνεται μεταξύ 0 και 1. Αν είναι μικρότερο του 0,3 τότε έχουμε μικρή σχέση μεταξύ των μεταβλητών, αν είναι κοντά στο 0,5 έχουμε μέτρια και αν είναι πάνω από 0,7 έχουμε ισχυρή. [35]

Παράδειγμα:

Έστω ότι έχουμε ένα δείγμα 1000 εφήβων από μια πόλη. Για την συγκεκριμένη πόλη λέγεται ότι 1 στους 3 εφήβους παρουσιάζουν νοητική υστέρηση. Εμείς θέλουμε να ελέγξουμε την υπόθεση αυτή. Δηλαδή, οι υποθέσεις μας είναι:

H_0 :Στην συγκεκριμένη πόλη το 1:3 των εφήβων παρουσιάζουν νοητική υστέρηση.
vs

H_1 :Στην συγκεκριμένη πόλη δεν παρουσιάζει το 1:3 των εφήβων νοητική υστέρηση.
Ελέγχοντας το δείγμα σχηματίζουμε τον παρακάτω πίνακα:

	Παρουσιάζουν νοητική υστέρηση.	Δεν παρουσιάζουν νοητική υστέρηση.
Παρατηρήσεις δείγματος. (έστω ότι πήραμε αυτές την μετρήσεις)	259	741
Αναμενόμενα αποτελέσματα.(1:3)	250	750

Επομένως χρησιμοποιώντας τον τύπο έχουμε:

$$\sum_{i=1}^2 \frac{(\text{Παρατηρήσεις δείγματος} - \text{Αναμενόμενα αποτελέσματα})^2}{\text{Αναμενόμενα αποτελέσματα}}$$

$$= \frac{(259 - 250)^2}{250} + \frac{(741 - 750)^2}{750} = \frac{81}{250} + \frac{81}{750} = 0,324 + 0,108 = 0,432$$

$$= \chi^2$$

Με επίπεδο σημαντικότητας $\alpha=0,05$, και βαθμούς ελευθερίας $k=2$ έχουμε:

$$\chi^2_{2-1}(1 - 0,05) = \chi^2_1(0,95) = 3,841. \text{ Βλέπουμε ότι } \chi^2 = 0,432 < \chi^2_1(0,95) = 3,841$$

,δηλαδή $d(\underline{N})=1(0,432 > 3,841)=0$. Άρα δεν απορρίπτουμε την H_0 .

3.2.2 Έλεγχος t (t-test):

Το t-test ασχολείται με ελέγχους που αφορούν την μέση τιμή του δείγματος και βασίζεται στην υπόθεση ότι το δείγμα προέρχεται από την Κανονική κατανομή

και είναι τυχαίο. Ευτυχώς, η εγκυρότητα του t-test δεν επηρεάζεται σημαντικά από μέτριες αποκλίσεις της παραπάνω υπόθεσης (αυτό λέγεται robustness και ο G.E.P. Box ήταν ο πρώτος που διατύπωσε τον όρο το 1953).[32]

Ο συγκεκριμένος έλεγχος παίρνει τις παρακάτω μορφές ανάλογα με την κάθε περίπτωση(σε όλες τις παραπάνω περιπτώσεις υποθέτουμε ότι οι διασπορές των δειγμάτων είναι ίσες) :

Α)Αμφίπλευρος μονο-δειγματικός έλεγχος t :

Έστω X_1, \dots, X_n ανεξάρτητες και ισόνομες τυχαίες μεταβλητές που ακολουθούν την $N(\mu, \sigma^2), \theta = (\mu, \sigma^2) \in \Theta = \mathbb{R} \times \mathbb{R}_+$. Τότε ο ΕΠΜΠ, με επίπεδο σημαντικότητας α , της υπόθεσης $H_0: \mu = \mu_0$ vs $H_1: \mu \neq \mu_0$ είναι ο ακόλουθος :

$$d(\underline{X}) = 1(|T_n| > t_{n-1}(1 - \frac{\alpha}{2}))$$

Όπου $T_n := \frac{\sqrt{n}(\bar{X}_n - \mu_0)}{S_n} \sim t_{n-1}$, με $S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$.

Απόδειξη:

Το $LR(\underline{X}) := \frac{L(\hat{\theta}_1 | \underline{X})}{L(\hat{\theta}_0 | \underline{X})}$ ονομάζεται πηλίκο μεγίστων πιθανοφανειών (για περισσότερες πληροφορίες βλέπε: Στατιστική Συμπερασματολογία, Τόμος 1-Εκτιμητική(κεφ.3)-Γεωργίου Γρ. Ρουσσά) και επειδή $\Theta_0 = \{\mu_0\} \times \mathbb{R}_+$, $\Theta_1 = \Theta \setminus \Theta_0$ και $\dim \Theta_0 = 1 < \dim \Theta = 2$ και η Κανονική κατανομή είναι συνεχής ισχύει ότι

$$LR(\underline{X}) := \frac{L(\hat{\theta}_1 | \underline{X})}{L(\hat{\theta}_0 | \underline{X})} = \frac{L(\hat{\theta} | \underline{X})}{L(\hat{\theta}_0 | \underline{X})}$$

Στην προκειμένη περίπτωση έχουμε: $LR(\underline{X}) := \frac{L(\hat{\mu}, \hat{\sigma}^2 | \underline{X})}{L(\hat{\mu}_0, \hat{\sigma}_0^2 | \underline{X})} = \frac{L(\bar{X}_n, \hat{\sigma}_n^2 | \underline{X})}{L(\hat{\mu}_0, \hat{\sigma}_n^2 | \underline{X})} = \lambda$ όπου

$$\hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu_0)^2, \hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2.$$

$$\text{Άρα } \lambda = \frac{(2\pi\hat{\sigma}_n^2)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2\hat{\sigma}_n^2} \sum_{i=1}^n (X_i - \bar{X}_n)^2\right\}}{(2\pi\hat{\sigma}_n^2)^{-\frac{n}{2}} \exp\left\{-\frac{1}{2\hat{\sigma}_n^2} \sum_{i=1}^n (X_i - \mu_0)^2\right\}} = \left(\frac{\hat{\sigma}_n^2}{\hat{\sigma}_n^2}\right)^{\frac{n}{2}} = \left\{ \frac{\sum_{i=1}^n (X_i - \bar{X}_n)^2 + n(\bar{X}_n - \mu_0)^2}{\sum_{i=1}^n (X_i - \bar{X}_n)^2} \right\}^{\frac{n}{2}} =$$

$$= \left(1 + \frac{1}{n-1} T_n^2\right)^{\frac{n}{2}} > c \Leftrightarrow T_n^2 > c \Leftrightarrow |T_n| > c$$

Άρα $d(\underline{X}) = 1(|T_n| > c_\alpha)$ όπου

$$\alpha = P_{\mu_0}(|T_n| > c_\alpha) = 1 - P_{\mu_0}(-c_\alpha \leq T_n \leq c_\alpha) = 1 - F_{T_n}(c_\alpha) + F_{T_n}(-c_\alpha) = 2[1 - F_{T_n}(c_\alpha)]$$

$$\Rightarrow c_\alpha = F_{T_n}^{-1}\left(1 - \frac{\alpha}{2}\right) =: t_{n-1}\left(1 - \frac{\alpha}{2}\right)$$

■

Τα $-t_{n-1}\left(1 - \frac{\alpha}{2}\right)$ και $t_{n-1}\left(1 - \frac{\alpha}{2}\right)$ είναι οι κρίσιμες τιμές του T_n . Αν το T_n είναι μεγαλύτερο από το $t_{n-1}\left(1 - \frac{\alpha}{2}\right)$ ή μικρότερο ή ίσο με το $-t_{n-1}\left(1 - \frac{\alpha}{2}\right)$, τότε απορρίπτουμε την H_0 .

Β) Μονόπλευρος μονο-δειγματικός έλεγχος t:

Έστω X_1, \dots, X_n ανεξάρτητες και ισόνομες τυχαίες μεταβλητές που ακολουθούν την $N(\mu, \sigma^2)$, $\theta = (\mu, \sigma^2) \in \Theta = \mathbb{R} \times \mathbb{R}_+$. Τότε ο ΕΠΜΠ, με επίπεδο σημαντικότητας α , της υπόθεσης $H_0: \mu \leq \mu_0$ vs $H_1: \mu > \mu_0$ είναι ο ακόλουθος:

$$d(\underline{X}) = 1(T_n > t_{n-1}(1 - \alpha))$$

$$\text{Όπου } T_n := \frac{\sqrt{n}(\bar{X}_n - \mu_0)}{S_n} \sim t_{n-1}, \text{ με } S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2.$$

Απόδειξη:

Σε αυτή την περίπτωση $\Theta_0 = [-\infty, \mu_0] \times \mathbb{R}_+$, $\Theta_1 = \Theta \setminus \Theta_0$ και $\dim \Theta_0 = \dim \Theta = 2$.

Τώρα

$LR(\underline{X}) := \frac{L(\hat{\theta}_1 | \underline{X})}{L(\hat{\theta}_0 | \underline{X})} = \frac{L(\hat{\mu}_1, \hat{\sigma}_1^2 | \underline{X})}{L(\hat{\mu}_0, \hat{\sigma}_0^2 | \underline{X})}$ την ύπαρξη της πιθανοφάνειας μας εξασφαλίζουν η κοιλότητα της πιθανοφάνειας και η κυρτότητα των Θ_0, Θ_1 .

