



**ΠΑΝΕΠΙΣΤΗΜΙΟ ΣΤΕΡΕΑΣ ΕΛΛΑΔΟΣ  
ΤΜΗΜΑ ΠΛΗΡΟΦΟΡΙΚΗΣ  
ΜΕ ΕΦΑΡΜΟΓΕΣ ΣΤΗ ΒΙΟΪΑΤΡΙΚΗ**

**Μεθοδολογίες μελέτης απλότυπων σε μελέτες γενετικής  
συσχέτισης**

**Λούκας Αλέξιος**

**ΠΤΥΧΙΑΚΗ ΕΡΓΑΣΙΑ**

**Επιβλέποντες**

**Μπάγκος Παντελής, Επίκουρος Καθηγητής**

**Ιωαννίδης Αναστάσιος, Επισκέπτης Επίκουρος Καθηγητής**

**Λαμία, 2010**



Εξεταστική επιτροπή

1. Μπάγκος Παντελής (επιβλέπων)
2. Αδάμ Μαρία
3. Ιωαννίδης Αναστάσιος (επιβλέπων)



## ΕΥΧΑΡΙΣΤΙΕΣ

Η εργασία αυτή έλαβε χώρα στην Λαμία, στο Πανεπιστήμιο Στερεάς Ελλάδος, στο Τμήμα Πληροφορικής με εφαρμογές στην βιοϊατρική. Η ολοκλήρωσή της οφείλεται όχι μόνο στον συγγραφέα, αλλά και στους ανθρώπους που ακολουθούν.

Θα ήθελα να ευχαριστήσω τους επιβλέποντες καθηγητές μου κ. Μπάγκο Παντελή και κ. Ιωαννίδη Αναστάσιο που με επέλεξαν για την εκπόνηση της εργασίας αυτής, καθώς και για τις συμβουλές τους πάνω στο θέμα και το κείμενό της. Επίσης θα ήθελα να ευχαριστήσω τους φίλους, τους συμφοιτητές, την οικογένειά μου, τους γονείς μου αλλά και τον Σεβασμιώτατο Μητροπολίτη Ηλείας και Ωλένης Γερμανό Παρασκευόπουλο για την υποστήριξη που μου παρείχαν κατά την διάρκεια των σπουδών μου.



# ΠΕΡΙΕΧΟΜΕΝΑ

Εξεταστική επιτροπή .....	3
ΕΥΧΑΡΙΣΤΙΕΣ .....	5
ΠΕΡΙΕΧΟΜΕΝΑ .....	7
ΠΕΡΙΛΗΨΗ .....	9
ABSTRACT .....	11
<b>ΚΕΦΑΛΑΙΟ 1 - ΕΙΣΑΓΩΓΗ.....</b>	<b>13</b>
1.1 Εισαγωγή στους απλότυπους .....	13
1.2 Το International HarMap project.....	13
1.3 Μονονουκλεοτιδικός πολυμορφισμός .....	18
1.4 Το πρόβλημα της εύρεσης του απλότυπου.....	20
1.4.1 Τα γενετικά μοντέλα .....	23
<b>ΚΕΦΑΛΑΙΟ 2 – ΥΛΙΚΑ ΚΑΙ ΜΕΘΟΔΟΙ .....</b>	<b>27</b>
2.1 Ορισμός του προβλήματος .....	27
2.2 Η εφαρμογή Stata.....	27
2.3 Εκτίμηση του Linkage disequilibrium.....	29
2.3.1 Είδη και μεθοδολογίες υπολογισμού του LD.....	29
2.4 Μελέτες με ασθενείς και μάρτυρες.....	36
2.4.1 Ο έλεγχος $\chi^2$ του Pearson για την ανεξαρτησία δυο δειγμάτων.....	40
2.4.2 Εκτίμηση μέγιστης πιθανοφάνειας.....	42
2.4.3 Λογιστική παλινδρόμηση .....	43
2.4.4 Πολυωνομική λογιστική παλινδρόμηση .....	52
2.4.5 Παλινδρόμηση Poisson .....	58
2.4.6 Η εφαρμογή Harlview .....	67
2.5 Μελέτες με μεταβλητή συνεχούς τιμής.....	73
2.5.1 Απλή γραμμική παλινδρόμηση .....	73
<b>ΚΕΦΑΛΑΙΟ 3 – ΑΠΟΤΕΛΕΣΜΑΤΑ ΚΑΙ ΣΥΖΗΤΗΣΗ.....</b>	<b>79</b>
3.1 ΑΠΟΤΕΛΕΣΜΑΤΑ.....	79
3.1.1 Μελέτες απλότυπων με τις τιμές της HDL χοληστερίνης (συνεχής μεταβλητή).....	79
3.1.2 Μελέτες με ασθενείς και μάρτυρες (case-control) .....	83
3.2 ΣΥΖΗΤΗΣΗ.....	95
Παράρτημα 1 – απόρριψη της μηδενικής υπόθεσης.....	99
Παράρτημα 2 – αποδοχή της μηδενικής υπόθεσης.....	111
<b>ΒΙΒΛΙΟΓΡΑΦΙΑ .....</b>	<b>139</b>





## ΠΕΡΙΛΗΨΗ

Η απλοτυπική ανάλυση αποτελεί πλέον αναπόσπαστο κομμάτι της γονιδιωματικής ανάλυσης. Ο απλότυπος είναι ένα σύνολο πολυμορφισμών που παρουσιάζονται στο DNA και είναι στατιστικά συνδεδεμένοι μεταξύ τους λόγω γειτνίασης. Από τις κλασικές μεθόδους αλληλούχισης, η γνώση του γονότυπου δεν μας δίνει μονοσήμαντα την πληροφορία για τους απλότυπους και για αυτό μας χρειάζονται αλγόριθμοι στατιστικής φύσεως (το πρόβλημα της φάσης των απλότυπων). Στις μελέτες γενετικής συσχέτισης χρησιμοποιούνται ολοένα και περισσότερο οι απλότυποι για την διερεύνηση της συσχέτισης κοινών πολυμορφισμών με ασθένειες. Η δημιουργία του «χάρτη» του ανθρώπινου γονιδιώματος (HarMap) ανοίγει νέες πόρτες με την μεγάλη αυτή βάση δεδομένων και ενθαρρύνει την επιστημονική κοινότητα για περαιτέρω ενασχόληση. Σε αυτήν εδώ την εργασία συνοψίζονται οι μέθοδοι που χρησιμοποιούνται σε μελέτες γενετικής συσχέτισης με απλότυπους, αναλύονται τα μαθηματικά μοντέλα για αυτές τις μεθόδους και παρουσιάζονται αναλυτικά παραδείγματα με την χρήση των λογισμικών Stata και HarMapView. Η εφαρμογή των μεθόδων γίνεται σε πραγματικά δεδομένα τα οποία πήραμε από δημοσιευμένες μελέτες.

Λέξεις-Κλειδιά: Απλότυπος, Συσχέτιση, Μονοσημειακός Πολυμορφισμός, HarMap, Ασθενείς-Μάρτυρες, Ανάλυση



## **ABSTRACT**

Haplotype analysis constitutes an integral part of genomic analysis. A Haplotype is referred to as a set of single nucleotide polymorphisms that we see in certain DNA regions, and they are statistically associated. From the classic sequencing methods, the knowledge alone of the genotype, doesn't give us all the information we need about haplotypes, that is why we need to use statistical inference algorithms (haplotype inference problem). In genetic association studies more and more we see haplotypes used for association between single nucleotide polymorphisms and diseases. The creation of the human genome "map" known as HapMap opens new doors with this large database and encourages the scientific community to further activity. In this thesis we summarize the methods that are used in genetical association with haplotypes and we analyze the mathematical models that these methods use. We present analytical examples with the use of Stata and Haploview software. These methods are applied in real data, taken from published studies.

Key-Words: Haplotype, Association, Single Nucleotide Polymorphism, HapMap, Case-Control, Analysis

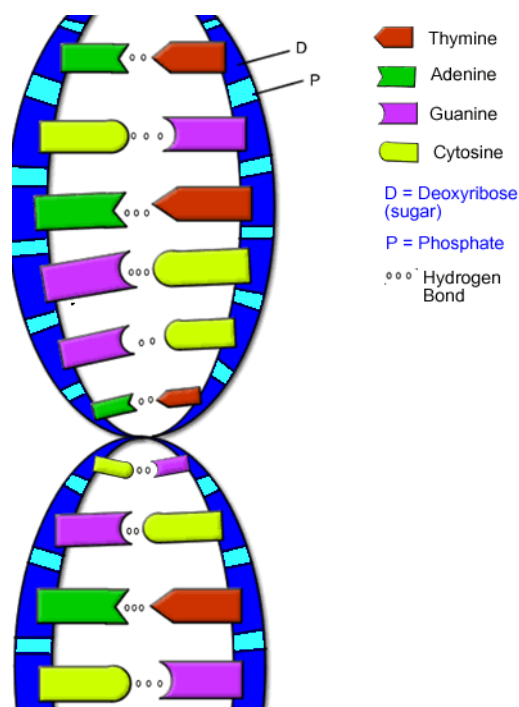


# ΚΕΦΑΛΑΙΟ 1 - ΕΙΣΑΓΩΓΗ

## 1.1 Εισαγωγή στους απλότυπους

Η λέξη απλότυπος προέρχεται από τα παράγωγα «απλοειδής» και «γονότυπος». Ουσιαστικά είναι ένα σύνολο πολυμορφισμών σε διάφορα σημεία του DNA (εικόνα 1), οι οποίοι είναι στατιστικά συνδεδεμένοι μεταξύ τους, δηλαδή η εμφάνιση του ενός κάνει πολύ πιθανή και την εμφάνιση άλλων. Μεταδίδονται μαζί σαν ένα ενιαίο κομμάτι πάνω στο γενετικό υλικό (HapMap 2005).

Έχοντας στην κατοχή την γνώση κάποιων απλότυπων για το γονιδίωμα του ανθρώπου, ίσως έχουμε στα χέρια μας πολύτιμες πληροφορίες που θα μας βοηθήσουν να διερευνήσουμε σε μεγαλύτερο βαθμό την γενετική συμπεριφορά κάποιων ασθενειών. Ήδη γίνεται μια παγκόσμια προσπάθεια για την καταγραφή ενός «χάρτη» απλότυπων του ανθρώπινου γονιδιώματος, ο οποίος θα περιγράφει κάποια κοινά πρότυπα στην γενετική διαφοροποίηση μεταξύ των ανθρώπων.



Εικόνα 1: Η δομή του DNA

## 1.2 Το International HapMap project

Ο χάρτης που αναφέρεται παραπάνω είναι ελεύθερα προσβάσιμος στους ερευνητές, ώστε να τους βοηθήσει να συσχετίσουν εκείνες τις γενετικές πληροφορίες που επηρεάζουν την υγεία, συμβάλουν στις ασθένειες, στην αλληλεπίδραση των οργανισμών με τα φάρμακα, αλλά και το ρόλο που παίζουν κάποιοι περιβαλλοντικοί παράγοντες. Μια τόσο μεγάλη μελέτη βέβαια διέπεται από κάποιους κανόνες μεθοδολογίας και ηθικής που δημοσιεύτηκαν από την ομάδα των ερευνητών (HapMap 2004).

Η έρευνα σήμερα στρέφεται κυρίως στους λεγόμενους «μονοσημειακούς πολυμορφισμούς» (Single nucleotide polymorphism ή SNP), οι οποίοι δηλώνουν μια σημειακή αλλαγή στο γενετικό υλικό. Υπάρχουν περίπου δέκα εκατομμύρια πολυμορφισμοί που παρατηρούνται στο ανθρώπινο γονιδίωμα. Από αυτούς, μια πολύ μικρή ομάδα είναι άμεσα συνδεδεμένη με ασθένειες, παρενέργειες κλπ.

Σκοπός αυτού του «χάρτη» είναι πρώτον η διεύρυνση αυτής της μικρής ομάδας, και δεύτερον να κάνει προσβάσιμη την γενετική πληροφορία στους ερευνητές, οι οποίοι δεν θα χρειάζεται πλέον να ψάχνουν για πολυμορφισμούς (HarMap 2003).

Σημαντική συμμετοχή σε αυτήν την προσπάθεια υπάρχει και από έναν Έλληνα επιστήμονα, τον Δρ. Πάνο Δελούκα ο οποίος μαζί με συναδέλφους του στο Sanger Center της Μ. Βρετανίας υλοποίησαν το HarMap project (εικόνα 2) (<http://www.sanger.ac.uk/research/faculty/pdeloukas/>). Όπως εξηγεί ο κύριος Δελούκας με δηλώσεις του «το διεθνές πρόγραμμα HarMap δεν είναι παρά ένα εργαλείο το οποίο θα βοηθήσει στην αξιοποίηση της πληροφορίας που έχει προκύψει από την ανάγνωση του ανθρώπινου γονιδιώματος. Στην πράξη, θα δίνει το στίγμα κάποιων πολυμορφισμών με ιδιαίτερο ενδιαφέρον πάνω στο γονιδίωμα».



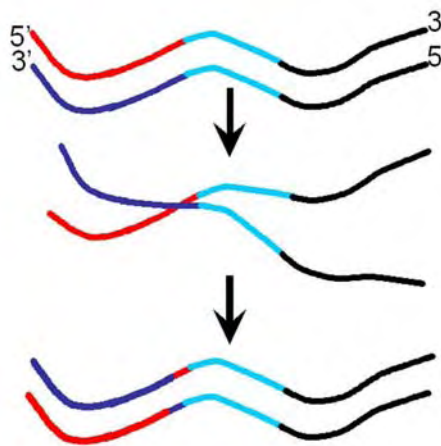
**Εικόνα 2: Ο Δρ. Πάνος Δελούκας**

Αυτός ο «χάρτης» θα χρησιμοποιεί τους πολυμορφισμούς ως σημεία αναφοράς πάνω στο γονιδίωμα, και ενδέχεται να επιταχύνει τον εντοπισμό γονιδίων που σχετίζονται με ασθένειες και ειδικότερα με ασθένειες όπου η περιβαλλοντική συνιστώσα παίζει καθοριστικό ρόλο στην έκφρασή τους. «Τα γονίδια των

μονογονιδιακών ασθενειών, αυτών δηλαδή που η αιτιολογία τους οφείλεται σε μεταλλάξεις ενός και μόνο γονιδίου, δεν αποτελούν αντικείμενο των δικών μας αναζητήσεων. Οι επιστήμονες έχουν αποτελεσματικούς τρόπους να τα εντοπίζουν. Αυτό που μας ενδιαφέρει είναι γονίδια τα οποία κάνουν κάποιους από μας πιο ευαίσθητους σε λοιμώξεις ή αυξάνουν τον κίνδυνο που διατρέχουμε να αναπτύξουμε διαβήτη ή άσθμα ή υπέρταση αν εκτεθούμε σε συγκεκριμένες συνθήκες. Σε τέτοιου είδους αναζητήσεις θα είναι πολύτιμο το HarMap. Για παράδειγμα, αν υποθεθεί ότι θα θέλαμε να αναζητήσουμε τη γενετική διαφορά που καθιστά κάποιους ανθρώπους πιο ευαίσθητους από άλλους στην εμφάνιση διαβήτη. Χωρίς το HarMap θα έπρεπε να πάρουμε δύο ομάδες ανθρώπων, ας πούμε 1.000 ασθενείς με διαβήτη και 1.000 υγιείς, και να εξετάσουμε ποιος από τους 10 εκατομμύρια SNP's είναι διαφορετικός στην ομάδα των ασθενών. Αντιλαμβάνεστε πόσο χρονοβόρα αλλά και ακριβή θα ήταν μια τέτοια διαδικασία, ειδικά καθώς θα έπρεπε να επαναληφθεί για όλες τις ασθένειες που μας ενδιαφέρουν...» μας λέει ο κ. Δελούκας (Deloukas and Bentley 2004).

Σημαντικό στοιχείο στο θέμα του χάρτη είναι η γνωστή διαδικασία του ανασυνδυασμού του DNA (εικόνα 3), που ως γνωστό λαμβάνει χώρα κατά την διάρκεια της μείωσης των γαμετικών κυττάρων. Αυτό έχει ως αποτέλεσμα το γεγονός ότι ο καθένας μας έχει στα γαμετικά του κύτταρα ένα νέο σύνολο γενετικού υλικού, και όχι απλά ένα μείγμα από το DNA του πατέρα και της μητέρας μας. Ο ανασυνδυασμός έχει ιδιαιτερότητες οι οποίες αξιοποιούνται για τη δημιουργία του HarMap (HarMap 2003).