Τότε μέσο της $\frac{\partial}{\partial \sigma^2} \text{Log} L(\mu, \sigma^2 | \underline{X}) = 0$ παίρνουμε $\sigma^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$, άρα

$$\hat{\sigma}_0^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu}_0)^2 \text{ και } \hat{\sigma}_1^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu}_1)^2$$

$$\text{Οπότε } L(\hat{\mu}_1, \hat{\sigma}_1^2 | \underline{X}) = (2\pi\hat{\sigma}_1^2)^{-\frac{n}{2}} \exp\left\{-\frac{n}{2}\right\} = (2\pi e)^{-\frac{n}{2}} \left\{\frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu}_1)^2\right\}^{-\frac{n}{2}} =$$

$$= c \left\{ \sum_{i=1}^n (X_i - \bar{X}_n)^2 + n(\bar{X}_n - \hat{\mu}_1)^2 \right\}^{\frac{-n}{2}}$$

$$\begin{aligned} \text{και } L(\hat{\mu}_0, \hat{\sigma}_0^2 | X) &= (2\pi\hat{\sigma}_0^2)^{\frac{-n}{2}} \exp\left\{-\frac{n}{2}\right\} = (2\pi e)^{\frac{-n}{2}} \left\{ \frac{1}{n} \sum_{i=1}^n (X_i - \hat{\mu}_0)^2 \right\}^{\frac{-n}{2}} = \\ &= c \left\{ \sum_{i=1}^n (X_i - \bar{X}_n)^2 + n(\bar{X}_n - \hat{\mu}_0)^2 \right\}^{\frac{-n}{2}} \end{aligned}$$

$$\text{όπου } c = \left(\frac{2\pi e}{n} \right)^{\frac{-n}{2}}$$

$$\begin{aligned} \text{Άρα } L(\hat{\mu}_0, \hat{\sigma}_0^2 | X) &= c \left\{ (n-1)S_n^2 + n(\bar{X}_n - \hat{\mu}_0)^2 \right\}^{\frac{-n}{2}} \leq \\ &\leq c \left\{ (n-1)S_n^2 + n(\bar{X}_n - \mu_0 \wedge \bar{X}_n)^2 \right\}^{\frac{-n}{2}} \end{aligned}$$

$$\begin{aligned} \text{και } L(\hat{\mu}_1, \hat{\sigma}_1^2 | X) &= c \left\{ (n-1)S_n^2 + n(\bar{X}_n - \hat{\mu}_1)^2 \right\}^{\frac{-n}{2}} \leq \\ &\leq c \left\{ (n-1)S_n^2 + n(\bar{X}_n - \mu_0 \vee \bar{X}_n)^2 \right\}^{\frac{-n}{2}} \end{aligned}$$

Έχουμε, λοιπόν:

$$\begin{aligned} LR(X) &= \left(\frac{\hat{\sigma}_0^2}{\hat{\sigma}_1^2} \right)^{\frac{n}{2}} = \frac{\left\{ \sum_{i=1}^n (X_i - \mu_0 \wedge \bar{X}_n)^2 \right\}^{\frac{n}{2}}}{\left\{ \sum_{i=1}^n (X_i - \mu_0 \vee \bar{X}_n)^2 \right\}^{\frac{n}{2}}} = \frac{\left\{ (n-1)S_n^2 + n(\bar{X}_n - \mu_0 \wedge \bar{X}_n)^2 \right\}^{\frac{n}{2}}}{\left\{ (n-1)S_n^2 + n(\bar{X}_n - \mu_0 \vee \bar{X}_n)^2 \right\}^{\frac{n}{2}}} = \\ &= \begin{cases} \left\{ \frac{(n-1)S_n^2}{(n-1)S_n^2 + n(\bar{X}_n - \mu_0)^2} \right\}^{\frac{n}{2}} & \text{αν } \bar{X}_n \leq \mu_0 \\ \left\{ \frac{(n-1)S_n^2 + n(\bar{X}_n - \mu_0)^2}{(n-1)S_n^2} \right\}^{\frac{n}{2}} & \text{αν } \bar{X}_n > \mu_0 \end{cases} = \\ &= \begin{cases} \left(1 + \frac{1}{n-1} T_n^2 \right)^{\frac{-n}{2}} > c & \text{αν } T_n \leq 0 \\ \left(1 + \frac{1}{n-1} T_n^2 \right)^{\frac{n}{2}} > c & \text{αν } T_n > 0 \end{cases} \Leftrightarrow \\ &\Leftrightarrow \begin{cases} |T_n| < c & \text{αν } T_n \leq 0 \\ |T_n| > c & \text{αν } T_n > 0 \end{cases} \Leftrightarrow T_n > c \end{aligned}$$

Άρα, $d(X) = 1(T_n > c_\alpha)$ όπου $\alpha = \sup\{P_\mu(T_n > c_\alpha) : \mu \leq \mu_0\} = P_{\mu_0}(T_n > c_\alpha)$, διότι

$$\begin{aligned} P_\mu(T_n > c_\alpha) &= P_\mu \left(\frac{\sqrt{n}(\bar{X}_n - \mu)}{S_n} + \frac{\sqrt{n}(\mu - \mu_0)}{S_n} > c_\alpha \right) = \\ &= P_\mu \left(t_{n-1} + \frac{\sqrt{n}(\mu - \mu_0)}{S_n} > c_\alpha \right) \leq P_\mu(t_{n-1} > c_\alpha) \end{aligned}$$

εφόσον $\mu \leq \mu_0$.

$$\text{Άρα } 1 - \alpha = P_\mu(t_{n-1} > c_\alpha) \Rightarrow c_\alpha = t_{n-1}$$

■

Γ) Αμφίπλευρος δι-δειγματικός έλεγχος t , ανεξάρτητων μεταβλητών:

Σε αυτή την περίπτωση έχουμε δύο δείγματα έστω X_1, \dots, X_n και Y_1, \dots, Y_m ανεξάρτητα μεταξύ τους τα οποία ακολουθούν τις $N(\mu_1, \sigma^2)$ και $N(\mu_2, \sigma^2)$ αντίστοιχα, δηλαδή $\underline{\theta} = (\mu_1, \mu_2, \sigma^2) \in \Theta = \mathbb{R} \times \mathbb{R} \times \mathbb{R}_+$ και ο αντίστοιχος ΕΠΜΠ, μεγέθους α , της υπόθεσης $H_0: \mu_1 = \mu_2$ vs $H_1: \mu_1 \neq \mu_2$ είναι ο ακόλουθος :

$$d(\underline{X}, \underline{Y}) = 1(|T(\underline{X}, \underline{Y})| > t_{m+n-2}(1 - \frac{\alpha}{2}))$$

Όπου $T(\underline{X}, \underline{Y}) := \frac{\sqrt{\frac{m+n}{mn}}(\bar{X}_n - \bar{Y}_m)}{S_p} \sim t_{m+n-2}$ με $S_p^2 = \frac{(n-1)S_x^2 + (m-1)S_y^2}{m+n-2}$ και

$$S_x^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \text{ και } S_y^2 = \frac{1}{m-1} \sum_{i=1}^m (Y_i - \bar{Y}_m)^2.$$

Απόδειξη:

Εφόσον $H_0: \mu_1 = \mu_2 = \mu$, η κοινή μέση, και $\Theta_0 = \mathbb{R} \times \mathbb{R}_+$, δηλαδή $\dim \Theta_0 = 2 < \dim \Theta = 3$

ισχύει όπως και πριν $LR(\underline{X}, \underline{Y}) := \frac{L(\hat{\theta}|\underline{X}, \underline{Y})}{L(\hat{\theta}_0|\underline{X}, \underline{Y})} = \lambda(\underline{X}, \underline{Y})$. Υπό την H_0 τα $\underline{X}, \underline{Y}$ έχουν την ίδια κατανομή $N(\mu, \sigma^2)$ άρα, $\hat{\mu}_0 = \hat{\mu}_1 = \hat{\mu}_2 = \frac{n\bar{X}_n + m\bar{Y}_m}{m+n}$ και $\hat{\sigma}_0^2 = \frac{\sum_{i=1}^n (X_i - \hat{\mu}_0)^2 + \sum_{i=1}^m (Y_i - \hat{\mu}_0)^2}{m+n}$ (όπου με το σύμβολο $\hat{\cdot}$ συμβολίζουμε την εκτιμήτρια της αντίστοιχης μεταβλητής).

$$\text{Άρα } L(\hat{\theta}_0|\underline{X}, \underline{Y}) = (2\pi\hat{\sigma}_0^2)^{-\frac{(m+n)}{2}} \exp\left\{-\frac{(m+n)}{2}\right\}.$$

Τώρα, για την γενική περίπτωση, $\hat{\mu}_1 = \bar{X}_n, \hat{\mu}_2 = \bar{Y}_m$ και $\hat{\sigma}^2 = \frac{\sum_{i=1}^n (X_i - \bar{X}_n)^2 + \sum_{i=1}^m (Y_i - \bar{Y}_m)^2}{m+n}$

$$\text{και άρα, } L(\hat{\theta}|\underline{X}, \underline{Y}) = (2\pi\hat{\sigma}^2)^{-\frac{(m+n)}{2}} \exp\left\{-\frac{(m+n)}{2}\right\}.$$

$$\text{Έχουμε λοιπόν } \lambda(\underline{X}, \underline{Y}) = \frac{\hat{\sigma}_0^2}{\hat{\sigma}^2}^{\frac{n+m}{2}} = \frac{\left\{\sum_{i=1}^n (X_i - \hat{\mu}_0)^2 + n(\bar{X}_n - \hat{\mu}_0)^2\right\}^{\frac{n+m}{2}}}{\left\{\sum_{i=1}^n (X_i - \bar{X}_n)^2 + n(\bar{X}_n - \bar{Y}_m)^2\right\}^{\frac{n+m}{2}}} =$$

$$= \left\{ \frac{(m+n-2)S_p^2 + n(\bar{X}_n - \hat{\mu}_0)^2 + m(\bar{Y}_m - \hat{\mu}_0)^2}{(m+n-2)S_p^2} \right\}^{\frac{m+n}{2}}$$

$$\text{Τώρα, } n(\bar{X}_n - \hat{\mu}_0)^2 = n\left(\bar{X}_n - \frac{n\bar{X}_n + m\bar{Y}_m}{m+n}\right)^2 = \frac{m^2 n}{(m+n)^2} (\bar{X}_n - \bar{Y}_m)^2 = m(\bar{Y}_m - \hat{\mu}_0)^2$$

$$\text{Άρα } \lambda(\underline{X}, \underline{Y}) = \left\{ 1 + \frac{1}{m+n-2} \sqrt{\frac{mn}{m+n}} \frac{(\bar{X}_n - \bar{Y}_m)^2}{S_p^2} \right\}^{\frac{(m+n)}{2}} = \left\{ 1 + \frac{1}{m+n-2} [T(\underline{X}, \underline{Y})]^2 \right\}^{\frac{(m+n)}{2}} > c$$

$$\Leftrightarrow |T(\underline{X}, \underline{Y})| > c$$

$$\text{Τώρα, } \frac{\bar{X}_n - \bar{Y}_m}{\sigma} = \frac{\bar{X}_n - \mu_1}{\sigma} - \frac{\bar{X}_n - \mu_2}{\sigma} \sim N\left(0, \frac{1}{n}\right) + N\left(0, \frac{1}{m}\right) = N\left(0, \frac{1}{n} + \frac{1}{m}\right) \Rightarrow \sqrt{\frac{m+n}{mn}} \frac{\bar{X}_n - \bar{Y}_m}{\sigma} \sim N(0,1)$$

(κατά κατανομή, υπό την H_0) και

$$\frac{(m+n-2)S_p^2}{\sigma^2} = \frac{\sum_{i=1}^n (X_i - \bar{X}_n)^2}{\sigma^2} + \frac{\sum_{i=1}^m (Y_i - \bar{Y}_m)^2}{\sigma^2} \sim X_{n-1}^2 + X_{m-1}^2 = X_{n+m-2}^2 \text{ (κατά κατανομή).}$$

Άρα, $T(\underline{X}, \underline{Y}) \sim t_{m+n-2}$ οπότε παίρνουμε την ακριβή μορφή του ΕΠΜΠ, μεγέθους α :

$$P(|T(\underline{X}, \underline{Y})| > t_{m+n-2}(1 - \frac{\alpha}{2})) = \alpha$$

■

Δ) Αμφίπλευρος δι-δειγματικός έλεγχος t , εξαρτημένος μεταβλητών:

Σε αντίθεση με τις δύο προηγούμενες περιπτώσεις, σε αυτή την περίπτωση έχουμε δύο δείγματα, X_1, \dots, X_n και Y_1, \dots, Y_m εξαρτημένα μεταξύ τους τα οποία ακολουθούν τις $N(\mu_1, \sigma^2)$ και $N(\mu_2, \sigma^2)$ αντίστοιχα, δηλαδή $\underline{\theta} = (\mu_1, \mu_2, \sigma^2) \in \Theta = \mathbb{R} \times \mathbb{R} \times \mathbb{R}_+$ και ο αντίστοιχος ΕΠΜΠ, μεγέθους α , της υπόθεσης $H_0: \delta = \mu_1 - \mu_2 = \delta_0$ vs $H_1: \delta \neq \delta_0$ είναι ισοδύναμος με αυτόν της $H_0: \mu_1^* = \mu_2$ vs $H_1: \mu_1^* \neq \mu_2$ (όπου $\mu_1^* = \mu_1 - \delta_0$, και $x_i^* := x_i - \delta_0, i = 1, 2, \dots, n$ είναι ανεξάρτητες και ισόνομες τυχαίες μεταβλητές της $N(\mu_1^*, \sigma^2)$), ο οποίος είναι ο ακόλουθος:

$$d(\underline{X}, \underline{Y}) = 1(|T(\underline{X}, \underline{Y})| > t_{m+n-2}(1 - \frac{\alpha}{2}))$$

$$\text{Όπου } T(\underline{X}, \underline{Y}) := \frac{\sqrt{\frac{m+n}{mn}}(\bar{X}_n^* - \bar{Y}_m)}{S_p} = \frac{\sqrt{\frac{m+n}{mn}}(\bar{X}_n - \bar{Y}_m - \delta_0)}{S_p} \sim t_{m+n-2} \text{ με } S_p^2 = \frac{(n-1)S_x^2 + (m-1)S_y^2}{m+n-2}$$

και

$$S_x^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \text{ και } S_y^2 = \frac{1}{m-1} \sum_{i=1}^m (Y_i - \bar{Y}_m)^2.$$

Απόδειξη:

Δεν θα αναφέρουμε ολόκληρη την απόδειξη διότι είναι παρόμοια με αυτήν του (B). Θα αναφέρουμε απλώς ότι το S_p^2 (συνεπώς και το S_p) δεν επηρεάζονται από το δ_0 διότι δεν επηρεάζεται το S_x^2 , αφού $\bar{X}_n^* = \frac{1}{n} \sum_{i=1}^n (X_i - \delta_0) = \bar{X}_n - \delta_0$.

■

Τέλος, για τις περιπτώσεις των (Γ) (Δ), τα $-t_{m+n-2}(1 - \frac{\alpha}{2})$ και $t_{m+n-2}(1 - \frac{\alpha}{2})$ είναι οι κρίσιμες τιμές του $T(\underline{X}, \underline{Y})$. Αν το $T(\underline{X}, \underline{Y})$ είναι μεγαλύτερο από το $t_{m+n-2}(1 - \frac{\alpha}{2})$ ή μικρότερο ή ίσο με το $-t_{m+n-2}(1 - \frac{\alpha}{2})$, τότε απορρίπτουμε την H_0 και αποδεχόμαστε την H_1 . [35]

3.2.3 P-value:

Έστω ότι για τον έλεγχο της $H_0: \theta \in \Theta_0$ vs $H_1: \theta \notin \Theta_0$ καταλήξαμε σε μια ελεγχοσυνάρτηση της μορφής: $d(\underline{X}) = 1(T(\underline{X}) > c)$, για κάποια στατιστική συνάρτηση T και κάποια-όχι συγκεκριμένη-σταθερά c . Εκτός από το να ακολουθήσουμε την άποψη των Neyman-Pearson, περί σταθερού-προκαθορισμένου-επιπέδου α για τον d , έχουμε και την ακόλουθη διέξοδο:

Για το συγκεκριμένο δείγμα \underline{x} που πήραμε, μπορούμε να υπολογίσουμε την πιθανότητα: $\alpha(T(\underline{x})) := \sup\{P_\theta(T(\underline{X}) > T(\underline{x})) : \theta \in \Theta_0\} \forall \underline{x} \in \mathcal{X}$, η οποία καλείται p-value του ελέγχου d στο δείγμα \underline{x} .

Παρατηρούμε ότι, η p-value του d στο \underline{x} , είναι ένα μέτρο του πόσο σύνηθες είναι να παρατηρηθεί τυχόν δείγμα \underline{X} πιο ακραίο απ'το \underline{x} που παρατηρήσαμε, εφόσον η H_0 είναι αληθής. Αν η p-value του d στο \underline{x} είναι μεγάλη, αυτό δείχνει ότι, υπό την H_0 , η παρατήρηση ενός δείγματος σαν το \underline{x} ή και πιο ακραίου, δεν είναι κάτι το ασύνηθες. Άρα το ότι παρατηρήσαμε το \underline{x} , αποτελεί στοιχείο που ενισχύει την αληθοφάνεια της H_0 .

Από την άποψη της μεθόδου των Neyman-Pearson (για περαιτέρω πληροφορίες σχετικά με την θεωρία των Neyman και Pearson δείτε [36]), $\forall \alpha < \alpha(T(\underline{x}))$, έστω c_α η αντίστοιχη κρίσιμη σταθερά του ελέγχου d , μεγέθους α ,

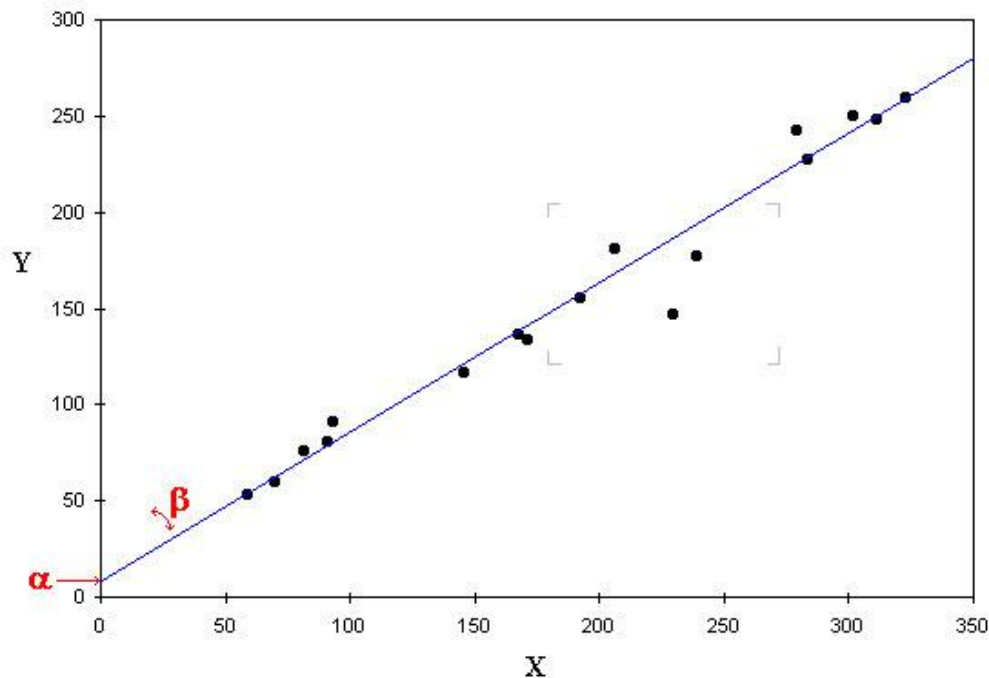
δηλαδή, $d(\underline{X})=d_{\alpha}(\underline{X}):=1(T(\underline{X})>c_{\alpha})$. Αν $T(\underline{x})>c_{\alpha}$ τότε: $P_{\theta}(T(\underline{X})>c_{\alpha}) \geq P_{\theta}(T(\underline{X})>T(\underline{x})) \forall \theta \in \Theta \Rightarrow \alpha \geq \sup\{P_{\theta}(T(\underline{X})>c_{\alpha}) : \theta \in \Theta_0\} \geq \sup\{P_{\theta}(T(\underline{X})>T(\underline{x})) : \theta \in \Theta_0\} = \alpha(T(\underline{x}))$, το οποίο είναι άτοπο. Άρα, $T(\underline{x}) \leq c_{\alpha}$, δηλαδή, δεχόμαστε την H_0 ως αληθή σε επίπεδο α , $\forall \alpha < \alpha(T(\underline{x}))$. (Γενικά η H_0 γίνεται δεκτή $\forall \alpha \leq \alpha(T(\underline{x}))$, βάση του δείγματος \underline{x} που παρατηρήσαμε) [35]