Σύμφωνα με τον κ. Δελούκα «τα σημεία πάνω στα χρωμοσώματα που συμβαίνει ο ανασυνδυασμός δεν είναι τυχαία. Επίσης οι άνθρωποι είμαστε σχετικά νέο είδος, πράγμα που σημαίνει ότι δεν έχουν γίνει τόσο πολλοί ανασυνδυασμοί κατά τη διάρκεια των αιώνων. Αυτό που αναζητούμε και χαρτογραφούμε εμείς είναι κομμάτια DNA τα οποία έχουν περάσει αυτούσια. Πάνω σε αυτά τα κομμάτια υπάρχουν SNP's οι οποίοι, όταν ολοκληρωθεί το HarMap, θα λειτουργούν ως δείκτες για να καθοδηγούν τους ερευνητές στα σημεία του γονιδιώματος που πρέπει να ψάξουν προκειμένου να βρουν αυτό που τους ενδιαφέρει».



**Εικόνα 3: Συνοπτικά ο ανασυνδιασμός του DNA**

Πρακτικά, αυτό που θα κάνουν οι ερευνητές όταν τελειώσει η δημιουργία του HarMap, θα είναι να αναζητούν αν οι άνθρωποι που φέρουν κάποια νόσο, φέρουν επίσης και κάποιους χαρακτηριστικούς πολυμορφισμούς. Αυτοί οι πολυμορφισμοί θα τους δείξουν σε ποιο σημείο πάνω στα χρωμοσώματα πρέπει να αναζητήσουν το γονίδιο που σχετίζεται με αυτή τη νόσο.

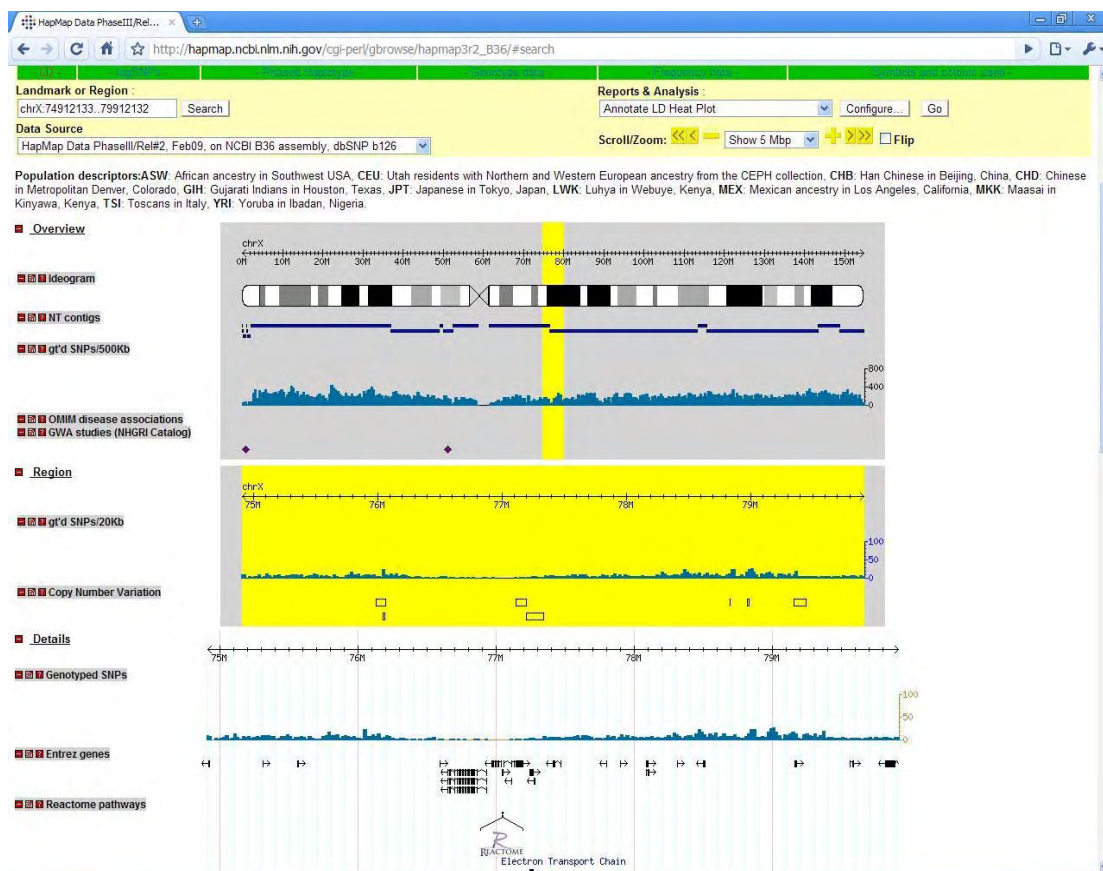
Ο οργανισμός ο οποίος ξεκίνησε την κατασκευή αυτού του χάρτη ονομάζεται «International HarMap Project» (HarMap 2003) και ουσιαστικά είναι μια συνεργασία μεταξύ διαφόρων ερευνητών ακαδημαϊκών κέντρων, μη-κερδοσκοπικών βιοϊατρικών οργανώσεων και ιδιωτικών εταιριών σε Καναδά, Κίνα, Ιαπωνία, Νιγηρία, Ηνωμένο βασίλειο και Ηνωμένες πολιτείες Αμερικής (HarMap 2005). Το λογότυπο του οργανισμού φαίνεται στην εικόνα 6.

Το project ξεκίνησε τον Οκτώβριο του 2002 και αρχικά εκτιμήθηκε ότι θα κρατούσε τρία χρόνια, ωστόσο βλέπουμε ότι συνεχίζεται ακόμη καθώς τον Μάιο του 2010 είχαμε ακόμη μια δημοσίευση γενετικών δεδομένων (HarMap3 Public Release #3). Οι ενδιαφερόμενοι ερευνητές μπορούν να βρίσκουν τα αποτελέσματα της χαρτογράφησης αυτής στον δικτυακό τόπο <http://www.hapmap.org/>.

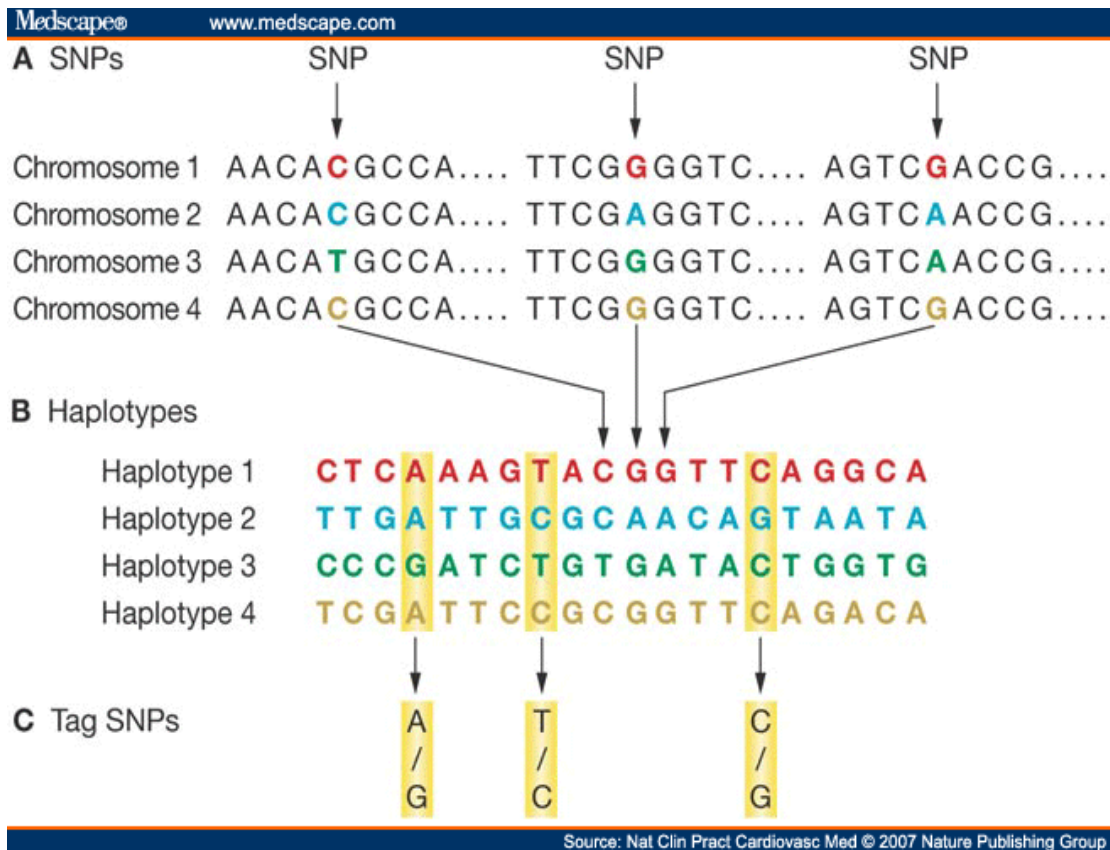
Η ιστοσελίδα επιτρέπει την περιήγηση σε όλες τις περιοχές του ανθρώπινου γονιδιώματος, προσφέροντας όλα τα καταγεγραμμένα στοιχεία όπως γνωστοί απλότυποι, πολυμορφισμοί, γονίδια, εσώνια-εξώνια κλπ. Σαν παράδειγμα στην εικόνα 4 που ακολουθεί, έχουμε επιλέξει μια περιοχή του χρωμοσώματος X προς επεξεργασία μέσω της εφαρμογής που παρέχεται από την ιστοσελίδα του HarMap. Η εφαρμογή μας δείχνει όλα τα δεδομένα που έχει καταχωρημένα και χωρίζει την



περιοχή του γονιδιώματος που ερευνούμε σε τρία επίπεδα μεγέθυνσης. Φυσικά μπορούμε να μεγεθύνουμε επιπλέον την εικόνα και να δούμε από πιο κοντά τις βάσεις. Όσο μεγεθύνουμε και μπαίνουμε σε λεπτομέρειες, τόσο περισσότερο φαίνονται οι καταχωρημένες πληροφορίες των πολυμορφισμών. Ο χρήστης μπορεί να διαλέξει και να εξειδικεύσει την αναζήτησή του με μια πληθώρα επιλογών, σε μια μεγάλη βάση βιολογικών δεδομένων, η οποία είναι εξαιρετικά οργανωμένη. Το μεγαλύτερο πλεονέκτημα είναι ότι είναι εντελώς ελεύθερη η προσπέλαση σε αυτήν, καθώς επίσης και ότι βρίσκεται φορτωμένη στο διαδίκτυο, πράγμα που την κάνει εύκολα προσβάσιμη από οποιοδήποτε μέρος.



Εικόνα 4: Το χρωμόσωμα X στον χάρτη του HapMap



Εικόνα 5: Πως δημιουργούνται οι απλότυποι από τους πολυμορφισμούς

### 1.3 Μονονουκλεοτιδικός πολυμορφισμός

Παρόλο που δυο οποιοδήποτε άνθρωποι έχουν κοινό το 99.5% του γονιδιώματός τους (HarMap 2003), κάποιος σε μια συγκεκριμένη θέση στο γενετικό τους υλικό μπορεί να έχουν την βάση A (αδενίνη) ενώ κάποιος άλλος την G (γουανίνη). Μια τέτοια αλλαγή (A->G) είναι γνωστή με το όνομα μονονουκλεοτιδικός πολυμορφισμός (single-nucleotide polymorphism ή SNP) και μπορεί να συμβεί σε διάφορα μέρη του γονιδιώματος κατά την διάρκεια ζωής ενός οργανισμού. Κάθε περίπτωση (A ή G) ονομάζεται αλληλόμορφο. Μια απεικόνιση των αλληλόμορφων φαίνεται στην εικόνα 7, όπου η πλειοψηφία των ατόμων έχουν την βάση C (κυτοσίνη) ενώ εμφανίζεται σπάνια η T (θυμίνη).

Το HarMap ασχολείται μόνο με κοινούς πολυμορφισμούς, δηλαδή με αυτούς όπου κάθε αλληλόμορφο υπάρχει τουλάχιστον στο 1% του πληθυσμού. Κάθε άνθρωπος έχει δυο αντίγραφα κάθε χρωμοσώματος, εκτός από τα χρωμοσώματα που καθορίζουν το φύλο. Για κάθε SNP, ο συνδυασμός των αλληλόμορφων που έχει ένα άτομο, ονομάζεται γονότυπος. Έτσι υπάρχει η διαδικασία εύρεσης του γονότυπου

ενός ατόμου, σε μια συγκεκριμένη περιοχή του γονιδιώματός του. Μια τέτοια διαδικασία γίνεται με δυο τρόπους. Σε εργαστήρια με την μέθοδο της ηλεκτροφόρησης, αλλά και με στατιστικές μεθόδους όπου χρησιμοποιούμε βέβαια την ήδη υπάρχουσα πληροφορία που έχουμε από το εργαστήριο. Σαν αποτέλεσμα



**Εικόνα 6: Το HapMap project και οι χώρες που το αποτελούν**

βρίσκουμε όλους τους πολυμορφισμούς και τα αλληλόμορφα τους για μια συγκεκριμένη περιοχή του DNA. Οι πολυμορφισμοί που βρίσκονται σε κοντινή θέση μεταξύ τους και πάνω στο ίδιο χρωμόσωμα, είναι συσχετισμένοι. Αυτό σημαίνει πως αν το αλληλόμορφο από έναν πολυμορφισμό είναι γνωστό, τότε τα αλληλόμορφα των διπλανών πολυμορφισμών, ίσως να μπορούμε να τα προβλέψουμε (HapMap 2003).

Αυτό συμβαίνει επειδή κάθε πολυμορφισμός, εμφανίστηκε εξελικτικά ως μια σημειακή μετάλλαξη, η οποία στη συνέχεια πέρασε στους απογόνους περιβαλλόμενη από άλλες, πιο παλιές σημειακές μεταλλάξεις. Συνήθως βέβαια, πολυμορφισμοί που χωρίζονται από μεγάλη απόσταση, δεν εμφανίζουν συσχέτιση. Κάτι τέτοιο συμβαίνει λόγω του ότι χωρίζονται με την διαδικασία ανασυνδιασμού του DNA σε κάθε γενιά, με αποτέλεσμα το ανακάτεμα των ακολουθιών των αλληλόμορφων των δυο χρωμοσωμάτων. Ωστόσο η πληροφορία που ψάχνουμε είναι μια ακολουθία αλληλόμορφων πάνω σε ένα συγκεκριμένο χρωμόσωμα, η οποία δεν έχει διαχωριστεί με τον ανασυνδιασμό του DNA. Μια τέτοια ακολουθία ονομάζεται απλότυπος. Επίσης απλότυπος θεωρείται και ένα σύνολο από μονοσημειακούς πολυμορφισμούς, όχι απαραίτητα σε ακολουθία, αλλά σε μικρή απόσταση μεταξύ τους σε μια περιοχή του DNA (HapMap 2003).

Για να βρούμε τους γενετικούς παράγοντες που εμπλέκονται σε μια ασθένεια, μπορούμε να ακολουθήσουμε την παρακάτω διαδικασία. Πρώτα πρέπει να βρεθεί μια περιοχή ενδιαφέροντος πάνω στο γονιδίωμα (π.χ. περιοχή γονιδίου), πιθανότατα γνωστή από προηγούμενες μελέτες. Έπειτα εντοπίζουμε σε αυτή την περιοχή ένα

σύνολο από πολυμορφισμούς οι οποίοι έχουν ισχυρή συσχέτιση μεταξύ τους. Το σύνολο των αλληλόμορφων από αυτούς τους πολυμορφισμούς θα μας δώσει τον απλότυπο αυτού του ατόμου. Η διαδικασία αυτή συνοψίζεται στην εικόνα 5. Στη συνέχεια βρίσκουμε τον απλότυπο αυτών των πολυμορφισμών σε ένα σύνολο ατόμων, κάποιων που να έχουν την ασθένεια, και κάποιων που να μην την έχουν. Με διάφορες συγκρίσεις μεταξύ των δυο αυτών ομάδων (ασθενείς – μη ασθενείς) μπορούμε να καταλάβουμε ποιες περιοχές και ποιοι απλότυποι παίζουν κάποιο ρόλο στην εμφάνιση και την συμπεριφορά της ασθένειας (Fallin, Cohen et al. 2001).



Εικόνα 7: Μια διαφορετική απεικόνιση του μονοσημειακού πολυμορφισμού

#### 1.4 Το πρόβλημα της εύρεσης του απλότυπου

Αντίθετα από τις σπάνιες κληρονομικές ασθένειες, συνδυασμοί διαφόρων γονιδίων αλλά και το περιβάλλον, παίζουν κάποιο ρόλο στην εμφάνιση πιο κοινών παθήσεων (όπως διαβήτης, καρκίνος, παθήσεις της καρδιάς, εγκεφαλικό, κατάθλιψη και άσθμα) αλλά και παρενεργειών στην λήψη κάποιων φαρμάκων (HarMap 2003).