3.2.4 Γραμμική Παλινδρόμηση

Ένα σημαντικό ερώτημα σε πάρα πολλά προβλήματα είναι αν μπορούμε να εκτιμήσουμε ή να προβλέψουμε την τιμή μιας ή περισσότερων «μεταβλητών» κάτω από ορισμένες συνθήκες. Οι δοσμένες συνθήκες περιγράφονται και αυτές από μεταβλητές, οι τιμές των οποίων είναι δυνατό να ελεγχθούν από τον ερευνητή. Έτσι για παράδειγμα η μεταβλητή Y που ζητούμε να εκτιμηθεί ή να προβλεφθεί, μπορεί να παριστάνει «ζήτηση κάποιου προϊόντος στην αγορά», «απόδοση μαθητού» κλπ. Ενώ οι μεταβλητές X_i που περιγράφουν τις συνθήκες και που μπορούν να ελεγχθούν, μπορεί να παριστάνουν «θερμοκρασία», «φύλο», και άλλα. Σε αυτές τις περιπτώσεις το μέγεθος των μεταβλητών Y (εξαρτημένες) υποθέτουμε πως καθορίζεται από το μέγεθος των μεταβλητών X_i (ανεξάρτητες), ενώ το αντίθετο δεν ισχύει. Οι όροι εξαρτημένες-ανεξάρτητες μεταβλητές δεν αναφέρονται απαραίτητα στην καθαυτή σχέση μεταξύ των μεταβλητών Y και X_i . Αυτή η διαδικασία δημιουργίας μοντέλου πρόγνωσης της εξαρτημένης μεταβλητής από τις ανεξάρτητες μεταβλητές λέγεται παλινδρόμηση (regression). Ο όρος απλή παλινδρόμηση (simple regression) αναφέρεται στην απλούστερη μορφή παλινδρόμησης, στην οποία εμπλέκονται μόνο δύο μεταβλητές.

Τα δεδομένα που υπόκειται σε απλή παλινδρόμηση αποτελούνται από ένα ζευγάρι δεδομένων. Είναι πολύ βολικό και χρήσιμο να απεικονίζουμε αυτά τα δεδομένα σε έναν πίνακα, όπου στον Y άξονα θα έχουμε τις εξαρτημένες μεταβλητές και στον X τις ανεξάρτητες.

Πίνακας 3.Απεικόνιση δεδομένων.



Όπου οι κουκίδες είναι οι παρατηρούμενες τιμές και η διαγώνιος είναι οι εκτιμώμενες.

Σε κάποιες περιπτώσεις, η σχέση μεταξύ των δύο μεταβλητών δεν είναι σχέση εξάρτησης. Σε αυτές τις περιπτώσεις, όταν αλλάζει το μέγεθος της μιας μεταβλητής επηρεάζεται το μέγεθος της άλλης, αλλά δεν μπορούμε να πούμε με σιγουριά ότι υπάρχει μια εξαρτημένη και μια ανεξάρτητη μεταβλητή. Τότε αντί της απλής παλινδρόμησης χρησιμοποιούμε την συσχέτιση (correlation), που θα αναλύσουμε αργότερα.

3.2.5 Απλή Γραμμική Παλινδρόμηση:

Η απλούστερη σχέση δύο μεταβλητών σε ένα πληθυσμό είναι αυτή της απλής γραμμικής παλινδρόμησης, η οποία γράφεται ως εξής:

$$Y_i = \alpha + \beta X_i$$

Όπου τα α και β σταθερές (παραμέτροι του πληθυσμού). Αυτή είναι η εξίσωση μιας ευθείας γραμμής, αλλά επειδή σε έναν πληθυσμό, συνήθως, τα δεδομένα δεν ακολουθούν μια ευθεία γραμμή, εισάγουμε το λεγόμενο σφάλμα στην εξίσωση μας, η οποία παίρνει την μορφή:

$$Y_i = \alpha + \beta X_i + \varepsilon_i$$

όπου το άθροισμα των ε_i είναι μηδέν.

Σκοπός μας είναι να εκτιμήσουμε τις τιμές της Y μεταβλητής με τον καλύτερο δυνατό τρόπο ώστε τα δεδομένα μας να μπορούν να απεικονιστούν σε μια ευθεία γραμμή. Ο τρόπος αυτός ονομάζεται *μέθοδος ελαχίστων τετραγώνων*. Κάθε τιμή της μεταβλητής X θα αντιστοιχεί, εκτός από την τιμή Y , σε μια άλλη τιμή \hat{Y} , η οποία θα βρίσκεται πάνω στην ευθεία γραμμή. Δηλαδή η \hat{Y} είναι η εκτιμήτρια της Y , την οποία και θέλουμε να υπολογίσουμε ούτως ώστε το $\sum_{i=1}^n (Y_i - \hat{Y}_i)^2$ (τετράγωνο των σφαλμάτων όπου n είναι το πλήθος του δείγματος), να είναι το ελάχιστο.

Για τον υπολογισμό των \hat{Y}_i θα πρέπει πρώτα να υπολογίσουμε τις παραμέτρους πληθυσμού α και β , και για να γίνει αυτό θα πρέπει να κατέχουμε όλα τα δεδομένα του πληθυσμού που μελετάμε. Επειδή αυτό όμως είναι σχεδόν πάντα αδύνατο, πρέπει να εκτιμήσουμε αυτές τις παραμέτρους οι οποίες προέρχονται από ένα δείγμα μεγέθους n , όπως έχουμε ήδη αναφέρει. Θα χρειαστούμε τις δυο παρακάτω εξισώσεις:

- $\sum_{i=1}^n (x_i)^2 = \sum_{i=1}^n (X_i - \bar{X})^2 = \sum_{i=1}^n (X_i)^2 - \frac{(\sum_{i=1}^n X_i)^2}{n}$
- $\sum_{i=1}^n x_i y_i = \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y}) = \sum_{i=1}^n X_i Y_i - \frac{(\sum_{i=1}^n X_i)(\sum_{i=1}^n Y_i)}{n}$

Οπότε η εκτιμήτρια (έστω b) για την παράμετρο β υπολογίζεται με την παρακάτω εξίσωση:

$$b = \frac{\sum_{i=1}^n x_i y_i}{\sum_{i=1}^n (x_i)^2}$$

όπου $-\infty \leq b \leq \infty$.

Η παράμετρος α υπολογίζεται λύνοντας την εξίσωση $\bar{Y} = \alpha + b\bar{X}$ ως προς α , δηλαδή:

$$\alpha = \bar{Y} - b\bar{X}$$

Τώρα μπορούμε να υπολογίσουμε τις εκτιμήτριες \hat{Y}_i από την εξίσωση:

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

Απαραίτητες προϋποθέσεις για την παλινδρόμηση:

1. Οι τιμές των μεταβλητών X και Y πρέπει να είναι ανεξάρτητες μεταξύ τους.
2. Για κάθε τιμή των X στον πληθυσμό, οι τιμές της Y πρέπει να ακολουθούν την Κανονική κατανομή.
3. Οι διασπορές των τιμών της Y πρέπει να είναι ίσες.
4. Η σχέση μεταξύ των X και Y πρέπει να είναι γραμμική.
5. Οι μετρήσεις των τιμών της X μεταβλητής δεν πρέπει συμπεριλαμβάνουν σφάλματα.