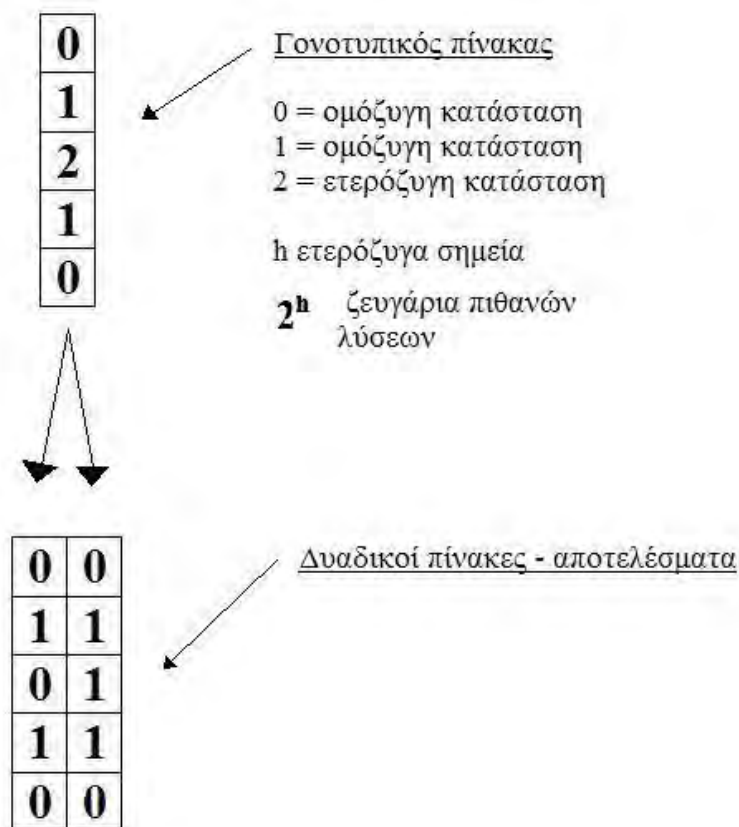
Για να μπορέσουμε να βρούμε τους γενετικούς παράγοντες που παίζουν ρόλο στην εμφάνιση των παραπάνω στοιχείων, θα πρέπει να είμαστε σε θέση να αποκτήσουμε στα χέρια μας την πλήρη γενετική ακολουθία κάποιων ατόμων. Θα χρειαστούμε κάποια άτομα που να έχουν την ασθένεια, αλλά και κάποια υγιή. Κατόπιν θα πρέπει να συγκρίνουμε τις ακολουθίες αυτών που έχουν την ασθένεια, και αυτών που είναι υγιείς, για διαφορές στο γονιδίωμά τους. Μία τέτοια μέθοδος είναι δύσκολο να επιτευχθεί διότι το κόστος του να αποκτήσουμε την πλήρη ακολουθία του γονιδιώματος ενός ατόμου είναι μεγάλο (HarMap 2003). Το HarMap

προτείνει μια πιο σύντομη μέθοδο η οποία περιγράφεται παρακάτω ως «το πρόβλημα της εύρεσης του απλότυπου» (Lin, Cutler et al. 2002).

Ένα βασικό πρόβλημα των απλότυπων, είναι ότι αν έχουμε τον γονότυπο ενός ατόμου, είναι δύσκολο να καταλήξουμε σε συμπέρασμα για το ποιος θα είναι ο απλότυπός του. Στους διπλοειδείς οργανισμούς, όπου ανήκει και ο άνθρωπος, υπάρχουν δυο αντίγραφα κάθε χρωμοσώματος και κατά συνέπεια και δυο περιοχές ενδιαφέροντος. Ο απλότυπος ουσιαστικά είναι τα στοιχεία της περιοχής ενδιαφέροντος του ενός μόνο αντίγραφου, ενώ ο γονότυπος περιλαμβάνει τα στοιχεία και των δυο χρωμοσωμάτων. Στις ασθένειες οι οποίες επηρεάζονται από παραπάνω από ένα γονίδιο, είναι συνήθως πιο χρήσιμη η πληροφορία που μας δίνει ο απλότυπος συγκριτικά με τον γονότυπο. Η πληροφορία από την οποία αποτελείται ο απλότυπος, όπως αναφέραμε παραπάνω, είναι είτε η πλήρης DNA ακολουθία της περιοχής ενδιαφέροντος, είτε το σύνολο των μονοσημειακών πολυμορφισμών σε εκείνη την περιοχή. Η δεύτερη εκδοχή είναι αυτή η οποία κυριαρχεί, αλλά και η πιο χρήσιμη, επειδή περιέχει ακριβώς την πληροφορία που πρέπει να αναλυθεί χωρίς περιττά τμήματα.

Το πρόβλημα αυτό της εύρεσης του απλότυπου μπορεί να αντιμετωπιστεί και ως υπολογιστικό πρόβλημα (Gusfield 2000), με την δημιουργία γονοτυπικών πινάκων, πινάκων δηλαδή που θα περιέχουν πληροφορία για περιοχές ενδιαφέροντος του γενετικού υλικού. Στις θέσεις των πινάκων αυτών θα υπάρχουν οι αριθμοί 0, 1 και 2. Για θέσεις του πίνακα όπου η αντίστοιχος περιοχές στα ομόλογα χρωμοσώματα είναι σε ομόζυγη κατάσταση θα αντιστοιχούν οι αριθμοί 0 ή 1 (κάθε SNP μπορεί να έχει δυο πιθανά αλληλόμορφα), ενώ για ετερόζυγη κατάσταση το 2.





**Εικόνα 8: Το υπολογιστικό πρόβλημα της εύρεσης του απλότυπου**

Για ένα σύνολο  $n$  πινάκων θα έχουμε ως αποτέλεσμα ένα σύνολο από  $n$  δυαδικούς πίνακες, όπου κάθε δυαδικός πίνακας αντιστοιχεί σε έναν από τους αρχικούς. Αυτός ο δυαδικός πίνακας θα πρέπει να έχει 0 ή 1 και στις δυο θέσεις όπου ο αρχικός πίνακας είχε 0 ή 1. Επίσης για τον πίνακα που είχε το 2, θα πρέπει ο δυαδικός να έχει 0 στην πρώτη θέση και 1 στην δεύτερη ή αντίστροφα. Αυτοί οι δυαδικοί πίνακες μπορούν να γραφούν με διάφορους τρόπους, και δηλώνουν έτσι τις πιθανές μορφές των απλότυπων που μπορούν να προκύψουν. Για ένα άτομο με  $h$  ετερόζυγα σημεία προκύπτουν  $2^{h-1}$  ζευγάρια που θα μπορούσαν να είναι η λύση του προβλήματος της εύρεσης του απλότυπου (εικόνα 8). Ωστόσο χωρίς επιπλέον επεξεργασία των δυαδικών πινάκων δεν μπορούμε να βρούμε ποιές από τις πιθανές λύσεις είναι οι σωστές. Συνεπώς υπάρχει η ανάγκη χρήσης κάποιων γενετικών μοντέλων (Stephens and Donnelly 2003) για να μας καθοδηγήσουν προς την λύση.

### 1.4.1 Τα γενετικά μοντέλα

Η εύρεση ενός απλότυπου μπορεί να γίνει και πειραματικά σε βιολογικά εργαστήρια (Clark 1990), αλλά αυτή η μέθοδος είναι αρκετά χρονοβόρα, απαιτεί πολύ εργασία και κοστίζει πολλά χρήματα. Για αυτό το λόγο χρησιμοποιούνται εναλλακτικά διαφόρων ειδών στατιστικές μέθοδοι, οι οποίες μπορούν αν χρησιμοποιηθούν σωστά να είναι σύντομες και έγκυρες, χωρίς μεγάλες χρηματικές απαιτήσεις (Lin and Fann 2009). Αυτές οι στατιστικές μέθοδοι μπορούν να χωριστούν σε τρεις κύριες κατηγορίες, τον αλγόριθμο του Clark, τον αλγόριθμο EM (expectation-maximization) και τις προσεγγίσεις του Bayes (Niu 2004). Επιπλέον έχουν προταθεί και άλλοι αλγόριθμοι πέρα από τους κλασικούς (Gusfield 2000).

Ο αλγόριθμος του Clark (Clark 1990), αρχίζει με μια λίστα γνωστών απλότυπων από τα ευκρινή άτομα (άτομα δηλαδή που οι γονότυποι τους δεν έχουν ή έχουν μόνο μια ετερόζυγη περιοχή). Στη συνέχεια η μέθοδος επιχειρεί να αναλύσει και τους γονότυπους από τα εναπομείναντα άτομα χρησιμοποιώντας τους απλότυπους που βρήκε αρχικά και το συμπλήρωμά τους, δηλαδή με χρήση του εναλλακτικού τους αλληλόμορφου σε κάθε θέση, καθώς και τις συχνότητες τους. Κάποιες φορές ο αλγόριθμος αποτυγχάνει να αναγνωρίσει τους απλότυπους για όλα τα άτομα, λόγω του ότι αν βρεθεί σε λάθος μορφή κάποιος απλότυπος, αυτό θα οδηγήσει σε σφάλμα στα επόμενα βήματα. Παρ όλα αυτά ο

αλγόριθμος του Clark είναι η πρώτη στατιστική μέθοδος για την απόκτηση πληροφορίας απλότυπων και έτσι θεμελιωδώς σημαντικός. Σαν είσοδο ο αλγόριθμος δέχεται  $n$  γονοτυπικούς πίνακες μήκους  $m$ . Όπως είπαμε παραπάνω αυτοί οι πίνακες περιέχουν τους αριθμούς 0, 1 ή 2. Οι θέσεις των πινάκων που έχουν τους αριθμούς 0 ή 1 θεωρούνται «αποφασισμένες» ενώ αυτές που έχουν 2 θεωρούνται «διφορούμενες». Ένας πίνακας που δεν έχει «διφορούμενες» θέσεις ονομάζεται και αυτός «αποφασισμένος», ενώ ένας πίνακας με «διφορούμενες» θέσεις ονομάζεται «διφορούμενος». Για δυο



**Εικόνα 9: Το λογότυπο της εφαρμογής Hap**

διαφορετικούς πίνακες A και B η συγχώνευση τους δημιουργεί τον πίνακα AB ο οποίος χαρακτηρίζεται «διφορούμενος» και οι θέσεις όπου και ο A και ο B είχαν 0 ή 1, γίνονται 0 ή 1 αντίστοιχα. Οι θέσεις του AB όπου οι A και B είχαν (1 και 0) ή (0 και 1) αντίστοιχα παίρνουν την τιμή 2.

Η μέθοδος του Clark αρχίζει με την ταυτοποίηση κάθε πίνακα με καμία ή μια «διφορούμενη» θέση, μιας και στην περίπτωση που δεν υπάρχει καμία «διφορούμενη» θέση ο απλότυπος που προκύπτει είναι ο ίδιος ο γονοτυπικός πίνακας απaráλλακτος. Στην περίπτωση που έχουμε μια «διφορούμενη» θέση τότε προκύπτουν δυο απλοτυπικοί πίνακες. Οι απλότυποι που προκύπτουν ονομάζονται «αρχικά αποφασισμένοι πίνακες».

Το κύριο κομμάτι του αλγορίθμου του Clark επεκτείνεται με βάση τους «αρχικά αποφασισμένους πίνακες» και αποφασίζει για τους απομένοντες «διφορούμενους» πίνακες. Ο Clark πρότεινε τον ακόλουθο κανόνα που συνάγει τον επόμενο «αποφασισμένο» πίνακα, από έναν γονοτυπικό πίνακα A και από έναν ήδη «αποφασισμένο» πίνακα B. Ο «αποφασισμένος» πίνακας B μπορεί να είναι είτε από τα δεδομένα εισόδου του αλγορίθμου, είτε από την προηγούμενη εφαρμογή αυτού του κανόνα.

Ο κανόνας αυτός ονομάζεται «κανόνας απόφασης» και είναι ο εξής: Αν υποθέσουμε ότι ο πίνακας A είναι «διφορούμενος» με h διφορούμενες θέσεις, και R ένας «αποφασισμένος» πίνακας ο οποίος αντιστοιχεί σε μια από τις  $2^h$  πιθανές αναλύσεις του πίνακα A, τότε ο A είναι η συγχώνευση του ενός αντιγράφου του «αποφασισμένου» πίνακα R και ενός άλλου «αποφασισμένου» πίνακα NR (ο οποίος καθορίζεται με συγκεκριμένο τρόπο). Όλες οι αποφασισμένες θέσεις του A αντιγράφονται στον πίνακα NR αντίθετα από ότι σε σχέση με τον R. Έτσι όταν ο πίνακας NR χαρακτηριστεί «αποφασισμένος», προστίθεται στο πλήθος των αποφασισμένων πινάκων, και ο πίνακας A διαγράφεται από το σύνολο των πινάκων.

Για παράδειγμα αν ο A είναι 0212 και ο R είναι 0110, τότε ο NR είναι 0011. Η ερμηνεία αυτού είναι ότι αν οι δυο απλότυποι σε ένα άτομο είναι 0110 και 0011 τότε ο γονότυπος που παρατηρείται θα είναι 0212. Ο κανόνας απόφασης βγάζει συμπέρασμα για τον πίνακα 0212 με το σκεπτικό ότι ο 0110 είναι ένας απλότυπος στον πληθυσμό, για να αποφασίσει ότι ο 0011 είναι επίσης ένας απλότυπος του πληθυσμού. Όταν ο κανόνας απόφασης μπορεί να εφαρμοστεί για την εξαγωγή συμπεράσματος στον πίνακα NR από τους πίνακες A και R, τότε μπορούμε να πούμε ότι ο R μπορεί να εφαρμοστεί στον A. Είναι εύκολο να αποφανθούμε αν ένας



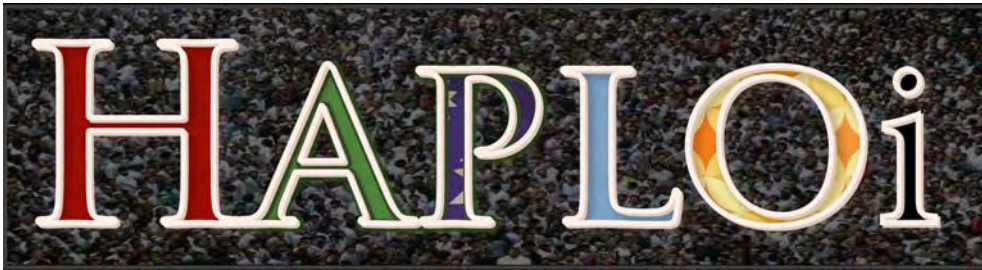
«αποφασισμένος» πίνακας R μπορεί να εφαρμοστεί σε έναν «διφορούμενο» πίνακα A: Ο R μπορεί να εφαρμοστεί στον A αν και μόνο αν ο A δεν περιέχει καμία θέση s, τέτοια ώστε οι αποφασισμένες τιμές του A και του R να διαφέρουν σε αυτή τη θέση s. Μια «αποφασισμένη» θέση s στον A της οποίας η τιμή διαφέρει από αυτήν στον R λέγεται ότι σταματάει την δυνατότητα εφαρμογής του R στον A. Για παράδειγμα ο 0110 δεν μπορεί να εφαρμοστεί στον 2012, γιατί η δεύτερη θέση εμποδίζει την εφαρμογή. Η διαδικασία αυτή από τον (Gusfield 2001), μας δίνει μια ιδέα για το πώς λειτουργεί ένας αλγόριθμος εύρεσης απλότυπων.

Ο αλγόριθμος αυτός για να βγάλει συμπέρασμα για ένα σύνολο γονότυπων, πρέπει πρώτα να αναγνωρίσει το αρχικά «αποφασισμένο» σύνολο, και στη συνέχεια εφαρμόζει τον κανόνα απόφασης επαναληπτικά, μέχρι να χαρακτηριστούν όλοι οι γονότυποι «αποφασισμένοι» ή να μην υπάρχουν άλλοι γονότυποι που να μπορούν να χαρακτηριστούν «αποφασισμένοι».

Ο αλγόριθμος EM αντί να καταλήξει σε έναν απλότυπο για κάθε άτομο, στοχεύει στην εκτίμηση της συχνότητας των απλότυπων με την μέγιστη πιθανοφάνεια. Αντιμετωπίζει τους γονότυπους σαν καταμετρημένα, αλλά ημιτελή δεδομένα με άγνωστους απλότυπους. Αρχίζοντας από μια αυθαίρετη συχνότητα απλότυπων, ενημερώνει αυτή την συχνότητα περιοδικά μέσα από κάθε EM κύκλο μέχρι να επιτευχθεί σύγκλιση.