3.2.6 Έλεγχοι Υποθέσεων Με Βάση Την Παλινδρόμηση:

Για να προχωρήσουμε στους ελέγχους υποθέσεων θα πρέπει πρώτα να υπολογίσουμε το σφάλμα των εκτιμητριών \hat{Y}_i . Θα χρειαστούμε τους παρακάτω τύπους:

- $\sum_{i=1}^n (y_i)^2 = \sum_{i=1}^n (Y_i - \bar{Y})^2 = \sum_{i=1}^n (Y_i)^2 - \frac{(\sum_{i=1}^n Y_i)^2}{n}$
- $\frac{(\sum_{i=1}^n x_i y_i)^2}{\sum_{i=1}^n (x_i)^2} = \frac{(\sum_{i=1}^n x_i Y_i - \frac{(\sum_{i=1}^n x_i)(\sum_{i=1}^n Y_i)}{n})^2}{\sum_{i=1}^n (x_i)^2 - \frac{(\sum_{i=1}^n x_i)^2}{n}} = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$

$$\begin{aligned}
 \bullet \quad r^2 &= \frac{\frac{(\sum_{i=1}^n X_i Y_i - \frac{(\sum_{i=1}^n X_i)(\sum_{i=1}^n Y_i)}{n})^2}{\sum_{i=1}^n (X_i)^2 - \frac{(\sum_{i=1}^n X_i)^2}{n}}}{\frac{\sum_{i=1}^n (Y_i)^2 - \frac{(\sum_{i=1}^n Y_i)^2}{n}}{\sum_{i=1}^n (Y_i - \bar{Y})^2}} = \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} \\
 \bullet \quad S_Y^2 &= \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y}_n)^2 \\
 \bullet \quad S_{YX}^2 &= S_Y^2 \frac{(1-r^2)(n-1)}{(n-2)}
 \end{aligned}$$

Το σφάλμα των εκτιμητριών \hat{Y}_i , λοιπόν, υπολογίζεται από τον παρακάτω τύπο:

$$S_{\hat{Y}_i} = \sqrt{S_{YX}^2 \left[\frac{1}{n} + \frac{(X_i - \bar{X})^2}{\sum_{i=1}^n (x_i)^2} \right]}$$

(προσοχή, υπολογίζουμε το σφάλμα για κάθε μια εκτιμήτρια ξεχωριστά).

Μετά τον υπολογισμό των εκτιμητριών \hat{Y}_i και των αντίστοιχων σφαλμάτων τους $S_{\hat{Y}_i}$ μπορούμε να προχωρήσουμε στους ελέγχους υποθέσεων (που αφορούν π.χ. την μέση τιμή $\mu_{\hat{Y}_i}$ της \hat{Y}_i), χρησιμοποιώντας τον μονόπλευρο ή αμφίπλευρο έλεγχο t. Εφαρμόζοντας την συγκεκριμένη ονοματολογία της παραγράφου θα πάρουμε τους αντίστοιχους ελέγχους:

A) Για τον αμφίπλευρο έλεγχο t, της υπόθεσης $H_0: \mu_{\hat{Y}_i} = \mu_0$ vs $H_1: \mu_{\hat{Y}_i} \neq \mu_0$:

$$d(\hat{Y}_i) = 1(|t| > t_{n-1}(1 - \frac{\alpha}{2}))$$

B) Για τον μονόπλευρο έλεγχο t, της υπόθεσης $H_0: \mu_{\hat{Y}_i} \leq \mu_0$ vs $H_1: \mu_{\hat{Y}_i} > \mu_0$:

$$d(\hat{Y}_i) = 1(t > t_{n-1}(1 - \alpha))$$

Όπου $t = \frac{\hat{Y}_i - \mu_0}{S_{\hat{Y}_i}}$. [32]

3.2.7 Απλή Γραμμική Συσχέτιση (συντελεστής συσχέτισης):

Σε αυτή την περίπτωση όπως αναφέραμε και πριν ασχολούμαστε με τις μεταβλητές X και Y για τις οποίες δεν απαιτούμε να έχουν κάποια σχέση εξάρτησης,

δηλαδή αυτού του είδους η ανάλυση δίνει τα ίδια αποτελέσματα ανεξάρτητα με το ποιά μεταβλητή θα βάλουμε ως X και ποιά ως Y .

Ο συντελεστής συσχέτισης (correlation coefficient) είναι ο ακόλουθος:

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

Όπου οι τιμές των $\sum_{i=1}^n x_i y_i$, $\sum_{i=1}^n (x_i)^2$ και $\sum_{i=1}^n (y_i)^2$ υπολογίζονται από τις προαναφερθείσες εξισώσεις.

Ο συντελεστής συσχέτισης r μεταβάλλεται στο διάστημα $[-1,1]$ διότι ο αριθμητής είναι πάντα μεγαλύτερος από τον παρανομαστή και παίρνει και αρνητικές και θετικές τιμές. Αν ο συντελεστής συσχέτισης είναι αρνητικός τότε η αύξηση του μεγέθους της μιας μεταβλητής συνεπάγεται την μείωση του μεγέθους της άλλης, αν είναι θετικός τότε η αύξηση του μεγέθους της μιας μεταβλητής συνεπάγεται την αύξηση του μεγέθους της άλλης, και αν είναι μηδέν τότε η αλλαγή του μεγέθους της μιας μεταβλητής δεν συνεπάγεται την αλλαγή του μεγέθους της άλλης.

Το σφάλμα του συντελεστή συσχέτισης είναι το:

$$S_r = \sqrt{\frac{1-r^2}{n-2}}$$

όπου n είναι το πλήθος του δείγματος.

Χρειαζόμαστε μόνο μια προϋπόθεση για την ανάλυση με συντελεστή συσχέτισης, η οποία είναι ότι για κάθε τιμή των X στον πληθυσμό, οι τιμές της Y πρέπει να ακολουθούν την Κανονική κατανομή και για κάθε τιμή των Y στον πληθυσμό, οι τιμές της X πρέπει να ακολουθούν την Κανονική κατανομή ομοίως.

Ο έλεγχος, λοιπόν, της υπόθεσης του αν υπάρχει συσχέτιση μεταξύ των X και Y στον πληθυσμό, ή ποιο συγκεκριμένα της υπόθεσης $H_0: \rho = 0$ vs $H_1: \rho \neq 0$, μεγέθους α , είναι ο ακόλουθος:

$$d(\underline{X}, \underline{Y}) = 1(t < t_{n-2}(1 - \frac{\alpha}{2}))$$

$$\text{όπου } t = \frac{r}{s_r}$$

Δηλαδή αν $t < t_{n-2}(1 - \frac{\alpha}{2})$ δεχόμαστε την H_0 , άρα δεν υπάρχει συσχέτιση μεταξύ των μεταβλητών.[32]

3.2.8 Λογιστική παλινδρόμηση:

Τα τελευταία χρόνια η μέθοδος της λογιστικής παλινδρόμησης αποτελεί την κύρια μέθοδο ανάλυσης σε πολλά πεδία, όπως η βιοστατιστική ή η οικονομία. [37]. Η λογιστική παλινδρόμηση είναι μια τεχνική σχεδιασμένη για την ανάλυση δεδομένων που αφορούν την μελέτη και την πρόβλεψη των τιμών κάποιας κατηγορικής εξαρτημένης μεταβλητής και χρησιμοποιεί ποσοτικές ή ποιοτικές ανεξάρτητες μεταβλητές. Δυστυχώς όμως, για την ανάλυση των δεδομένων μας μέσω της λογιστικής παλινδρόμησης, δεν μπορούμε να χρησιμοποιήσουμε τον αλγόριθμο της απλής γραμμικής παλινδρόμησης διότι οι τιμές που παίρνει η μεταβλητή μας με το μοντέλο της λογιστικής παλινδρόμησης είναι 0 ή 1 (σε αντίθεση με αυτές του μοντέλου της απλής γραμμικής παλινδρόμησης που είναι μεγαλύτερες του 1 ή μικρότερες του 0) και αυτό διότι στην ουσία υπολογίζουμε την πιθανότητα με την οποία η εξαρτημένη μεταβλητή θα λάβει κάποια συγκεκριμένη τιμή [38].

Η εξίσωση της λογιστικής παλινδρόμησης για ένα δείγμα n παρατηρήσεων είναι η $g(X_i) = \beta_0 + \beta_1 X_i$, όπου το $g(X_i)$, το οποίο συχνά γράφεται και ως $\text{Ln}(odds)$ [Norušis, 1997 #909], είναι ίσο με $\ln(\frac{\pi(X_i)}{1-\pi(X_i)})$ όπου $\pi(X_i) = \frac{e^{\beta_0 + \beta_1 X_i}}{1 + e^{\beta_0 + \beta_1 X_i}}$. Ακόμα, ως Y_i συμβολίζουμε την μεταβλητή που παίρνει τις τιμές 0 ή 1 και αντιστοιχεί στην παρατήρηση X_i .

Για να μπορέσουμε να λύσουμε την εξίσωση της λογιστικής παλινδρόμησης θα πρέπει πρώτα να εκτιμήσουμε τις τιμές των β_i . Πριν προχωρήσουμε στην εκτίμηση των β_i όμως, θα θέλαμε πρώτα να αναφερθούμε σε μια διαφορά της

λογιστικής και της απλής γραμμικής παλινδρόμησης. Στην περίπτωση της απλής γραμμικής παλινδρόμησης, είδαμε πως μπορούμε να γράψουμε την εξίσωση μας στην μορφή $Y_i = \alpha + \beta X_i + \varepsilon_i$ (ή αλλιώς $Y_i = \dot{Y}_i + \varepsilon_i$) όπου τα ε_i τα ονομάσαμε errors. Στην περίπτωση της λογιστικής παλινδρόμησης η εξίσωση μας γίνεται $Y_i = \pi(X_i) + \varepsilon_i$ όπου αν $Y_i = 1$ τότε $\varepsilon_i = 1 - \pi(X_i)$ και αν $Y_i = 0$ τότε $\varepsilon_i = \pi(X_i)$. Το ε_i ακολουθεί λοιπόν μια κατανομή με μέση τιμή 0 και διασπορά $\pi(X_i)(1 - \pi(X_i))$, δηλαδή ακολουθεί την Διωνυμική κατανομή με $p = \pi(X_i)$, σε αντίθεση με το error της απλής γραμμικής που ακολουθεί την Κανονική κατανομή.

Για την εκτίμηση των β_i τώρα, θα εργαστούμε με παρόμοιο τρόπο με αυτόν που εργαστήκαμε στην απλή γραμμική παλινδρόμηση. Στην απλή γραμμική παλινδρόμηση χρησιμοποιήσαμε την μέθοδο ελαχίστων τετραγώνων, ενώ εδώ θα χρησιμοποιήσουμε μια άλλη μέθοδο διότι η μέθοδος των ελαχίστων τετραγώνων δεν δίνει καλούς εκτιμητές αν εφαρμοστεί στην λογιστική παλινδρόμηση. Η μέθοδος που θα χρησιμοποιήσουμε βασίζεται στην μέθοδο των μεγίστων πιθανοφανειών.

Εφόσον οι μεταβλητές Y_i παίρνουν τις τιμές 0 ή 1, τότε το $\pi(X_i)$ εκφράζει την δεσμευμένη πιθανότητα του $Y_i = 1$ δεδομένου του X_i , δηλαδή $P(Y_i = 1|X_i) = \pi(X_i)$ και άρα $P(Y_i = 0|X_i) = 1 - \pi(X_i)$. Επομένως η συν-κατανομή των ζευγών (Y_i, X_i) είναι η $(\pi(X_i))^{Y_i} (1 - \pi(X_i))^{1-Y_i}$ και γι' αυτόν τον λόγο η πιθανοφάνεια είναι η ακόλουθη:

$$l(\beta) = \prod_{i=1}^n (\pi(X_i))^{Y_i} (1 - \pi(X_i))^{1-Y_i}$$

όπου $\beta = (\beta_0, \beta_1)$.

Η Log- πιθανοφάνεια είναι η:

$$\text{Log}(l(\beta)) = \sum_{i=1}^n [Y_i \text{Ln}(\pi(X_i)) + (1 - Y_i) \text{Ln}(1 - \pi(X_i))].$$

Παραγωγίζοντας, λοιπόν, την Log- πιθανοφάνεια ως προς β_0 και β_1 και εξισώνοντας τα αποτελέσματα με το 0, παίρνουμε αντίστοιχα:

$$\sum_{i=1}^n [Y_i - \pi(X_i)] = 0 (\alpha)$$

και

$$\sum_{i=1}^n X_i [Y_i - \pi(X_i)] = 0 (\beta)$$

Αυτές οι δύο εξισώσεις δεν είναι γραμμικές με τα β_0 και β_1 , και για να επιλυθούν επιβάλλεται η χρήση κατάλληλων προγραμμάτων στον υπολογιστή (logistic regression softwares).

Λύνοντας λοιπόν τις (α) και (β), βρίσκουμε τις εκτιμήτριες $\hat{\beta}_0$ και $\hat{\beta}_1$ των β_0 και β_1 αντίστοιχα (ή αλλιώς την εκτιμήτρια $\hat{\beta}$ του β), τώρα μπορούμε να προχωρήσουμε στον έλεγχο υποθέσεων. Παρατηρούμε ότι αν $\hat{\pi}(X_i)$ είναι η εκτιμήτρια του $\pi(X_i)$, τότε $\sum_{i=1}^n Y_i = \sum_{i=1}^n \hat{\pi}(X_i)$. Δηλαδή το άθροισμα των παρατηρούμενων τιμών του Y είναι ίσο με το άθροισμα των εκτιμώμενων. [37]

3.2.9 Έλεγχοι Υποθέσεων Με Βάση Την Λογιστική Παλινδρόμηση:

Η αρχή της λογιστικής παλινδρόμησης είναι ότι συγκρίνουμε τις παρατηρούμενες μεταβλητές της εξαρτημένης μεταβλητής για να προβλέψουμε τις τιμές του μοντέλου που έχει την ανεξάρτητη μεταβλητή και αυτού που δεν την έχει. Η σύγκριση των παρατηρούμενων με των προβλεπόμενων μεταβλητών με την βοήθεια της πιθανοφάνειας βασίζεται στην παρακάτω έκφραση:

$$D = -2 \ln \left[\frac{\text{πιθανοφάνεια κανονικού μοντέλου}}{\text{πιθανοφάνεια κορεσμένου μοντέλου}} \right]$$

$$= -2 \sum_{i=1}^n \left[y_i \ln \left(\frac{\hat{\pi}(x_i)}{y_i} \right) + (1 - y_i) \ln \left(\frac{1 - \hat{\pi}(x_i)}{1 - y_i} \right) \right]$$

Όπου το D παίζει τον ίδιο ρόλο με αυτόν του $\sum_{i=1}^n (y_i)^2$ στην απλή γραμμική παλινδρόμηση. Ακόμα με τον όρο "κορεσμένο μοντέλο" εννοούμε το μοντέλο που έχει τόσες παραμέτρους όσο το πλήθος του δείγματος, δηλαδή για το δικό μας

μοντέλο που έχει δύο παραμέτρους β_0 και β_1 , θα είχαμε κορεσμό μόνο για $n=2$. Τέλος, επειδή εμάς μας ενδιαφέρει κυρίως η περίπτωση όπου η μεταβλητή Y παίρνει τιμές 0 ή 1, η πιθανοφάνεια του κορεσμένου μοντέλου είναι 1, επομένως το D γίνεται:

$$D = -2Ln[\text{πιθανοφάνεια κανονικού μοντέλου}]$$

Πρίν την διατύπωση του ελέγχου για το μοντέλο της λογιστικής παλινδρόμησης, πρέπει να υπολογίσουμε την στατιστική συνάρτηση G , στην οποία θα βασιστούμε για να διατυπώσουμε τον έλεγχο. Η G υπολογίζεται από τον τύπο:

$$G = D(1) - D(2) = -2Ln \left[\frac{\text{πιθανοφάνεια χωρίς την ανεξάρτητη μεταβλητή στην εξίσωση}}{\text{πιθανοφάνεια με την ανεξάρτητη μεταβλητή στην εξίσωση}} \right] =$$

$$= -2Ln \left[\frac{\left(\frac{\sum_{i=1}^n y_i}{n} \right)^{\sum_{i=1}^n y_i} \left(\frac{\sum_{i=1}^n (1 - y_i)}{n} \right)^{\sum_{i=1}^n (1 - y_i)}}{\prod_{i=1}^n (\hat{\pi}(x_i))^{y_i} (1 - \hat{\pi}(x_i))^{(1 - y_i)}} \right]$$

Όπου $D(1)=D(\text{Μοντέλο χωρίς την ανεξάρτητη μεταβλητή στην εξίσωση})$ και

$D(2)=D(\text{Μοντέλο με την ανεξάρτητη μεταβλητή στην εξίσωση})$. Κάτω από την υπόθεση $H_0:\beta_1 = 0$ vs $H_1:\beta_1 \neq 0$, η στατιστική συνάρτηση G ακολουθεί κατανομή χ_1^2 . Τέλος, ο έλεγχος της υπόθεσης μας είναι ο ακόλουθος:

$$d(Y)=1(G > \chi_1^2(1 - \alpha)).$$

Επομένως απορρίπτουμε την H_0 όταν $G > \chi_1^2(1 - \alpha)$. [37]

3.3 Επιλογή κατάλληλης στατιστικής μεθοδολογίας

Ο κύριος σκοπός της παρούσας πτυχιακής εργασίας ήταν η διερεύνηση της επίδρασης της ενεργητικής ανοσοποίησης (εμβολιασμών) στην νευροανάπτυξη των παιδιών ηλικίας έως 18 μηνών η οποία έγινε στα πλαίσια της μελέτης μητέρας-παιδιού ΡΕΑ.

Η νοητική και ψυχοκινητική ανάπτυξη των παιδιών εκτιμήθηκε στους 18 μήνες (± 6 εβδομάδες) χρησιμοποιώντας τις κλίμακες του Bayley για τα βρέφη και την νηπιακή ανάπτυξη [39]. Η μέθοδος Bayley-III αξιολογεί την αναπτυξιακή λειτουργία των βρεφών και των μικρών παιδιών ηλικίας 1 μηνός έως και 42. Ο πρωταρχικός σκοπός είναι ο προσδιορισμός των παιδιών με αναπτυξιακή καθυστέρηση και η παροχή πληροφοριών για το σχεδιασμό παρεμβάσεων.

Η μέθοδος Bayley-III αξιολογεί την ανάπτυξη του βρέφους και του παιδιού σε πέντε τομείς: Η αξιολόγηση της ανάπτυξης του βρέφους και του παιδιού, μέσω της μεθόδου Bayley-III, γίνεται σε 5 κλίμακες: (i) Η Γνωσιακή Κλίμακα (COG) περιλαμβάνει στοιχεία τα οποία εκτιμούν την αισθητικοκινητική ανάπτυξη, την διερεύνηση και την χειραγώγηση, τη συνάφεια του αντικειμένου, τον σχηματισμό έννοιας, τη μνήμη, και άλλες πτυχές της νοητικής επεξεργασίας. (ii) Η Κλίμακα Γλώσσας αποτελείται από δύο υπό-test, το test της Δεκτικής Επικοινωνίας (RC) και αυτό της Εκφραστικής Επικοινωνίας (EK). Το RC test περιλαμβάνει τα στοιχεία που αξιολογούν συμπεριφορές παιδιών που δεν έχουν μιλήσει ακόμα, την ανάπτυξη του λεξιλογίου, το λεξιλόγιο που σχετίζεται με την μορφολογική ανάπτυξη, την κατανόηση των μορφολογικών δεικτών, την κοινωνική και προφορική κατανόηση των παιδιών. Το EK test περιλαμβάνει τα στοιχεία που αξιολογούν την επικοινωνία των παιδιών που δεν έχουν μιλήσει ακόμα, την ανάπτυξη του λεξιλογίου, και την μορφο-συντακτική ανάπτυξη. (iii) Η κλίμακα της κίνησης διαιρείται στα test της Λεπτής Κινητικότητας [Fine Motor(FM)] και της Αδρής Κινητικότητας [Gross Motor(GM)]. Στο FM test περιλαμβάνονται κινητικές δεξιότητες που συνδέονται με την σύλληψη και αντίληψη της ολοκλήρωσης της κίνησης, τον κινητικό σχεδιασμό, και τη ταχύτητα της κίνησης. Το GM test μετράει κατά κύριο λόγο την κίνηση των άκρων και του κορμού. (iv) Η κοινωνικο-συναισθηματική κλίμακα (SE), αξιολογεί την απόκτηση κοινωνικών και συναισθηματικών οροσήμων σε βρέφη και μικρά παιδιά. Προσδιορίζει τα μεγάλα αναπτυξιακά ορόσημα που πρέπει να επιτευχθούν σε ορισμένες ηλικίες. Περιλαμβάνει στοιχεία που αξιολογούν την ικανότητα του παιδιού στις συναισθηματικές δεξιότητες, στις ανάγκες επικοινωνίας, στην ικανότητα δημιουργίας σχέσεων με άλλα άτομα, στη χρήση συναισθημάτων με διαδραστικό χαρακτήρα, και στη χρήση συναισθηματικών σημάτων ή χειρονομιών