Οι προσεγγίσεις του Bayes περιλαμβάνουν μια ομάδα από πολύπλοκες διαδικασίες που βασίζονται στις Monte Carlo μεθόδους (Sun, Greenwood et al. 2007). Αυτές οι μέθοδοι είναι ένα είδος υπολογιστικών αλγόριθμων, οι οποίοι βασίζονται στην επανάληψη τυχαίας δειγματοληψίας για να υπολογίσουν αποτελέσματα. Χρησιμοποιούνται συνήθως στην προσομοίωση φυσικών, μαθηματικών, αλλά και βιολογικών συστημάτων.

Την υλοποίηση όλων αυτών των μεθόδων καλείται να πραγματοποιήσει ο ηλεκτρονικός υπολογιστής, που ως γνωστόν, μπορεί να εκτελέσει μαθηματικές πράξεις με πολύ μεγάλη ταχύτητα. Έχουν δημιουργηθεί πολλές εφαρμογές οι οποίες βρίσκουν τον απλότυπο χρησιμοποιώντας τις παραπάνω μεθόδους που περιγράψαμε (Xu, Wu et al. 2004).



**Εικόνα 10:** Η εφαρμογή Harlo-i χρησιμοποιεί ένα μοντέλο του Bayes για την επίλυση του προβλήματος εύρεσης απλότυπων

Τα πιο γνωστά προγράμματα τα οποία χρησιμοποιούν μεθόδους του Bayes (Sun, Greenwood et al. 2007) είναι τα: HAPLOTYPER, HAP(εικόνα 9), PHASE και HAPLOi (εικόνα 10). Και τα τρία αυτά προγράμματα δίνουν σαν αποτέλεσμα απλότυπους για κάθε άτομο. Οι διαφορές τους βρίσκονται κυρίως στις υποθέσεις για την προηγούμενη πιθανοτική κατανομή και στην υλοποίηση του υπολογιστικού αλγορίθμου. Για παράδειγμα το PHASE υποθέτει μια προηγούμενη κατανομή που προήλθε από ένα προσεγγιστικό μίγμα, ενώ το HAP και το HAPLOTYPER υποθέτουν μια Dirichlet προηγούμενη κατανομή (Xing, Jordan et al. 2007) η οποία ανήκει σε μια οικογένεια πολυπαραγοντικών κατανομών.

Το PHASE βρέθηκε να έχει τα πιο ακριβή αποτελέσματα στην εύρεση του απλότυπου, τον υπολογισμό των συχνοτήτων και στις μελέτες συσχέτισης. Επιπλέον τα αποτελέσματά του δεν είχαν κανένα λάθος συγκριτικά με τους απλότυπους που προήλθαν από βιολογικά πειράματα. Οι άλλες δυο μέθοδοι του HAP και HAPLOTYPER είχαν επίσης χαμηλό ποσοστό λαθών. Συμπερασματικά και οι τρεις μέθοδοι έδωσαν συνεπή με την εργαστηριακή προσέγγιση αποτελέσματα, σε λιγότερο χρόνο και κόστος (Xu, Wu et al. 2004).

## ΚΕΦΑΛΑΙΟ 2 – ΥΛΙΚΑ ΚΑΙ ΜΕΘΟΔΟΙ

### 2.1 Ορισμός του προβλήματος

Σκοπός αυτής της εργασίας είναι αρχικά να διερευνήσουμε την συμβολή των πολυμορφισμών στην συμπεριφορά κοινών παθήσεων (όπως διαβήτης, καρκίνος, παθήσεις της καρδιάς κλπ). Αν δηλαδή κάποια σύνολα απλότυπων ή ακόμη και μεμονωμένοι απλότυποι που βρίσκονται στο γονιδίωμα των ανθρώπων, παίζουν ή όχι σημαντικό ρόλο στην εμφάνιση κάποιας τέτοιας ασθένειας. Επίσης θέλουμε να προτείνουμε κάποιες απλές μεθόδους για την ανάλυση αυτού του τύπου βιολογικών δεδομένων. Τα δεδομένα των απλότυπων δίνονται από μελέτες που χωρίζονται σε δυο κύριες κατηγορίες. Μελέτες με ασθενείς και μάρτυρες, και μελέτες με μεταβλητή συνεχούς τιμής.

### 2.2 Η εφαρμογή Stata

Το Stata είναι ένα στατιστικό πακέτο γενικής φύσεως το οποίο δημιουργήθηκε το 1985 από την εταιρία StataCorp και είναι ευρέως διαδεδομένο τόσο σε εταιρίες, όσο και σε ακαδημαϊκά κέντρα σε όλο τον κόσμο. Οι περισσότεροι χρήστες το προτιμούν για έρευνα, ειδικά σε τομείς όπως της οικονομίας, της κοινωνιολογίας και της επιδημιολογίας.

Οι δυνατότητες της εφαρμογής αυτής περιλαμβάνουν:

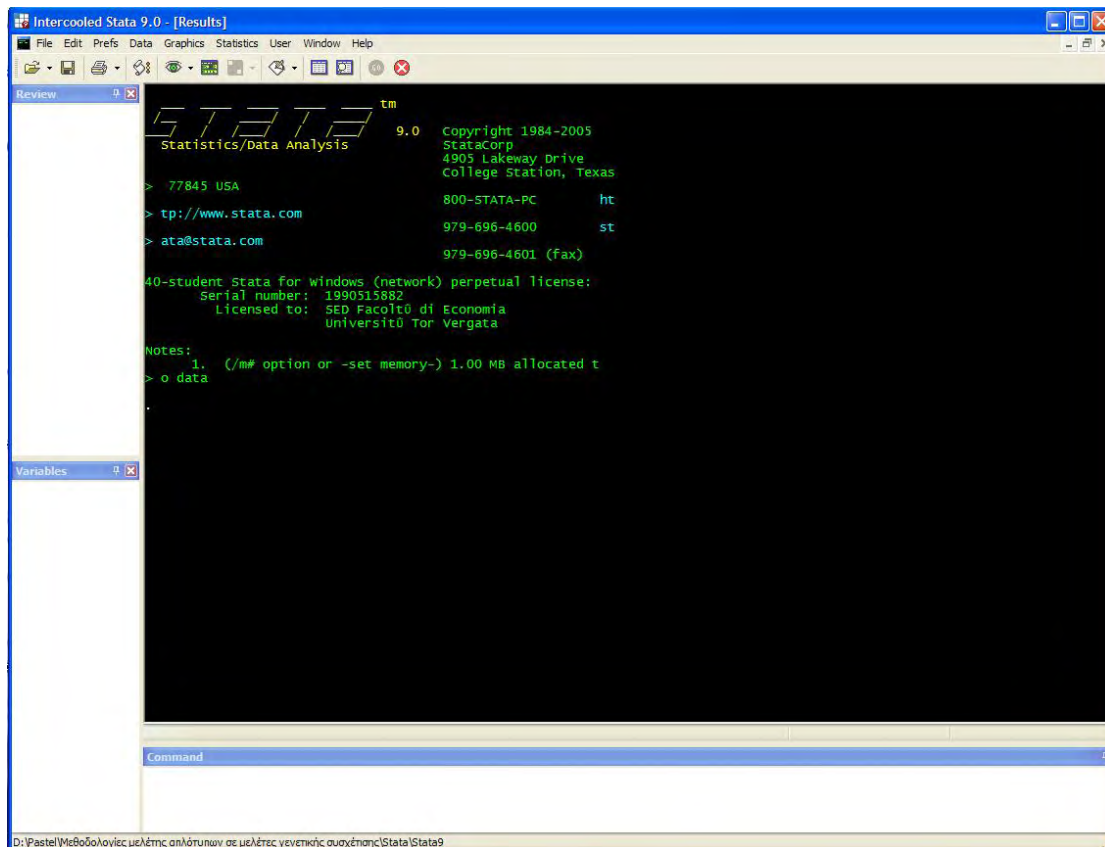
- Διαχείριση των δεδομένων
- Στατιστική ανάλυση
- Γραφικές απεικονίσεις
- Προσομοιώσεις
- Προγραμματισμός από τον χρήστη

Το όνομα Stata είναι μια συγχώνευση δυο λέξεων, statistics που σημαίνει στατιστική, και data, που σημαίνει δεδομένα. Υπάρχουν συνολικά τέσσερις εκδοχές του προγράμματος. Η πρώτη είναι κατασκευασμένη για πολυ-πύρινους υπολογιστές, η δεύτερη είναι ειδικά σχεδιασμένη για εξαιρετικά μεγάλες βάσεις δεδομένων, η τρίτη είναι η βασική έκδοση(και η πιο διαδεδομένη πιθανότατα) και τέλος η τέταρτη η οποία είναι για μικρού μεγέθους δεδομένα, κυρίως για εκπαιδευτική χρήση.

Στην κύρια οθόνη του προγράμματος υπάρχουν τέσσερα παράθυρα όπως φαίνεται και στην εικόνα 11, στο μεγαλύτερο από αυτά, το κεντρικό παράθυρο, εμφανίζονται τα αποτελέσματα των διεργασιών που εκτελεί η εφαρμογή (Results). Κάτω από το κύριο παράθυρο, είναι το παράθυρο των εντολών (Command), εκεί ο χρήστης πληκτρολογεί τις προς εκτέλεση εντολές. Τα δυο υπολειπόμενα παράθυρα έχουν βοηθητικό και πληροφοριακό ρόλο. Στο πάνω αριστερά, το επονομαζόμενο Review, καταγράφονται οι εντολές που έχουν εκτελεστεί από τον χρήστη μέχρι πρότινος. Και τέλος στο κάτω αριστερά, το επονομαζόμενο Variables βλέπουμε μια λίστα με τα ονόματα των υπαρχόντων μεταβλητών.

Παρόλο που υπάρχει το παράθυρο των εντολών και οι χρήστες εισάγουν τα δεδομένα και πραγματοποιούν στατιστικούς ελέγχους με την μορφή εντολών, η εφαρμογή παρέχει και εναλλακτικό τρόπο, χρησιμοποιώντας το κλασικό μενού επιλογών στο πάνω μέρος της εικόνας 11 (File Edit Prefs κ.λ.π.).

Θα πρέπει να σημειωθεί πως για την εισαγωγή δεδομένων από το excel θα πρέπει οι μεταβλητές να είναι μια λέξη, ή εφόσον είναι δύο η περισσότερες να μην χωρίζονται με κενό, αλλά με το σύμβολο της κάτω παύλας «\_».



Εικόνα 11: Η κύρια οθόνη του Stata

## 2.3 Εκτίμηση του Linkage disequilibrium

Στην πληθυσμιακή γενετική, ο όρος linkage disequilibrium που από εδώ και πέρα θα αναφερόμαστε ως LD, περιγράφει μη τυχαία συσχέτιση αλληλόμορφων σε δύο ή περισσότερα σημεία, τα οποία δεν βρίσκονται απαραίτητα στο ίδιο χρωμόσωμα. Δεν είναι το ίδιο με την γενετική συσχέτιση, η οποία περιγράφει την σχέση δύο ή περισσότερων θέσεων σε ένα χρωμόσωμα. Δύο περιοχές που βρίσκονται κοντά πάνω στην αλυσίδα του DNA είναι γενετικά συσχετισμένες με την έννοια ότι έχουν πολλές πιθανότητες μετά την διαδικασία του ανασυνδυασμού και της μείωσης να παραμείνουν και να κληρονομηθούν μαζί, στο ίδιο χρωμόσωμα.

Έτσι ο όρος LD περιγράφει μια κατάσταση στην οποία κάποιοι συνδυασμοί αλληλόμορφων ή γενετικών δεικτών παρατηρούνται λιγότερο ή περισσότερο σε έναν πληθυσμό από ότι θα ήταν αναμενόμενο από μια τυχαία διάταξη απλότυπων κάποιων αλληλόμορφων με βάση τις συχνότητές τους. Η μη-τυχαία συσχέτιση μεταξύ πολυμορφισμών σε διαφορετικές θέσεις μετριέται από τον βαθμό του LD. Αριθμητικά είναι η διαφορά μεταξύ των παρατηρηθέντων και των αναμενόμενων (υποθέτουμε τυχαία κατανομή) συχνοτήτων των αλληλόμορφων (Devlin and Risch 1995).

Ένα παράδειγμα είναι η εξάπλωση δύο σπάνιων ασθενειών στην Φιλανδία, όπου εκεί συγκριτικά με την υπόλοιπη Ευρώπη, η κυστική ίνωση εξαπλώνεται με αργό ρυθμό. Αντίθετα η έμφυτη χλωριούχος διάρροια (congenital chloride diarrhea) εξαπλώνεται γρήγορα. Και οι δύο ασθένειες οφείλονται σε μεταλλάξεις στο χρωμόσωμα 7, σε γειτονικά γονίδια. Κάτι τέτοιο μπορεί να σημαίνει διαφορές μεταξύ γενετικού υλικού Φιλανδών και υπόλοιπων Ευρωπαίων.

Ο βαθμός του LD επηρεάζεται από διάφορους παράγοντες όπως η γενετική συσχέτιση, ο ρυθμός του ανασυνδυασμού, ο ρυθμός των μεταλλάξεων και τη δομή του πληθυσμού. Για παράδειγμα, κάποιοι οργανισμοί όπως τα βακτήρια, τα οποία αναπαράγονται χωρίς μείωση καθότι απλοειδή, παρουσιάζουν LD επειδή απουσιάζει ο ανασυνδυασμός και συνεπώς οι περιοχές που συνδέονται παραμένουν σταθερές.

### 2.3.1 Είδη και μεθοδολογίες υπολογισμού του LD

Υπάρχουν διάφοροι τρόποι υπολογισμού της συσχέτισης όσον αφορά το LD. Για να κάνουμε τους υπολογισμούς, πρέπει πρώτα να ορίσουμε έναν πίνακα 2X2. Σύμφωνα με τους Devlin and Risch υποθέτουμε ότι έχουμε δύο περιοχές

ενδιαφέροντος στο γενετικό υλικό, και κάθε περιοχή έχει δύο αλληλόμορφα. Ένα αλληλόμορφο για την ασθένεια, και ένα κανονικό αλληλόμορφο τα οποία απομονώνονται στην πρώτη περιοχή ενδιαφέροντος. Τα άλλα δύο αλληλόμορφα τα οποία χρησιμοποιούνται ως γενετικοί δείκτες, απομονώνονται στην δεύτερη περιοχή (Devlin and Risch 1995). Σε αυτήν την περίπτωση τα δεδομένα μας μπορούν να πάρουν την μορφή του πίνακα 1. Παρατηρήστε τα σύνολα που καταγράφονται δεξιά και κάτω από τις απλοτυπικές συχνότητες, καθώς παίζουν σημαντικό ρόλο στους παρακάτω υπολογισμούς. Το  $n_{11}$  είναι ο αριθμός των απλότυπων στο δείγμα που φέρουν το αλληλόμορφο της ασθένειας για το A1. Το  $n_{A1}$  είναι ο συνολικός αριθμός των απλότυπων που φέρουν το A1 αλληλόμορφο,  $n_{\alpha}$  είναι ο αριθμός των απλότυπων που φέρουν το αλληλόμορφο της ασθένειας και  $n$  είναι ο συνολικός αριθμός των απλότυπων του δείγματος.

Δείκτης	Αλληλόμορφο για την ασθένεια	Κανονικό αλληλόμορφο	
A1	$n_{11}$	$n_{12}$	$n_{A1}$
A2	$n_{21}$	$n_{22}$	$n_{A2}$
	$n_{\alpha}$	$n_{\kappa}$	$n$

**Πίνακας 1: Διάταξη συχνοτήτων των απλότυπων σε 2X2 πίνακα**

Διαιρώντας αυτές τις ποσότητες με το  $n$  παίρνουμε τις αναμενόμενες συχνότητες και τις πιθανότητες  $p$  του δείγματος, (πίνακας 2). Αυτό μας βοηθάει σε διάφορους υπολογισμούς όπως για παράδειγμα, να βρούμε την πιθανότητα του να υπάρχει το αλληλόμορφο A1 στον απλότυπο, δεδομένου ότι έχουμε αλληλόμορφο της ασθένειας. Αυτό σημαίνει ότι  $P_{A1/ασθένεια} = P_{11} / P_{\alpha}$ . Ομοίως η πιθανότητα να έχουμε κανονικό αλληλόμορφο στον απλότυπο δεδομένου ότι υπάρχει μόνο το αλληλόμορφο A2 είναι  $P_{κανονικό/A2} = P_{22} / P_{A2}$ .

Βέβαια αυτές οι ποσότητες ( $p$ ) είναι μόνο μια εκτίμηση κάποιων άγνωστων παραμέτρων, οι οποίες θα συμβολιστούν με  $\pi$ . Χρησιμοποιούμε τα  $\pi$  στους ακόλουθους ορισμούς, με την προϋπόθεση ότι αυτές οι άγνωστες ποσότητες υπολογίζονται από τις παρατηρούμενες ποσότητες του δείγματος.

Δείκτης	Αλληλόμορφο για την ασθένεια	Κανονικό αλληλόμορφο	
A1	$p_{11}$	$p_{12}$	$p_{A1}$
A2	$p_{21}$	$p_{22}$	$p_{A2}$
	$p_{\alpha}$	$p_{\kappa}$	1

**Πίνακας 2:** Κάθε πιθανότητα προκύπτει από τα στοιχεία του πίνακα 1 διαιρούμενα με το  $n$

Ο κοινός παρονομαστής για κάθε μέτρο του LD είναι, όπως αναφέραμε και πριν, η διαφορά μεταξύ παρατηρηθέντων και αναμενόμενων (υπό ανεξαρτησία) αριθμών των απλότυπων που έχουν το αλληλόμορφο της ασθένειας και βρίσκονται στον δείκτη A1, ή τις ισοδύναμες εκφράσεις της σχέσης 1.

$$\begin{aligned}
 D &= \pi_{11} - \pi_{A1}\pi_{\alpha} = \pi_{22} - \pi_{A2}\pi_{\kappa} \\
 &= -\pi_{12} + \pi_{A1}\pi_{\kappa} = -\pi_{21} + \pi_{A2}\pi_{\alpha} \\
 &= \pi_{11}\pi_{22} - \pi_{12}\pi_{21}
 \end{aligned}$$

(1)

Σύμφωνα με τους Hill and Weir, το πιο ευρέως χρησιμοποιούμενο μέτρο για τον υπολογισμό του LD είναι το τετράγωνο του μέτρου  $\Delta$ .

$$\Delta = \frac{\pi_{11}\pi_{22} - \pi_{12}\pi_{21}}{(\pi_{A1}\pi_{A2}\pi_{\alpha}\pi_{\kappa})^{1/2}}$$

(2)

ή αλλιώς το  $\Delta^2$  (Hill and Weir 1994).

Άλλο ένα μέτρο που χρησιμοποιείται συχνά είναι το  $D'$ .

$$D' = \frac{\pi_{11}\pi_{22} - \pi_{12}\pi_{21}}{\min(\pi_{A1}\pi_{\kappa}, \pi_{\alpha}\pi_{A2})} D > 0$$

$$D' = \frac{\pi_{11}\pi_{22} - \pi_{12}\pi_{21}}{\min(\pi_{A1}\pi_{\alpha}, \pi_{\kappa}\pi_{A2})} D < 0$$
(3)

όπου η ποσότητα στον παρονομαστή είναι το μέγιστο  $D$  που θα μπορούσε να επιτευχθεί με τα δεδομένα του πίνακα. Αυτά τα δυο μέτρα ( $D', \Delta$ ) σχετίζονται με στατιστικές μεθόδους υπολογισμού της συσχέτισης. Συγκεκριμένα το  $\Delta$  είναι ο συντελεστής συσχέτισης για έναν πίνακα 2X2. Επίσης το  $\Delta$  είναι ανάλογο του  $r$  (σχέση 4) για έναν πίνακα 2X2.

$$r_{ij} = \frac{n_{ij} - \hat{m}_{ij}}{(\hat{m}_{ij}(1 - p_{Ai})(1 - p_{\alpha/\kappa}))^{1/2}}$$
(4)

όπου το  $p_{\alpha/\kappa}$  υπολογίζεται ανάλογα με το αν οι δείκτες  $i, j$  αντιστοιχούν σε στήλη με το αλληλόμορφο της ασθένειας ή με το κανονικό αλληλόμορφο, και το  $m_{ij}$  είναι ο αναμενόμενος αριθμός στο κελί  $ij$ .

Άλλο ένα μέτρο συσχέτισης που χρησιμοποιείται στην επιδημιολογία για τον υπολογισμό του LD είναι το  $\delta^*$ , το οποίο ορίζεται ως:

$$\delta^* = \frac{\pi_{A1} + (\varphi - 1)}{1 + \pi_{A1}(\varphi - 1)}$$
(5)

όπου  $\varphi = \{\pi_{11} / \pi_{A1}\} / \{\pi_{21} / \pi_{A2}\}$  είναι ο σχετικός κίνδυνος ή αλλιώς odds ratio.

Το  $\delta^*$  μπορούμε να το βρούμε μετά από κάποιες αλλαγές και στην μορφή της σχέσης 6.



$$\delta = \frac{\pi_{11}\pi_{22} - \pi_{12}\pi_{21}}{\pi_{\alpha}\pi_{22}} \quad (6)$$

η οποία είναι μια προσέγγιση του κινδύνου του πληθυσμού σε μια μεγαλύτερη κλίμακα από αυτή των πληθυσμιακών συχνοτήτων τους.

Αυτό το μέτρο  $\delta^*$  δεν είναι καινούργιο καθώς, όταν η ασθένεια είναι σπάνια και οι απλότυποι έχουν δειγματοληφθεί τυχαία, τότε ισχύει ότι  $\delta^* = D'$ :

$$\delta^* = \frac{\pi_{11}\pi_{22} - \pi_{12}\pi_{21}}{\pi_{\alpha}\pi_{A2}} \quad (7)$$

όμως

$$\begin{aligned} D' &= \frac{\pi_{11}\pi_{22} - \pi_{12}\pi_{21}}{\min(\pi_{A1}\pi_{\kappa}, \pi_{\alpha}\pi_{A2})} \\ &= \frac{\pi_{11}\pi_{22} - \pi_{12}\pi_{21}}{\pi_{\alpha}\pi_{A2}} \end{aligned} \quad (8)$$

βλέπουμε ότι τα δύο κλάσματα των σχέσεων (7) και (8) συμπίπτουν.

Άλλο ένα μέτρο το οποίο ενδείκνυται για χρήση σε μελέτες ασθενών-μαρτύρων είναι το  $d$ , το οποίο ορίζεται με την παρακάτω σχέση:

$$d = \frac{\pi_{11}}{\pi_{\alpha}} - \frac{\pi_{12}}{\pi_{\kappa}} = \frac{\pi_{11}\pi_{22} - \pi_{12}\pi_{21}}{\pi_{\alpha}\pi_{\kappa}} \quad (9)$$

Επίσης υπάρχουν και άλλα φυσικά μέτρα για το LD, όπως είναι ο σχετικός λόγος των πιθανοτήτων (odds ratio) αλλά και το Q του Yule. (Clegg, Kidwell et al. 1976; Nei and Li 1980; Olson and Wijsman 1994).

Ο λόγος των πιθανοτήτων ορίζεται ως:

$$OR = \frac{\pi_{11}\pi_{22}}{\pi_{12}\pi_{21}} \quad (10)$$

με εύρος τιμών  $[0, +\infty]$ , ενώ το Q ως:

$$Q = \frac{OR - 1}{OR + 1} = \frac{\pi_{11}\pi_{22} - \pi_{12}\pi_{21}}{\pi_{11}\pi_{22} + \pi_{12}\pi_{21}} \quad (11)$$

με εύρος τιμών  $[-1, +1]$ . Το τελευταίο κλάσμα της εξίσωσης 11 μας δείχνει την σχέση του Q με το δ. Στην πραγματικότητα οι αριθμητές των Δ, D', δ, d και Q είναι όλοι ίσοι με το D, καθώς αυτά τα μέτρα διαφέρουν μόνο κατά τους παρονομαστές τους, οι οποίοι έχουν την χρήση της τυποποίησης του D. Στον πίνακα 3 βλέπουμε συνοπτικά τα προαναφερθέντα μέτρα.

Σύμβολο	Τύπος
$\Delta$	$\frac{\pi_{11}\pi_{22} - \pi_{12}\pi_{21}}{(\pi_{A1}\pi_{A2}\pi_{\alpha}\pi_{\kappa})^{1/2}}$
$D'$	$\frac{\pi_{11}\pi_{22} - \pi_{12}\pi_{21}}{\min(\pi_{A1}\pi_{\kappa}, \pi_{\alpha}\pi_{A2})}$
$\delta$	$\frac{\pi_{11}\pi_{22} - \pi_{12}\pi_{21}}{\pi_{\alpha}\pi_{22}}$
$d$	$\frac{\pi_{11}\pi_{22} - \pi_{12}\pi_{21}}{\pi_{\alpha}\pi_{\kappa}}$
$Q$	$\frac{\pi_{11}\pi_{22} - \pi_{12}\pi_{21}}{\pi_{11}\pi_{22} + \pi_{12}\pi_{21}}$

Πίνακας 3: Μια σύνοψη των μέτρων υπολογισμού του LD

## 2.4 Μελέτες με ασθενείς και μάρτυρες

Οι μελέτες με ασθενείς και μάρτυρες είναι ένα είδος επιδημιολογικών μελετών. Χρησιμοποιούνται για την αναγνώριση κάποιων παραγόντων που επηρεάζουν ασθένειες, είτε θετικά είτε αρνητικά. Η μελέτη γίνεται πάνω σε έναν πληθυσμό ατόμων τα οποία αποτελούνται από άτομα που έχουν κάποια ασθένεια και από άτομα υγιή, χωρίς όμως να παρουσιάζουν σημαντικές διαφορές από τους ασθενείς (Fallin, Cohen et al. 2001), έχουν δηλαδή τις ίδιες συνθήκες ζωής και έχουν επιλεγεί με τα ίδια κριτήρια.

Αυτού του είδους οι μελέτες είναι ένα πολύ φθηνό εργαλείο που χρησιμοποιείται ευρέως και με μεγάλη ευκολία, από μικρές ομάδες ή ερευνητές. Το βασικό πλεονέκτημα είναι ότι δεν απαιτείται ιδιαίτερος εργαστηριακός εξοπλισμός καθώς βασίζονται στην στατιστική και όχι στην εργαστηριακή ανάλυση. Έχουν οδηγήσει σε σημαντικές ανακαλύψεις, αλλά έχουν το αρνητικό χαρακτηριστικό ότι μπορεί να υπάρχουν παράγοντες που δεν ξέρουμε, οι οποίοι επηρεάζουν κατά ένα ποσοστό το αποτέλεσμα.

Παρακάτω ακολουθεί μια εικόνα ενός πίνακα ασθενών-μαρτύρων (εικόνα 12). Στις στήλες του πίνακα διαχωρίζονται οι ασθενείς και οι μάρτυρες, ενώ στις γραμμές βλέπουμε τους διάφορους παράγοντες που μπορεί να σχετίζονται με την ασθένεια. Κάποιοι βασικοί παράγοντες είναι το κάπνισμα, ο τόπος κατοικίας, η ηλικία και η ύπαρξη άλλων μορφών καρκίνου. Κάτω από την γραμμή κάθε παράγοντα, παρατηρούμε και το p-value, το οποίο δείχνει αν απορρίπτουμε ή δεχόμαστε την μηδενική υπόθεση, η οποία είναι ότι δεν υπάρχει συσχέτιση παράγοντα-ασθένειας.

**Table 1.** Selected Characteristics of Lung Cancer Cases and Controls With Interview Data Available, the EAGLE Study, Lombardy, Italy, 2002–2005<sup>a,b</sup>

	Women				Men			
	Cases		Controls		Cases		Controls	
	No.	%	No.	%	No.	%	No.	%
Total participants enrolled	448		500		1,652		1,620	
Interviewed	406	100.0	499	100.0	1,537	100.0	1,617	100.0
Area of residence								
Milan	288	70.9	349	69.9	987	64.2	1,089	67.3
Monza	24	5.9	23	4.6	109	7.1	94	5.8
Brescia	47	11.6	53	10.6	203	13.2	194	12.0
Pavia	21	5.2	37	7.4	107	7.0	92	5.7
Varese	26	6.4	37	7.4	131	8.5	148	9.2
				$P = 0.55$				$P = 0.17$
Age, years (mean (SD))	64.8 (10.1)		64.1 (10.1)		66.8 (7.9)		65.8 (8.1)	
				$P = 0.32$				$P < 0.001$
Educational level								
None	21	5.2	24	4.8	91	5.9	66	4.1
Elementary school	128	31.5	143	28.7	625	40.7	431	26.7
Middle school	134	33.0	158	31.7	424	27.6	455	28.1
High school	104	25.6	135	27.1	314	20.4	441	27.3
University	19	4.7	39	7.8	83	5.4	224	13.9
				$P = 0.35$				$P < 0.001$
No. of jobs held								
1	166	40.9	168	33.7	375	24.4	370	22.9
2	96	23.7	158	31.7	404	26.3	356	22.0
3	77	19.0	82	16.4	305	19.8	356	22.0
4	30	7.4	49	9.8	194	12.6	226	14.0
≥5	37	9.1	42	8.4	259	16.9	309	19.1
				$P = 0.03$				$P = 0.02$

	Women				Men			
	Cases		Controls		Cases		Controls	
	No.	%	No.	%	No.	%	No.	%
Cigarette smoking								
Never	103	25.4	282	56.5	29	1.9	397	24.6
Former (quit >6 months ago)	116	28.6	110	22.0	723	47.0	799	49.4
Current	187	46.1	107	21.4	785	51.1	420	26.0
Unknown	0	0.0	0	0.0	0	0.0	1	0.1
				$P < 0.001$				$P < 0.001$
Cigarette pack-years (mean (SD))	24.3 (23.1)		7.2 (13.5)		50.9 (28.7)		22.1 (23.2)	
				$P < 0.001$				$P < 0.001$
Other cancer(s) <sup>c</sup>								
No	336	82.8	448	89.8	1,306	85.0	1,473	91.1
Yes	70	17.2	51	10.2	231	15.0	144	8.9
				$P = 0.002$				$P < 0.001$
Lung cancer morphology								
Adenocarcinoma	220	54.2			582	37.9		
Squamous cell carcinoma	45	11.1			459	29.9		
Large cell carcinoma	28	6.9			61	4.0		
Non-small-cell carcinoma NOS	34	8.4			142	9.2		
Small cell carcinoma	38	9.4			157	10.2		
Other	26	6.4			65	4.2		
Not available	15	3.7			71	4.6		
								$P < 0.001$

Abbreviations: EAGLE, Environment And Genetics in Lung cancer Etiology; NOS, not otherwise specified; SD, standard deviation.

<sup>a</sup>  $P$  values were derived from the chi-square test (categorical variables) or Student's  $t$  test (continuous variables).

<sup>b</sup> Percentages may not add to 100.0 because of rounding.

<sup>c</sup> Primary cancer(s) (previously or newly diagnosed) other than lung cancer.

**Εικόνα 12: Χαρακτηριστικά ασθενών-μαρτύρων από μια μελέτη για την σχέση μεταξύ καρκίνου του πνεύμονα και επαγγέλματος (Consonni, De Matteis et al. 2010)**

Σε αυτήν εδώ την εργασία τα δεδομένα έχουν την ίδια μορφή με την παραπάνω μελέτη. Η μόνη διαφορά είναι ότι οι παράγοντες κινδύνου δεν είναι ένα σύνολο από διατροφικές συνήθειες, ηλικίες κ.λ.π., αλλά από ένα σύνολο διαφορετικών απλότυπων. Έτσι η έρευνα επικεντρώνεται στο να διαπιστώσουμε κάποιες διαφορές ανάμεσα σε άτομα με διαφορετικό απλότυπο ώστε να μπορέσουμε να κάνουμε μια εκτίμηση για τον βαθμό που επηρεάζεται η ασθένεια από αυτόν. Συνεπώς για έναν πίνακα ασθενών και μαρτύρων θα έχουμε στην πρώτη στήλη τους ασθενείς (cases) και στην δεύτερη στήλη τους μάρτυρες (controls). Ενώ στην πρώτη γραμμή, τον απλότυπο τύπου 1, στην δεύτερη τον απλότυπο τύπου 2 κ.ο.κ.

	ασθενείς (cases)	μάρτυρες (controls)
απλότυπος 1	252	248
απλότυπος 2	289	234
...		
...		
...		
...		
...		
...		
απλότυπος n	234	654

**Πίνακας 4: Η κύρια μορφή των δεδομένων στις μελέτες με ασθενείς και μάρτυρες**

Με βάση τον πίνακα 4 μπορούμε να πραγματοποιήσουμε διάφορες στατιστικές διεργασίες με το πρόγραμμα Stata, οι οποίες θα μας βοηθήσουν να εξαγάγουμε κάποια συμπεράσματα. Τα δεδομένα σε αυτήν την εργασία έχουν την παρακάτω μορφή (εικόνα 13):

Microsoft Excel - testing

Αρχείο Επεξεργασία Προβολή Εισαγωγή Μορφή Εργαλεία Δεδομένα Παραβύρο Βοήθεια

A8 haplotype

	A	B	C	D	E	F	G	H	I	J	K
8	haplotype	haplotype_nr	frequency	frequency_percent	case	cases_sum	x1	x2	x3	x4	x5
9	GGGCA	1	1147	0,621	1	1827	0	0	0	0	0
10	GGTTG	2	298	0,1613	1	1827	0	0	1	1	1
11	GGGCG	3	163	0,0881	1	1827	0	0	0	0	1
12	AGGCA	4	114	0,0675	1	1827	1	0	0	0	0
13	GAGCA	5	105	0,0568	1	1827	0	1	0	0	0
14	GGGCA	1	529	0,6347	0	0	0	0	0	0	0
15	GGTTG	2	132	0,1589	0	0	0	0	1	1	1
16	GGGCG	3	73	0,088	0	0	0	0	0	0	1
17	AGGCA	4	51	0,0615	0	0	1	0	0	0	0
18	GAGCA	5	44	0,0526	0	0	0	1	0	0	0
19											

Εικόνα 13: Η μορφή των δεδομένων που εισάγονται στο Stata

Αφού αντιγράψουμε τα δεδομένα από τον πίνακα του excel, για να τα εισάγουμε στο Stata πληκτρολογούμε την εντολή «edit» στην κύρια οθόνη του, και στην συνέχεια στο παράθυρο που εμφανίζεται κάνουμε επικόλληση.