για την επίλυση προβλημάτων. (v) Η κλίμακα Προσαρμογής Συμπεριφοράς αξιολογεί τις καθημερινές λειτουργικές ικανότητες του παιδιού, μετρώντας αυτό που πράγματι κάνει το παιδί σε σχέση με αυτό που είναι ικανό να κάνει. Οι τομείς που μετρώνται στο πλαίσιο αυτής της κλίμακας είναι η επικοινωνία, η χρήση της κοινότητας, της υγείας και της ασφάλειας, η αναψυχή, η αυτοφροντίδα, η αυτοκατεύθυνση, οι συνθήκες διαβίωσης στο σπίτι, η κοινωνικότητα, και τέλος η κίνηση.

Η διαπολιτισμική προσαρμογή της κοινωνικο-συναισθηματικής κλίμακας του Bayley-III πραγματοποιήθηκε σύμφωνα με την διεθνώς προτεινόμενη μεθοδολογία, χρησιμοποιώντας τις ακόλουθες οδηγίες: μετάφραση προς τα εμπρός (αγγλικά προς ελληνικά), μετάφραση προς τα πίσω (από το ήδη μεταφρασμένο στα ελληνικά ερωτηματολόγιο σε αγγλικά), γνωστική διαδικασία απολογισμού, και προέλεγχος.

Για κάθε κλίμακα, η βαθμολογία του παιδιού έχει καθοριστεί από τον αριθμό των αντικειμένων για τα οποία έχει ληφθεί βαθμός. Αναλύσαμε πρωταρχικές βαθμολογίες αντί κλιμακωτές και σύνθετες βαθμολογίες, διότι το δείγμα των Ηνωμένων Πολιτειών μπορεί να μην είναι κατάλληλο για τα παιδιά έξω από της Ηνωμένες Πολιτείες. Τα πρωταρχικά αποτελέσματα ήταν τυποποιημένα με μέση τιμή 100 και τυπική απόκλιση 15, ώστε να ομογενοποιηθούν όλες τις κλίμακες.

Οι νευροαναπτυξιακές εκτιμήσεις έγιναν από τρεις εκπαιδευμένους ψυχολόγους, οι οποίοι ολοκλήρωσαν την επίσημη εκπαίδευση της μεθόδου Bayley-III. Όλα τα test έγιναν στην Ιατρική Σχολή του Πανεπιστημίου Κρήτης (UOC), σε δύο δημόσια νοσοκομεία στο Ηράκλειο και σε Ιατρικά Κέντρα Υγείας στην επαρχία, πάντα με την παρουσία της μητέρας. Ο εξεταστής 1 αξιολόγησε 253 βρέφη, ο εξεταστής 2 αξιολόγησε 312 και ο εξεταστής 3 αξιολόγησε 40. Οι εξεταστές, επίσης, σημείωναν κριτικά σχόλια σχετικά με τις δυσκολίες ή τις ειδικές συνθήκες της εκτίμησης της νευροανάπτυξης για να αξιολογηθεί η "ποιότητα της αξιολόγησης": δεν υπάρχουν δυσκολίες, τα βρέφη συμμετείχαν στη μελέτη της αξιοπιστίας, δυσκολίες εξαιτίας σωματικών προβλημάτων (π.χ. κόπωση, ασθένεια, ύπνος, κλπ.), δυσκολίες που οφείλονται σε προβλήματα συμπεριφοράς (π.χ. νευρικότητα, ντροπαλός/η, κλπ.).

Οι μητέρες ενημερώθηκαν μέσα σε ένα μήνα για τα αποτελέσματα των παιδιών τους, μέσω ηλεκτρονικού ταχυδρομείου [1].

Οι 5 κλίμακες της μεθόδου Bayley-III παίρνουν τιμές από 0 μέχρι 100 αναλόγως με την βαθμολογία του αντίστοιχου test. Δηλαδή οι μεταβλητές μας είναι συνεχείς. Ακόμα έχει γίνει κοινωνικοποίηση όλων το μεταβλητών ούτως ώστε να ακολουθούν την Κανονική κατανομή.

Μία κατάλληλη μέθοδος ανάλυσης των δεδομένων μας είναι αυτή της γραμμικής παλινδρόμησης. Όπου τα X_j για $j=1, \dots, 5$ αντιπροσωπεύουν τις 5 κλίμακες της μεθόδου Bayley-III και παίρνουν τιμές από 0 μέχρι 100, και η εξαρτημένη μεταβλητή Y μετράει την σχέση ενεργητικής ανοσοποίησης (εμβολιασμών) και νευροανάπτυξης των παιδιών. Δηλαδή το μοντέλο μας είναι το:

$$Y_i = \alpha + \beta X_1 + \gamma X_2 + \delta X_3 + \varepsilon X_4 + \zeta X_5$$

όπου το i είναι το μέγεθος του δείγματος. Το μοντέλο αυτό γράφεται στην μορφή πινάκων

$$\underline{Y} = \underline{X}\underline{\beta} + \underline{\varepsilon}$$

Όπου $\underline{Y} = \begin{bmatrix} Y_1 \\ \dots \\ Y_n \end{bmatrix}$ και n είναι το μέγεθος του δείγματος, $\underline{\beta} = \begin{bmatrix} \alpha \\ \dots \\ \zeta \end{bmatrix}$ και

$$\underline{X} = \begin{bmatrix} 1 & X_{11} & \dots & X_{51} \\ \dots & \dots & \dots & \dots \\ 1 & X_{1n} & \dots & X_{5n} \end{bmatrix}$$

Τέλος ο πίνακας $\underline{\beta}$ υπολογίζεται από την εξίσωση:

$$\underline{\beta} = (\underline{X}'\underline{X})^{-1}\underline{X}'\underline{Y}$$

Εναλλακτικά, γνωρίζοντας ότι η μέση τιμή ενός δείγματος, του οποίου οι μεταβλητές ακολουθούν την Κανονική κατανομή, είναι μία καλή εκτιμήτρια θέσης του δείγματος, μπορούμε να δημιουργήσουμε μία δίτιμη μεταβλητή και να εφαρμόσουμε Λογιστική παλινδρόμηση.

Για τον λόγο αυτό βρίσκουμε την διάμεσο της κάθε κλίμακας αντίστοιχα. Ορίζουμε, λοιπόν, την κάθε διάμεσο ως εκείνο το σημείο πάνω στην αντίστοιχη κλίμακα, για το οποίο όλες οι τιμές που βρίσκονται άνωθεν του θεωρούνται "επιτυχία", και όλες οι τιμές που βρίσκονται κάτωθεν του θεωρούνται "αποτυχία". Δηλαδή, εκείνες οι μεταβλητές που έχουν τιμή μεγαλύτερη της διαμέσου θα πάρουν την τιμή 1 και θεωρούνται "επιτυχίες", και αντίστοιχα εκείνες οι μεταβλητές οι οποίες έχουν τιμή μικρότερη της διαμέσου θα πάρουν την τιμή 0 και θεωρούνται "αποτυχίες". Αυτό σημαίνει ότι έχουμε κατασκευάσει ένα μοντέλο 2 μεταβλητών (Binary model) του οποίου οι μεταβλητές παίρνουν τις τιμές 0 και 1.

Κεφάλαιο 4

Συζήτηση:

Η νευροανάπτυξη των παιδιών επηρεάζεται από διάφορους παράγοντες όπως η διατροφή και το περιβάλλον. Τα τελευταία χρόνια επικρατεί η άποψη ότι τα εμβόλια προσχολικής ηλικίας ίσως είναι ένας σημαντικός παράγοντας που θα μπορούσε να επηρεάσει την διαμόρφωση του νευρικού συστήματος.

Στόχος της παρούσας πτυχιακής εργασίας ήταν η ανάλυση μερικών παραμετρικών test, και των αντίστοιχων ελέγχων τους, και η επιλογή του καταλληλότερου για τον έλεγχο της σχέσης ενεργητικής ανοσοποίησης (εμβολιασμών) και νευροανάπτυξης σε παιδιά προσχολικής ηλικίας στα πλαίσια της προοπτικής μελέτης μητέρας – παιδιού Ρέα, η οποία πραγματοποιείται στο Ηράκλειο Κρήτης. Τα test που αναλύθηκαν είναι: το χ^2 test, το t-test, ο συντελεστής γραμμικής συσχέτισης, η απλή γραμμική παλινδρόμηση και η λογιστική παλινδρόμηση.