### 2.4.1 Ο έλεγχος $\chi^2$ του Pearson για την ανεξαρτησία δυο δειγμάτων

Υπάρχουν δύο έλεγχοι του Pearson. Ο πρώτος αναφέρεται σαν τεστ καλής προσαρμογής, και ο δεύτερος (ο οποίος χρησιμοποιείται και αναλύεται εδώ) αναφέρεται ως τεστ ανεξαρτησίας. Με αυτό τον έλεγχο, έχοντας δυο σύνολα δεδομένων, μπορούμε να ελέγξουμε κατά πόσο αυτά τα σύνολα είναι στατιστικώς συνδεδεμένα. Σαν μηδενική υπόθεση έχουμε το γεγονός ότι τα δυο σύνολα είναι στατιστικώς ανεξάρτητα. Τα δεδομένα οργανώνονται σε έναν πίνακα με  $r$  γραμμές και  $c$  στήλες, έπειτα υπολογίζεται για κάθε κελί του πίνακα η ακόλουθη θεωρητική συχνότητα:

$$E_{i,j} = \frac{\sum_{k=1}^c O_{i,k} \sum_{k=1}^r O_{k,j}}{N} \quad (12)$$

Αυτή η προσαρμογή μειώνει τους βαθμούς ελευθερίας κατά  $p = r + c - 1$ .

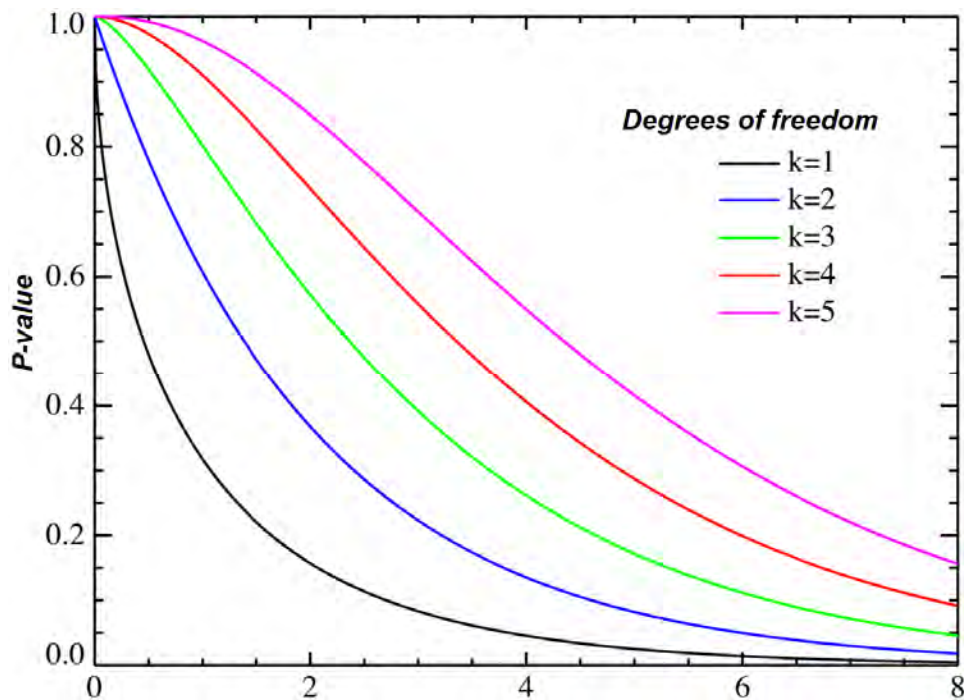
Τέλος υπολογίζουμε την τιμή του  $X^2$  χρησιμοποιώντας την σχέση:

$$X^2 = \sum_{i=1}^r \sum_{j=1}^c \frac{(O_{i,j} - E_{i,j})^2}{E_{i,j}} \quad (13)$$

Ο αριθμός των βαθμών ελευθερίας είναι ίσος με τον αριθμό των κελιών του πίνακα, μείον τον αριθμό  $p$ .

Αν από αυτό το τεστ έχουμε σαν αποτέλεσμα πιθανότητα  $X^2$  ίση ή μικρότερη από 0.05, τότε η μηδενική υπόθεση απορρίπτεται και μπορούμε να πούμε ότι υπάρχει τυχαία συσχέτιση μεταξύ των δυο συνόλων. Αυτήν την πιθανότητα την υπολογίζουμε από την γραφική παράσταση της κατανομής  $X^2$  (εικόνα 14).





Εικόνα 14: Η κατανομή  $\chi^2$

Παραδείγματος χάρη για να πραγματοποιήσουμε αυτόν τον έλεγχο σε δεδομένα όπως αυτά της εικόνας 13, αφού πρώτα τα εισάγουμε στο Stata, εκτελούμε την εντολή:

```
tabulate haplotype_nr case [fweight=frequency], row col chi2
```

Η εντολή `tabulate` δημιουργεί πίνακα συχνοτήτων μεταξύ των μεταβλητών `haplotype_nr` και `case`. Επιπλέον ορίζουμε ότι θέτουμε ως βάρος την συχνότητα, με την παράμετρο `[fweight=frequency]`, και ζητάμε να γίνει ο έλεγχος  $\chi^2$  με αναφορά ταυτόχρονα των σχετικών συχνοτήτων για κάθε γραμμή και κάθε στήλη (`row col`).

Τα αποτελέσματα της οποίας είναι τα ακόλουθα (εικόνα 15):

key			
frequency			
row percentage			
column percentage			
haplotype_	case		Total
nr	0	1	
1	3,283 57.65 44.08	2,412 42.35 48.70	5,695 100.00 45.92
2	2,623 61.73 35.22	1,626 38.27 32.83	4,249 100.00 34.26
3	674 59.80 9.05	453 40.20 9.15	1,127 100.00 9.09
4	504 64.86 6.77	273 35.14 5.51	777 100.00 6.27
5	161 75.94 2.16	51 24.06 1.03	212 100.00 1.71
6	48 59.26 0.64	33 40.74 0.67	81 100.00 0.65
7	155 59.62 2.08	105 40.38 2.12	260 100.00 2.10
Total	7,448 60.06 100.00	4,953 39.94 100.00	12,401 100.00 100.00

Pearson chi2(6) = 48.6233 Pr = 0.000

**Εικόνα 15: Αποτελέσματα του  $\chi^2$  του Pearson**

Στα αποτελέσματα αυτά παρατηρούμε ότι ο  $\chi^2$  έχει τιμή περίπου 48.6 και το p-value είναι κάτω από 0.05, αυτό δείχνει ότι υπάρχει τυχαία συσχέτιση μεταξύ του συνόλου των απλότυπων, και της ασθένειας. Ο έλεγχος αυτός έχει εισαγωγικό χαρακτήρα, καθώς όταν βλέπουμε πως υπάρχει τυχαία συσχέτιση μεταξύ των στοιχείων, μας προτρέπει να ασχοληθούμε παραπάνω με την συγκεκριμένη μελέτη. Στην αντίθετη περίπτωση κάτι τέτοιο θα ήταν πιθανότατα περιττό.

## 2.4.2 Εκτίμηση μέγιστης πιθανοφάνειας

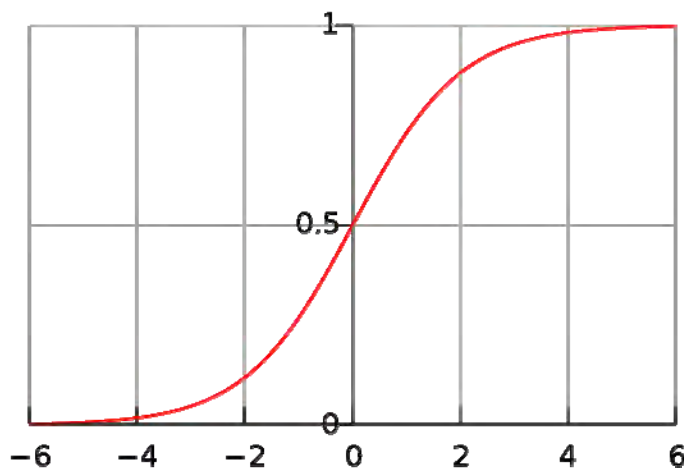
Η εκτίμηση μέγιστης πιθανοφάνειας είναι μια μέθοδος που χρησιμοποιούμε για την προσαρμογή ενός στατιστικού μοντέλου σε κάποια δεδομένα, καθώς επίσης και να κάνουμε μια εκτίμηση για τις παραμέτρους του μοντέλου αυτού. Για παράδειγμα ας υποθέσουμε ότι θέλουμε να ασχοληθούμε με το ύψος ενός πληθυσμού. Έχουμε ένα δείγμα ατόμων, αλλά όχι ολόκληρο των πληθυσμό (της χώρας, της κοινότητας κλπ), και καταγράφουμε τα ύψη τους. Θα πρέπει να υποθέσουμε ότι τα ύψη ακολουθούν κανονική κατανομή με κάποια άγνωστη μέση τιμή και διακύμανση. Τότε ο μέσος όρος του δείγματος είναι η εκτιμήτρια μέγιστης

πιθανοφάνειας του μέσου όρου ολόκληρου του πληθυσμού, όπως και η διακύμανση του δείγματος είναι η εκτιμήτρια μέγιστης πιθανοφάνειας της διακύμανσης του πληθυσμού.

### 2.4.3 Λογιστική παλινδρόμηση

Στη στατιστική η λογιστική παλινδρόμηση χρησιμοποιείται για την πρόβλεψη της πραγματοποίησης ενός γεγονότος, με την προσαρμογή κάποιων δεδομένων σε μια λογιστική συνάρτηση. Είναι ένα γενικευμένο γραμμικό μοντέλο που χρησιμοποιείται για την εξαγωγή αποτελεσμάτων διωνυμικού τύπου. Δηλαδή ισχύει-δεν ισχύει. Για την πρόβλεψη αυτή γίνεται χρήση διαφόρων μεταβλητών, είτε κατηγορικών, είτε αριθμητικών (Agresti 2002).

Για παράδειγμα η πιθανότητα για κάποιον να πάθει έμφραγμα μέσα σε μια συγκεκριμένη περίοδο, μπορεί να προβλεφθεί από τις γνώσεις που έχουμε για την ηλικία, το φύλο και τον δείκτη μάζας σώματος. Αυτό το μοντέλο χρησιμοποιείται σε μεγάλο βαθμό στις ιατρικές και κοινωνικές επιστήμες, αλλά και σε εφαρμογές του marketing όπως την πρόβλεψη των τάσεων του πελάτη για την αγορά ενός προϊόντος. Για να εξηγήσουμε την λογιστική παλινδρόμηση, πρέπει πρώτα να εξηγήσουμε την λογιστική συνάρτηση (εικόνα 16).



Εικόνα 16: Η συνάρτηση  $f(z)$  με διάφορες τιμές του  $z$  στον άξονα των  $x$

Σαν είσοδο της συνάρτησης έχουμε τον αριθμό  $z$ , και σαν έξοδο έχουμε το  $f(z)$ . Ο  $z$  παίρνει τιμές από το  $-\infty$  έως το  $+\infty$  ενώ το πεδίο τιμών της  $f(z)$  είναι το  $(0,1)$ . Η μεταβλητή  $z$  απεικονίζει την έκθεση σε ένα σύνολο από ανεξάρτητες

μεταβλητές, ενώ η  $f(z)$  απεικονίζει την πιθανότητα ενός συγκεκριμένου αποτελέσματος, με βάσει τις παραπάνω μεταβλητές.

Ο τύπος της συνάρτησης είναι ο εξής:

$$f(z) = \frac{1}{1 + e^{-z}} = \frac{e^z}{e^z + 1} \quad (14)$$

Επιπλέον το  $z$  εκφράζει την συνολική συνεισφορά όλων των ανεξάρτητων μεταβλητών που χρησιμοποιήθηκαν στο μοντέλο, και είναι γνωστό και ως logit.

$$z = \text{logit}(p) \quad (15)$$

Το logit ενός αριθμού  $p$  από 0 μέχρι 1 δίνεται από τον παρακάτω τύπο:

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right) = \log(p) - \log(1-p) \quad (16)$$

Η μεταβλητή  $z$  ορίζεται συνήθως σύμφωνα με τον παρακάτω τύπο:

$$z = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \dots + \beta_k x_k \quad (17)$$

όπου το στοιχείο  $\beta_0$  ονομάζεται στοιχείο αναφοράς και τα  $\beta_1, \beta_2, \beta_3, \dots, \beta_k$  ονομάζονται συντελεστές παλινδρόμησης των  $x_1, x_2, x_3, \dots, x_k$  αντίστοιχα.

Το στοιχείο αναφοράς είναι η τιμή του  $z$  όταν η τιμή όλων των ανεξάρτητων μεταβλητών είναι μηδέν (για παράδειγμα η τιμή του  $z$  όταν απεικονίζει την πιθανότητα καρδιακής προσβολής για ένα άτομο, χωρίς επιβαρυντικούς παράγοντες).

Κάθε ένας από τους συντελεστές παλινδρόμησης παριστάνει το μέγεθος της συνεισφοράς του αντίστοιχου επιβαρυντικού παράγοντα.

Ένας θετικός συντελεστής παλινδρόμησης δείχνει ότι ο αντίστοιχος παράγοντας αυξάνει την πιθανότητα πραγματοποίησης του αποτελέσματος, ενώ ένας αρνητικός συντελεστής παλινδρόμησης σημαίνει ότι ο αντίστοιχος παράγοντας μειώνει την πιθανότητα πραγματοποίησης του αποτελέσματος.

Επιπλέον ένας μεγάλος συντελεστής παλινδρόμησης σημαίνει ότι ο αντίστοιχος παράγοντας επηρεάζει κατά πολύ την πιθανότητα πραγματοποίησης του αποτελέσματος, ενώ ένας συντελεστής παλινδρόμησης κοντά στο μηδέν έχει μικρή επιρροή στην πιθανότητα πραγματοποίησης του αποτελέσματος.

Για να προσαρμόσουμε αυτό το μοντέλο στα δεδομένα μας, πρέπει να σκεφτούμε πως έχουμε μια κατηγορική μεταβλητή (case) για τους ασθενείς/μάρτυρες, η οποία παίζει το ρόλο της εξαρτημένης μεταβλητής, και μετά έχουμε τους απλότυπους που επηρεάζουν την κατάσταση της μεταβλητής αυτής θετικά ή αρνητικά, συνεπώς χρησιμοποιούνται ως ανεξάρτητες μεταβλητές. Έτσι μπορούμε να δημιουργήσουμε μια πιθανότητα εμφάνισης της ασθένειας, λαμβάνοντας υπόψη την επιρροή των απλότυπων. Έχουμε λοιπόν την πιθανότητα  $\pi_j = P(y_j = 1 | j)$  του να είναι κάποιος ασθενής (δηλαδή  $y=1$ ), έχοντας τον  $j$  απλότυπο. Το καλύτερο σενάριο θα ήταν να διαλέξουμε τον πιο κοινό απλότυπο (δηλαδή τον απλότυπο με την μεγαλύτερη συχνότητα) ως το στοιχείο αναφοράς, και να φτιάξουμε μια μεταβλητή για κάθε έναν από τους υπόλοιπους. Η μεταβλητή αυτή θα παίρνει τιμές  $z_j = 1$  για τον  $j$  απλότυπο, και 0 σε όλες τις άλλες περιπτώσεις. Στην εικόνα 17 βλέπουμε ένα παράδειγμα τέτοιας μεταβλητής για τον απλότυπο 3, τον τρίτο δηλαδή κατά σειρά απλότυπο. Παρατηρούμε ότι μόνο η τρίτη θέση έχει τον αριθμό 1. Αν αναφερόμασταν πχ στον απλότυπο 6, το 1 θα ήταν στην έκτη κατά σειρά θέση. Εδώ οι απλότυποι είναι 7 στο σύνολο και απλά η διαδικασία επαναλαμβάνεται πολλές φορές. Αυτό το μοντέλο λοιπόν μπορεί να πραγματοποιηθεί ως εξής, (Wallenstein, Hodge et al. 1998):

$$\log it(\pi_j) = \log it \left[ P(y_j = 1 | j) \right] = \beta_0 + \sum_{j=2}^r \beta_j z_j \quad (18)$$

The screenshot shows the Stata 9.0 Results window. The command window contains the following commands: `help vwls`, `edit`, `xi: logit case i.haplotype_nr`, `list _Ihap1~3`, `list _Ihap1~2`, and `list _Ihap1~2`. The Results window displays the following data:

Obs	_Ihap1~3
1.	0
2.	0
3.	1
4.	0
5.	0
6.	0
7.	0
8.	0
9.	0
10.	1
11.	0
12.	0
13.	0
14.	0
15.	0
16.	0
17.	1
18.	0
19.	0
20.	0
21.	0
22.	0
23.	0
24.	1
25.	0
26.	0
27.	0
28.	0
29.	0
30.	0

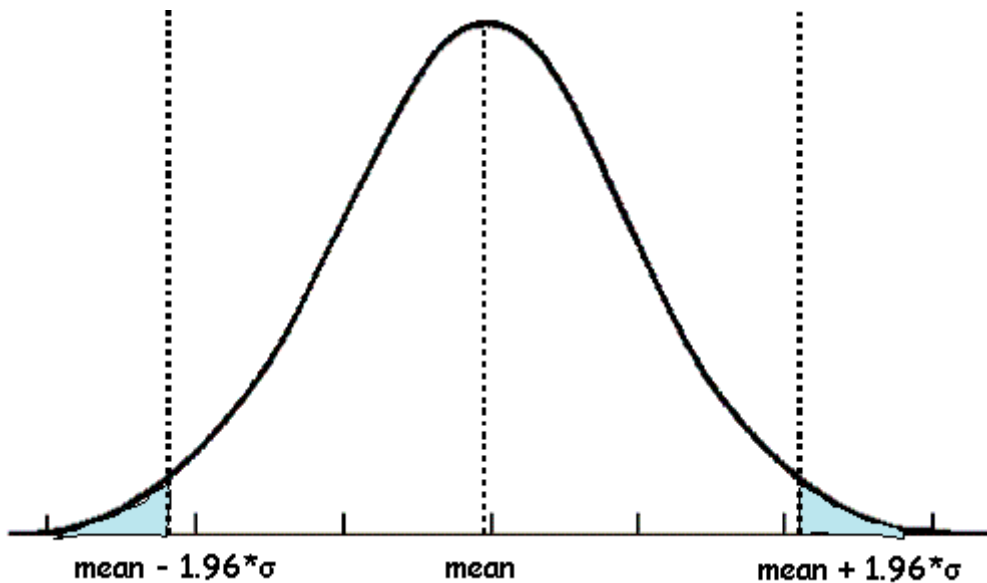
Εικόνα 17: Η μορφή της μεταβλητής που κωδικοποιεί τον απλότυπο 3

Για να εκτελέσουμε αυτόν τον έλεγχο στο Stata, αφού εισάγουμε τα δεδομένα, πρέπει να πληκτρολογήσουμε την παρακάτω εντολή:

```
xi : logit case i.haplotype_nr [fweight=frequency]
```

Με το `xi` στην αρχή της εντολής δηλώνουμε ότι θέλουμε να χρησιμοποιήσουμε μια κατηγορική μεταβλητή η οποία θα έχει έναν δείκτη. Αυτή θα είναι η `haplotype_nr` όπου βάζοντας μπροστά της το `i.` δηλώνουμε στο Stata ότι αυτή θα είναι η μεταβλητή μας. Πέρα από αυτά με την εντολή `logit` δηλώνουμε πως θέλουμε να κάνουμε λογιστική παλινδρόμηση με εξαρτημένη μεταβλητή την `case` και ανεξάρτητη την `i.haplotype_nr`. Επιπλέον θέλουμε να συμπεριλάβουμε τις συχνότητες των δεδομένων με την παράμετρο `[fweight=frequency]`. Και τα αποτελέσματα που εμφανίζονται στην οθόνη είναι τα εξής:





Εικόνα 19: Το διάστημα εμπιστοσύνης και τα όριά του

Υπάρχουν και άλλοι τρόποι απεικόνισης των αποτελεσμάτων. Όπως ξέρουμε, η παλινδρόμηση χρησιμοποιεί κατά την διάρκεια των υπολογισμών τον λόγο των πιθανοτήτων (σχέση 10). Για κάθε συντελεστή παλινδρόμησης ισχύει ότι  $\exp(\beta_i) = OR$ .

Για παράδειγμα στην συγκεκριμένη μελέτη, ο απλότυπος 1 έχει επιλεγθεί από το Stata ως στοιχείο αναφοράς. Αυτό σημαίνει πως οι επόμενοι απλότυποι θα εξεταστούν ως προς την απόκλισή τους από το στοιχείο αναφοράς, πράγμα που γίνεται με χρήση του πίνακα 5.

	case	control
haplotype 2	1626	2623
haplotype 1	2412	3283

Πίνακας 5: Απλότυποι σε σχέση με ασθένεια

Υπολογίζοντας τον λόγο των πιθανοτήτων του πίνακα προκύπτει ότι:

$$OR = 0.843754$$

Όταν ο λόγος είναι ίσος με την μονάδα, σημαίνει πως τα άτομα που έχουν τον απλότυπο 1, έχουν τον ίδιο κίνδυνο να εμφανίσουν την ασθένεια με τα άτομα που έχουν τον απλότυπο 2. Όταν ο λόγος είναι μεγαλύτερος από την μονάδα, τότε τα άτομα που έχουν τον απλότυπο 2 έχουν μεγαλύτερη πιθανότητα να εμφανίσουν την ασθένεια σε σύγκριση με αυτούς που έχουν τον απλότυπο 1. Τέλος όταν ο λόγος είναι



μικρότερος από την μονάδα, τότε τα άτομα που έχουν τον απλότυπο 2 έχουν μικρότερη πιθανότητα να εμφανίσουν την ασθένεια σε σύγκριση με αυτούς που έχουν τον απλότυπο 1.

Η ίδια διαδικασία θα επαναληφθεί για όλους τους απλότυπους της μελέτης και συνεπώς προκύπτουν και οι αντίστοιχοι λόγοι πιθανοτήτων.

Αξίζει να σημειώσουμε πως κάθε φορά συγκρίνεται κάθε απλότυπος με τον πρώτο, ο οποίος έχει και την μεγαλύτερη συχνότητα, και συνεπώς έχουμε μια εκτίμηση του πόσο αλλάζουν οι πιθανότητες ασθένειας για διαφορετικούς απλότυπους σε σύγκριση με τον πιο κοινό.

Με την ακόλουθη εντολή, το Stata μπορεί να μας εμφανίσει τους σχετικούς λόγους πιθανοτήτων (Odds ratio) που έχει υπολογίσει, κάτι που μπορεί να βοηθήσει στην καλύτερη ερμηνεία των αποτελεσμάτων:

```
xi : logit case i.haplotype_nr [fweight=frequency], or
```

Δεν αλλάζει κάτι στην λειτουργία της λογιστικής παλινδρόμησης, απλά προσθέτουμε σαν επιπλέον επιλογή το `or`.

Το αποτέλεσμα φαίνεται στην εικόνα 20.

```
i.haplotype_nr      _Ihaplotype_1-7      (naturally coded; _Ihaplotype_1 omitted)
Iteration 0:      log likelihood = -8343.008
Iteration 1:      log likelihood = -8318.0226
Iteration 2:      log likelihood = -8317.8958
Iteration 3:      log likelihood = -8317.8958

Logistic regression                                Number of obs   =      12401
                                                    LR chi2(6)      =        50.22
                                                    Prob > chi2     =        0.0000
                                                    Pseudo R2      =        0.0030

Log likelihood = -8317.8958
```

case	Odds Ratio	Std. Err.	z	P> z	[95% Conf. Interval]
_Ihaplotyp~2	.843754	.0349465	-4.10	0.000	.7779664 .9151047
_Ihaplotyp~3	.9148121	.0607533	-1.34	0.180	.8031617 1.041983
_Ihaplotyp~4	.7372685	.0588261	-3.82	0.000	.6305347 .8620698
_Ihaplotyp~5	.4311594	.0702382	-5.16	0.000	.3133087 .5933396
_Ihaplotyp~6	.9357639	.2130904	-0.29	0.771	.5988696 1.462178
_Ihaplotyp~7	.922043	.119135	-0.63	0.530	.7157633 1.187772

**Εικόνα 20: Αποτελέσματα της λογιστικής παλινδρόμησης με Odds Ratio**

Παρατηρούμε ότι η στήλη με τους συντελεστές παλινδρόμησης έχει αντικατασταθεί από τους λόγους πιθανοτήτων.

Μετά την εκτέλεση της παλινδρόμησης, μπορούμε να εκτελέσουμε έναν έλεγχο πάνω στους συντελεστές (coefficients) για να δούμε ποιος από αυτούς

διαφέρει από τους υπόλοιπους, και να έχουμε μια ποσοτική εκτίμηση για την διαφορά αυτή. Ο έλεγχος αυτός έχει τις εξής υποθέσεις:

$$H_0 = \text{Όλοι οι συντελεστές είναι ίσοι με } 0$$

$$H_1 = \text{Κάποιοι ή όλοι οι συντελεστές δεν είναι } 0$$

Για να εκτελέσουμε αυτή τη διαδικασία στο Stata πρέπει να πληκτρολογήσουμε την παρακάτω εντολή:

```
testparm _Ihaplotype_*
```

Έτσι δηλώνουμε ότι θέλουμε να κάνουμε τον έλεγχο αυτό στους συντελεστές της μεταβλητής `_Ihaplotype_*` η οποία προέκυψε από την παλινδρόμηση. Το σύμβολο `*` δηλώνει πως θέλουμε όλα τα περιεχόμενα της μεταβλητής αυτής.

και τα αποτελέσματα φαίνονται στην εικόνα 21.

```
( 1)  _Ihaplotype_2 = 0
( 2)  _Ihaplotype_3 = 0
( 3)  _Ihaplotype_4 = 0
( 4)  _Ihaplotype_5 = 0
( 5)  _Ihaplotype_6 = 0
( 6)  _Ihaplotype_7 = 0

      chi2( 6) =    47.57
      Prob > chi2 =    0.0000
```

**Εικόνα 21: Τα αποτελέσματα του testparm**

Ουσιαστικά συγκρίνουμε όλους τους συντελεστές με τον ίδιο αριθμό (εδώ το 0) και μετράμε την απόστασή τους από αυτόν. Στην περίπτωση αποδοχής της  $H_0$  οι διαφορές των συντελεστών δεν είναι μεγάλες και μπορούμε να συμπεράνουμε ότι δεν έχουμε κάποιον απλότυπο που να επηρεάζει κατά πολύ την ασθένεια. Αντίθετα στην  $H_1$  βλέπουμε ότι υπάρχουν διαφορές, και αυτό μας προτρέπει να ερευνήσουμε ποιος ή ποιοι απλότυποι είναι που διαφέρουν. Επιπλέον υπολογίζεται και ο  $X^2$  του Pearson, που δηλώνει το πόσο ισχυρή είναι η διαφορά και η ισότητα των συντελεστών αντίστοιχα.

Αυτός ο έλεγχος μπορεί να γίνει για διάφορους συνδυασμούς απλότυπων και δεν είναι υποχρεωτικό να έχουμε μια εκτίμηση μόνο για το σύνολό τους.

Παραδείγματος χάρη μπορούμε να ελέγξουμε ένα ζευγάρι απλότυπων, μια τριάδα κλπ. Σε αυτήν την περίπτωση η σύνταξη της εντολής αλλάζει και γίνεται έτσι:

```
test _Ihaplotype_3 _Ihaplotype_7
```

Έτσι ελέγχουμε την ομοιότητα των συντελεστών του απλότυπου 3 και του 7, βάζοντας απλά τον αντίστοιχο αριθμό στην θέση του '\*' που είδαμε πιο πριν.

Γίνεται σαφές ότι αυτός ο επιπλέον έλεγχος είναι απαραίτητος, διότι στα αποτελέσματα της εικόνας 21 βλέπουμε πως στο σύνολό τους οι απλότυποι διαφέρουν όλοι μεταξύ τους. Στην συνέχεια όμως θα δούμε πως ο 3 και ο 7 δεν έχουν στατιστικώς μεγάλη διαφορά (εικόνα 22) καθώς δεχόμαστε την μηδενική υπόθεση.

```
( 1)  _Ihaplotype_3 = 0
( 2)  _Ihaplotype_7 = 0

      chi2( 2) =      2.07
      Prob > chi2 =      0.3560
```

**Εικόνα 22: Αποτελέσματα του test για δυο μεταβλητές**

Βασικό κριτήριο πάντα για την αποδοχή ή απόρριψη της υπόθεσης είναι το p-value που όπως βλέπουμε στον πρώτο έλεγχο (εικόνα 21) είναι μικρότερο από 0.05 και συνεπώς δεν υπάρχει ισότητα, ενώ στον δεύτερο έλεγχο βλέπουμε τελικά ότι με p-value= 0.3560 οι απλότυποι 3 και 7 είναι στατιστικά ισοδύναμοι.

Φυσικά μπορούμε να κάνουμε αυτόν τον έλεγχο με όσες από τις μεταβλητές που έχουν αντίστοιχο συντελεστή παλινδρόμησης επιθυμούμε, πχ:

```
test _Ihaplotype_3 _Ihaplotype_7 _Ihaplotype_6
test _Ihaplotype_2 _Ihaplotype_5 _Ihaplotype_7 _Ihaplotype_4
```

Με αποτέλεσμα την εικόνα 23:

```

( 1)  _Ihaplotype_3 = 0
( 2)  _Ihaplotype_7 = 0
( 3)  _Ihaplotype_6 = 0

      chi2( 3) =    2.11
      Prob > chi2 =  0.5494

```

```

( 1)  _Ihaplotype_2 = 0
( 2)  _Ihaplotype_5 = 0
( 3)  _Ihaplotype_7 = 0
( 4)  _Ihaplotype_4 = 0

      chi2( 4) =   47.53
      Prob > chi2 =  0.0000

```

**Εικόνα 23: αποτελέσματα του test για 3 και 4 μεταβλητές**

#### **2.4.4 Πολυωνυμική λογιστική παλινδρόμηση**

Η πολυωνυμική λογιστική παλινδρόμηση είναι κατάλληλη για χρήση σε κατηγορικές, εξαρτημένες μεταβλητές, χρησιμοποιείται δηλαδή για την ανάλυση των σχέσεων μεταξύ μη – μετρικών εξαρτημένων μεταβλητών με μετρικές ή διχοτομικές ανεξάρτητων μεταβλητών. Συγκρίνει πολλές ομάδες μέσα από ένα συνδιασμό δυαδικών λογιστικών παλινδρομήσεων. Αυτό συμβαίνει επειδή όταν οι κατηγορίες είναι δυο, η διαδικασία είναι ίδια ακριβώς με την δυαδική λογιστική παλινδρόμηση που αναφέρουμε παραπάνω. Το πολυωνυμικό αυτό μοντέλο, είναι μια άμεση επέκταση των λογιστικών μοντέλων, και ο τρόπος υπολογισμού του έχει κάποια κοινά στοιχεία με αυτά (Chen and Kao 2006).

Ας υποθέσουμε ότι μια εξαρτημένη κατηγορική μεταβλητή έχει  $M$  κατηγορίες. Μια από τις τιμές της (συνήθως η πρώτη, η τελευταία, ή αυτή με την μεγαλύτερη συχνότητα) ορίζεται ως η κατηγορία αναφοράς. Αυτή η κατηγορία στην συνέχεια θα αποτελέσει κοινό παρονομαστή για όλες τις υπόλοιπες. Η πιθανότητα της σχέσης με άλλες κατηγορίες, συγκρίνεται με την πιθανότητα σχέσης με την κατηγορία αναφοράς.

Για μια εξαρτημένη μεταβλητή με  $M$  κατηγορίες, απαιτείται ο υπολογισμός  $M-1$  εξισώσεων, μια για κάθε κατηγορία εκτός της κατηγορίας αναφοράς, για να περιγράψουμε την σχέση της εξαρτημένης μεταβλητής με τις ανεξάρτητες.

Έτσι, αν υποθέσουμε ότι για κατηγορία αναφοράς επιλέγεται η πρώτη, τότε για  $m = 2, \dots, M$ , έχουμε:

$$\ln \frac{P(Y_i = m)}{P(Y_i = 1)} = a_m + \sum_{k=1}^K \beta_{mk} X_{ik} = Z_{mi} \quad (19)$$

Συνεπώς για κάθε περίπτωση θα έχουμε M-1 υπολογισμένες λογαριθμικές ποσότητες, μια για κάθε κατηγορία, σε σχέση με την κατηγορία αναφοράς.