Τα δεδομένα της μελέτης συλλέγονται με βάση την μέθοδο Bayley-III. Η μέθοδος Bayley-III αποτελείται από 5 test που γίνονται στα παιδιά και στις μητέρες. Τα αποτελέσματα των 5 test σημειώνονται πάνω σε κλίμακες, τις κλίμακες του Bayley, οι οποίες παίρνουν τις τιμές από 0 μέχρι 100. Καταλήξαμε ότι μία κατάλληλη μέθοδος είναι αυτή της Γραμμικής παλινδρόμησης, ενώ μετά από κανονικοποίηση των δεδομένων που πήραμε από τις κλίμακες και κατασκευάζοντας ένα μοντέλο δύο μεταβλητών οι οποίες παίρνουν τις τιμές 0 και 1, καταλήξαμε ότι μία άλλη καλή μέθοδος είναι εκείνη της Λογιστικής παλινδρόμησης.

Η στατιστική ανάλυση είναι πολύ σημαντική για διάφορες βιο-ιατρικές μελέτες που πραγματοποιούνται στον κόσμο. Ο κλάδος της στατιστικής που ασχολείται με τέτοιου είδους μελέτες ονομάζεται βιοστατιστική και είναι ευρέως διαδεδομένος τα τελευταία χρόνια. Η επιλογή του κατάλληλου ελέγχου – test γίνεται αναλόγως τα δεδομένα και τον στόχο της εκάστοτε μελέτης.

ΠΗΓΕΣ

1. Koutra K, Chatzi L, Roumeliotaki T, Vassilaki M, Giannakopoulou E, Batsos C, Koutis A, Kogevas M: **Socio-demographic determinants of infant neurodevelopment at 18 months of age: Mother-Child Cohort (Rhea Study) in Crete, Greece.** *Infant Behav Dev* 2012, **35**(1):48-59.
2. Chen RT, Mootrey G, DeStefano F: **Safety of routine childhood vaccinations. An epidemiological review.** *Paediatr Drugs* 2000, **2**(4):273-290.
3. Currenti SA: **Understanding and determining the etiology of autism.** *Cell Mol Neurobiol* 2010, **30**(2):161-171.
4. **Εθνικό Πρόγραμμα Εμβολιασμών-Η Ιατρική Σήμερα.** 2007, **50**:23.
5. Rask-Nissila L, Jokinen E, Terho P, Tammi A, Lapinleimu H, Ronnema T, Viikari J, Seppanen R, Korhonen T, Tuominen J *et al*: **Neurological development of 5-year-old children receiving a low-saturated fat, low-cholesterol diet since infancy: A randomized controlled trial.** *JAMA* 2000, **284**(8):993-1000.
6. Dauncey MJ, Bicknell RJ: **Nutrition and neurodevelopment: mechanisms of developmental dysfunction and disease in later life.** *Nutr Res Rev* 1999, **12**(2):231-253.
7. Bourre JM: **Roles of unsaturated fatty acids (especially omega-3 fatty acids) in the brain at various ages and during ageing.** *J Nutr Health Aging* 2004, **8**(3):163-174.
8. Neuringer M, Reisbick S, Janowsky J: **The role of n-3 fatty acids in visual and cognitive development: current evidence and methods of assessment.** *J Pediatr* 1994, **125**(5 Pt 2):S39-47.
9. Kretchmer N, Beard JL, Carlson S: **The role of nutrition in the development of normal cognition.** *Am J Clin Nutr* 1996, **63**(6):997S-1001S.

10. JJ. Strain MPB, E. M. Duffy, J. MW. Wallace, P. J. Robson, T. W. Clarkson and C. Shamlaye **Nutrition and neurodevelopment: the search for candidate nutrients in the Seychelles Child Development Nutrition Study.** *Seychelles Medical and Dental Journal* 2004, **7**:77-83
11. Bourre JM: **Effects of nutrients (in food) on the structure and function of the nervous system: update on dietary requirements for brain. Part 1: micronutrients.** *J Nutr Health Aging* 2006, **10**(5):377-385.
12. Safety IPoC: **Polychlorinated biphenyls and terphenyls.** 1993, **2**.
13. Ribas-Fito N, Sala M, Kogevinas M, Sunyer J: **Polychlorinated biphenyls (PCBs) and neurological development in children: a systematic review.** *J Epidemiol Community Health* 2001, **55**(8):537-546.
14. Carpenter DO: **Polychlorinated biphenyls and human health.** *Int J Occup Med Environ Health* 1998, **11**(4):291-303.
15. Dewailly E, Mulvad G, Pedersen HS, Ayotte P, Demers A, Weber JP, Hansen JC: **Concentration of organochlorines in human brain, liver, and adipose tissue autopsy samples from Greenland.** *Environ Health Perspect* 1999, **107**(10):823-828.
16. Geier DA, Geier MR: **A meta-analysis epidemiological assessment of neurodevelopmental disorders following vaccines administered from 1994 through 2000 in the United States.** *Neuro Endocrinol Lett* 2006, **27**(4):401-413.
17. Miller DL, Ross EM, Alderslade R, Bellman MH, Rawson NS: **Pertussis immunisation and serious acute neurological illness in children.** *Br Med J (Clin Res Ed)* 1981, **282**(6276):1595-1599.

18. Miller E, Andrews N, Stowe J, Grant A, Waight P, Taylor B: **Risks of convulsion and aseptic meningitis following measles-mumps-rubella vaccination in the United Kingdom.** *Am J Epidemiol* 2007, **165**(6):704-709.
19. Claudio L, Kwa WC, Russell AL, Wallinga D: **Testing methods for developmental neurotoxicity of environmental chemicals.** *Toxicol Appl Pharmacol* 2000, **164**(1):1-14.
20. McMahon AW, Eidex RB, Marfin AA, Russell M, Sejvar JJ, Markoff L, Hayes EB, Chen RT, Ball R, Braun MM *et al*: **Neurologic disease associated with 17D-204 yellow fever vaccination: a report of 15 cases.** *Vaccine* 2007, **25**(10):1727-1734.
21. Grandjean P, Landrigan PJ: **Developmental neurotoxicity of industrial chemicals.** *Lancet* 2006, **368**(9553):2167-2178.
22. Dorea JG: **Making sense of epidemiological studies of young children exposed to thimerosal in vaccines.** *Clin Chim Acta* 2010, **411**(21-22):1580-1586.
23. Ray P, Hayward J, Michelson D, Lewis E, Schwalbe J, Black S, Shinefield H, Marcy M, Huff K, Ward J *et al*: **Encephalopathy after whole-cell pertussis or measles vaccination: lack of evidence for a causal association in a retrospective case-control study.** *Pediatr Infect Dis J* 2006, **25**(9):768-773.
24. Sejvar JJ, Labutta RJ, Chapman LE, Grabenstein JD, Iskander J, Lane JM: **Neurologic adverse events associated with smallpox vaccination in the United States, 2002-2004.** *JAMA* 2005, **294**(21):2744-2750.
25. Monteiro SA, Takano OA, Waldman EA: **Surveillance for adverse events after DTwP/Hib vaccination in Brazil: sensitivity and factors associated with reporting.** *Vaccine* 2010, **28**(18):3127-3133.
26. Geier DA, Geier MR: **A two-phased population epidemiological study of the safety of thimerosal-containing vaccines: a follow-up analysis.** *Med Sci Monit* 2005, **11**(4):CR160-170.

27. Camargo M, Soto-De Leon SC, Sanchez R, Perez-Prados A, Patarroyo ME, Patarroyo MA: **Frequency of human papillomavirus infection, coinfection, and association with different risk factors in Colombia.** *Ann Epidemiol* 2011, **21**(3):204-213.
28. Firnhaber C, Van Le H, Pettifor A, Schulze D, Michelow P, Sanne IM, Lewis DA, Williamson AL, Allan B, Williams S *et al*: **Association between cervical dysplasia and human papillomavirus in HIV seropositive women from Johannesburg South Africa.** *Cancer Causes Control* 2010, **21**(3):433-443.
29. Ong BA, Forester J, Fallot A: **Does influenza vaccination improve pediatric asthma outcomes?** *J Asthma* 2009, **46**(5):477-480.
30. Mikaeloff Y, Caridade G, Suissa S, Tardieu M: **Hepatitis B vaccine and the risk of CNS inflammatory demyelination in childhood.** *Neurology* 2009, **72**(10):873-880.
31. Χαραλαμπίδη ΧΑ: **Θεωρία Πιθανοτήτων Και Εφαρμογές**, vol. 2: Εκδόσεις Συμμετρία
1999.
32. H.Zar J: **Biostatistical Analysis**: Pearson Education.
33. Κλωνιάς Β: **Παραμετρικής Στατιστικής**: Πανεπιστημιακές Εκδόσεις Κρήτης-media; 2007.
34. Ρουσσά ΓΓ: **Στατιστική Συμπερασματολογία**: Εκδόσεις ΖΗΤΗ; 1994.
35. Κλωνιάς Β: **Εφηρμοσμένη Στατιστική**: Πανεπιστημιακές Εκδόσεις Κρήτης-media; 2007.
36. Lehmann E: **Fisher, neyman, and the creation of classical statistics.** New York: Springer; 2011.
37. Hosmer DW, Lemeshow S: **Applied logistic regression**, 2nd edn. New York: Wiley; 2000.

38. Norušis MJ, SPSS Inc.: **SPSS professional statistics**. In., [] 7.5 for Windows. edn. Chicago, Ill.: SPSS,; 1997.
39. Milne S, McDonald J, Comino EJ: **The use of the Bayley Scales of Infant and Toddler Development III with clinical populations: a preliminary exploration**. *Phys Occup Ther Pediatr* 2012, **32**(1):24-33.