Αξίζει να σημειωθεί πως όταν m=1, έχουμε σαν αποτέλεσμα στην σχέση 19 το  $\ln(1) = 0 = Z_{11}$ .

Όταν έχουμε εξαρτημένη μεταβλητή με παραπάνω από δυο κατηγορίες, ο υπολογισμός των πιθανοτήτων γίνεται κάπως πιο πολύπλοκος από την λογιστική παλινδρόμηση. Για m = 2.....M, έχουμε:

$$P(Y_i = m) = \frac{\exp(Z_{mi})}{1 + \sum_{m=2}^M \exp(Z_{mi})} \quad (20)$$

και για την κατηγορία αναφοράς:

$$P(Y_i = 1) = \frac{1}{1 + \sum_{m=2}^M \exp(Z_{mi})} \quad (21)$$

Συνοπτικά η διαδικασία έχει ως εξής: παίρνουμε τους M-1 λογάριθμους που υπολογίσαμε, και τους υψώνουμε ως εκθέτες με βάση τον αριθμό e. Μόλις γίνει αυτό, ο υπολογισμός των πιθανοτήτων είναι εύκολη υπόθεση. Αυτή η μέθοδος βασίζεται στην πιθανότητα επιρροής από παράγοντες, δεδομένης της κατάστασης (ασθενής ή μάρτυρας), όπως εφαρμόζεται στις μελέτες ασθενών και μαρτύρων. Οι απλότυποι χρησιμοποιούνται ως εξαρτημένες μεταβλητές και η κατάσταση (ασθενής/μάρτυρας) λαμβάνεται ως κατηγορία αναφοράς (base outcome) σε αυτήν την μέθοδο (McCullagh 1999). Είναι εύκολο να δούμε ότι οι συντελεστές της παλινδρόμησης ( $\beta_j$ ) είναι προσεγγίσεις του σχετικού λόγου των πιθανοτήτων καθώς  $e^{\beta_j} = OR$ .

Φυσικά ο πρώτος συντελεστής  $\beta_1=0$  στην περίπτωση που χρησιμοποιείται ως κατηγορία αναφοράς. Αυτό το μοντέλο προτάθηκε για την ανάλυση δεδομένων απλότυπων πρώτα από τους (Chen and Kao 2006).

Για να εκτελέσουμε αυτόν τον έλεγχο στο Stata, αφού εισάγουμε τα δεδομένα, πρέπει να πληκτρολογήσουμε την παρακάτω εντολή:

```
xi: mlogit case i.haplotype_nr [fweight=frequency]
```

Με το xi στην αρχή της εντολής δηλώνουμε ότι θέλουμε να χρησιμοποιήσουμε μια κατηγορική μεταβλητή η οποία θα έχει έναν δείκτη. Αυτή θα είναι η haplotype\_nr όπου βάζοντας μπροστά της το i. δηλώνουμε στο Stata ότι αυτή θα είναι η μεταβλητή μας. Πέρα από αυτά με την εντολή mlogit δηλώνουμε πως θέλουμε να κάνουμε πολυωνυμική λογιστική παλινδρόμηση με εξαρτημένη μεταβλητή την case και ανεξάρτητη την i.haplotype\_nr. Επιπλέον θέλουμε να συμπεριλάβουμε τις συχνότητες των δεδομένων με την παράμετρο [fweight=frequency].

Και τα αποτελέσματα εμφανίζονται στην εικόνα 24.

```
i.haplotype_nr    _Ihaplotype_1-7    (naturally coded; _Ihaplotype_1 omitted)
Iteration 0:    log likelihood = -8343.008
Iteration 1:    log likelihood = -8318.0226
Iteration 2:    log likelihood = -8317.8958
Iteration 3:    log likelihood = -8317.8958

Multinomial logistic regression                Number of obs =      12401
                                                LR chi2(6)      =       50.22
                                                Prob > chi2     =       0.0000
Log likelihood = -8317.8958                    Pseudo R2      =       0.0030
```

case	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
1					
_Ihaplotyp~2	-.1698943	.0414179	-4.10	0.000	-.2510719 -.0887167
_Ihaplotyp~3	-.0890366	.0664107	-1.34	0.180	-.2191992 .0411259
_Ihaplotyp~4	-.3048031	.0797892	-3.82	0.000	-.4611872 -.1484191
_Ihaplotyp~5	-.8412774	.1629055	-5.16	0.000	-1.160566 -.5219884
_Ihaplotyp~6	-.0663921	.2277182	-0.29	0.771	-.5127115 .3799273
_Ihaplotyp~7	-.0811634	.1292076	-0.63	0.530	-.3344057 .1720789
_cons	-.3083014	.0268178	-11.50	0.000	-.3608632 -.2557395

(case==0 is the base outcome)

**Εικόνα 24: Αποτελέσματα της πολυωνυμικής παλινδρόμησης**

Παρατηρούμε ότι τα αποτελέσματα είναι ίδια με αυτά της προηγούμενης μεθόδου. Αυτό είναι αναμενόμενο καθώς υπάρχουν πολλές ομοιότητες.



Και σε αυτήν την περίπτωση μπορούμε να απεικονίσουμε τους σχετικούς λόγους των πιθανοτήτων με την παρακάτω εντολή:

```
xi: mlogit case i.haplotype_nr [fweight=frequency], rrr
```

Με το rrr (relative risk ratio) σαν πρόσθετη επιλογή, δηλαδή ο λόγος σχετικού κινδύνου, έχουμε το αποτέλεσμα της εικόνας 25.

```
i.haplotype_nr    _Ihaplotype_1-7    (naturally coded; _Ihaplotype_1 omitted)
Iteration 0:    log likelihood = -8343.008
Iteration 1:    log likelihood = -8318.0226
Iteration 2:    log likelihood = -8317.8958
Iteration 3:    log likelihood = -8317.8958

Multinomial logistic regression                Number of obs =      12401
LR chi2(6) = 50.22
Prob > chi2 = 0.0000
Pseudo R2 = 0.0030

Log likelihood = -8317.8958
```

case	RRR	Std. Err.	z	P> z	[95% Conf. Interval]	
1						
_Ihaplotyp~2	.843754	.0349465	-4.10	0.000	.7779664	.9151047
_Ihaplotyp~3	.9148121	.0607533	-1.34	0.180	.8031617	1.041983
_Ihaplotyp~4	.7372685	.0588261	-3.82	0.000	.6305347	.8620698
_Ihaplotyp~5	.4311594	.0702382	-5.16	0.000	.3133087	.5933396
_Ihaplotyp~6	.9357639	.2130904	-0.29	0.771	.5988696	1.462178
_Ihaplotyp~7	.922043	.119135	-0.63	0.530	.7157633	1.187772

(case==0 is the base outcome)

**Εικόνα 25: Αποτελέσματα πολωνομικής παλινδρόμησης με relative risk ratio**

Και πάλι βλέπουμε την στήλη με τους συντελεστές να γίνεται λόγος σχετικού κινδύνου για κάθε απλότυπο.

Στη συνέχεια μπορούμε να εκτελέσουμε τον έλεγχο ισότητας των συντελεστών με την παρακάτω εντολή:

```
test [1]
```

Και τα αποτελέσματα στην εικόνα 26 είναι ακριβώς ίδια με της λογιστικής παλινδρόμησης, απλώς αντί για την εντολή testparm χρησιμοποιούμε την test [1] και δεν είναι απαραίτητο να συμπεριλάβουμε την μεταβλητή με τους συντελεστές στην σύνταξη.

```

( 1)  _Ihaplotype_2 = 0
( 2)  _Ihaplotype_3 = 0
( 3)  _Ihaplotype_4 = 0
( 4)  _Ihaplotype_5 = 0
( 5)  _Ihaplotype_6 = 0
( 6)  _Ihaplotype_7 = 0

      chi2( 6) =    47.57
      Prob > chi2 =    0.0000

```

**Εικόνα 26: Αποτελέσματα του test [1]**

Επίσης μπορούμε όπως αναφερθήκαμε και παραπάνω να εκτελέσουμε αυτόν τον έλεγχο όχι μόνο για το σύνολο, αλλά και για κάποιες ομάδες απλότυπων. Η μόνη διαφορά φαίνεται στην σύνταξη της εντολής, όπου βάζουμε άνω και κάτω τελεία πριν τις μεταβλητές που επιθυμούμε, σε αντίθεση με την λογιστική παλινδρόμηση:

```
test [1]: _Ihaplotype_6 _Ihaplotype_7
```

```

( 1)  [1]_Ihaplotype_6 = 0
( 2)  [1]_Ihaplotype_7 = 0

      chi2( 2) =    0.47
      Prob > chi2 =    0.7902

```

**Εικόνα 27: test[1] για δυο μεταβλητές**

Βλέπουμε ότι ο απλότυπος 7 είναι ίδιος με τον 6 (εικόνα 27), και μπορούμε να πούμε ότι οι απλότυποι 6 και 7 επηρεάζουν στην ίδια περίπου ένταση το αποτέλεσμα της παλινδρόμησης και κατ' επέκταση την ασθένεια. Όσο αυξάνεται ο αριθμός των ελέγχων που κάνουμε, τόσο καλύτερη εικόνα έχουμε για το τι συμβαίνει.

Φυσικά μπορούμε να κάνουμε αυτόν τον έλεγχο με όσες από τις μεταβλητές που έχουν αντίστοιχο συντελεστή παλινδρόμησης επιθυμούμε, πχ:

```

test [1]: _Ihaplotype_3 _Ihaplotype_7 _Ihaplotype_6
test [1]: _Ihaplotype_2 _Ihaplotype_5 _Ihaplotype_7 _Ihaplotype_4

```

Με τα ακόλουθα αποτελέσματα (εικόνα 28):



```
( 1) [1]_Ihaplotype_3 = 0
( 2) [1]_Ihaplotype_7 = 0
( 3) [1]_Ihaplotype_6 = 0

      chi2( 3) =    2.11
      Prob > chi2 = 0.5494
```

```
( 1) [1]_Ihaplotype_2 = 0
( 2) [1]_Ihaplotype_5 = 0
( 3) [1]_Ihaplotype_7 = 0
( 4) [1]_Ihaplotype_4 = 0

      chi2( 4) =   47.53
      Prob > chi2 = 0.0000
```

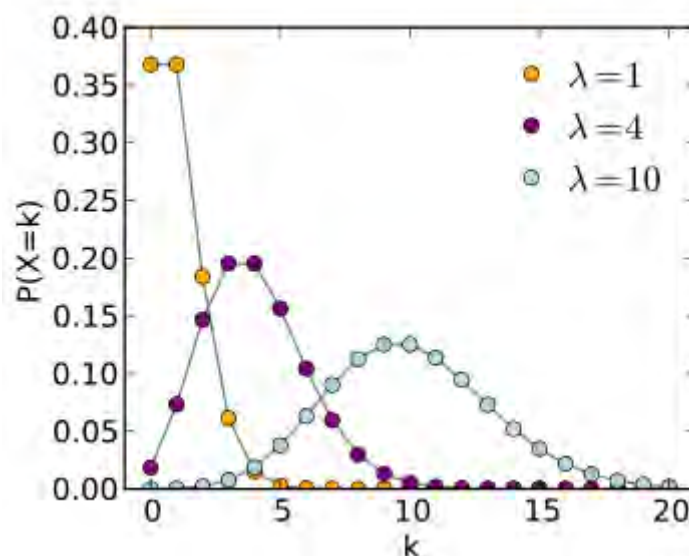
**Εικόνα 28:** Αποτελέσματα του test [1] για 3 και 4 μεταβλητές

## 2.4.5 Παλινδρόμηση Poisson

Αυτού του είδους η παλινδρόμηση σχετίζεται με την κατανομή Poisson (εικόνα 29), η οποία είναι μια κατανομή που εκφράζει την πιθανότητα εμφάνισης ενός αριθμού γεγονότων σε συγκεκριμένο χρόνο. Τα γεγονότα αυτά είναι ανεξάρτητα το ένα με το άλλο. Αν ο αναμενόμενος αριθμός των γεγονότων είναι  $\lambda$ , τότε η πιθανότητα να έχουμε ακριβώς  $n$  γεγονότα ( $n = 0, 1, 2, \dots$ ) είναι:

$$f(n, \lambda) = \frac{\lambda^n e^{-\lambda}}{n!}$$

(22)



Εικόνα 29: Η συνάρτηση μάζας πιθανότητας της κατανομής Poisson

Η παλινδρόμηση Poisson είναι μια μορφή παλινδρόμησης που χρησιμοποιείται για να μοντελοποιήσει απαριθμήσιμα δεδομένα και πίνακες ενδεχομένων. Απαριθμήσιμα δεδομένα, είναι τα δεδομένα στα οποία οι παρατηρήσεις παίρνουν μη-μηδενικές ακέραιες τιμές, και προκύπτουν από μέτρημα, όχι κατηγοριοποίηση. Οι πίνακες ενδεχομένων χρησιμοποιούνται για να καταγράψουν και να αναλύσουν τη σχέση μεταξύ δυο ή περισσότερων μεταβλητών.

Το μοντέλο χρησιμοποιεί μια μεταβλητή  $Y$  η οποία ακολουθεί κατανομή Poisson, καθώς και τον λογάριθμο της αναμενόμενης τιμής της μεταβλητής  $Y$ . Τα παραπάνω στοιχεία μπορούν να μοντελοποιηθούν από έναν γραμμικό συνδυασμό από

άγνωστες παραμέτρους. Το μοντέλο αυτό ονομάζεται και λογαριθμικό-γραμμικό (log-linear) μοντέλο, ειδικά όταν χρησιμοποιείται με πίνακες ενδεχομένων.

Η πιο απλή περίπτωση με μια ανεξάρτητη μεταβλητή  $x$ , έχει την εξής μορφή:

$$\log(E(Y)) = a + bx \quad (23)$$

Αν  $y_i$  είναι ανεξάρτητες παρατηρήσεις με αντίστοιχες τιμές  $x_i$  για την μεταβλητή  $x$  τότε τα  $a$  και  $b$  μπορούν να υπολογιστούν από την μέγιστη πιθανοφάνεια αν ο αριθμός των ευδιάκριτων  $x$  είναι τουλάχιστον 2.

Η παλινδρόμηση Poisson είναι κατάλληλη για διαβαθμιζόμενα δεδομένα, όπου η διαβάθμιση είναι ένα μέτρημα γεγονότων που λαμβάνουν χώρα σε μια συγκεκριμένη μονάδα παρατηρήσεων, διαιρούμενη από μια μονάδα μέτρησης της «έκθεσης» αυτής (exposure). Για παράδειγμα οι βιολόγοι μετρούν τον αριθμό των ειδών των δέντρων σε ένα δάσος όπου η διαβάθμιση θα είναι ο αριθμός των ειδών ανά τετραγωνικό χιλιόμετρο ή αυτοί που ασχολούνται με την δημογραφική επιστήμη μπορούν να μοντελοποιούν ρυθμούς θανάτων ανά γεωγραφικές περιοχές, ως τον αριθμό των θανάτων διαιρούμενο από άτομα/χρόνο.

Γενικότερα οι ρυθμοί των γεγονότων μπορούν να υπολογιστούν ως γεγονότα ανά μονάδα χρόνου. Στα παραπάνω παραδείγματα η «έκθεση» είναι αντίστοιχα το τετραγωνικό χιλιόμετρο, τα άτομα/χρόνο και ο χρόνος. Στην παλινδρόμηση Poisson αυτό αντιμετωπίζεται ως ένα αντιστάθμισμα, όπου η μεταβλητή της «έκθεσης» εισάγεται στην δεξιά πλευρά της εξίσωσης και με μια παραμετροποίηση στην λογαριθμική ποσότητα ( $\log(\text{exposure})$ ) περιορίζεται μέχρι το 1.

$$\log(E(Y)) = \log(\text{exposure}) + a + bx \quad (24)$$

το οποίο γίνεται:

$$\log(E(Y)) - \log(\text{exposure}) = \log\left(\frac{E(Y)}{\text{exposure}}\right) = a + bx \quad (25)$$

Ένα χαρακτηριστικό της κατανομής Poisson είναι ότι η μέση τιμή της είναι ίση με την διακύμανση. Κάποιες φορές όμως παρατηρούμε ότι η διακύμανση είναι μεγαλύτερη από την μέση τιμή. Κάτι τέτοιο είναι γνωστό με τον όρο υπερδιασπορά (Berk and MacDonald 2007) και δηλώνει ότι το μοντέλο δεν είναι κατάλληλο για τα δεδομένα. Συχνά ο λόγος που εμφανίζεται η υπερδιασπορά είναι η παράλειψη κάποιων επεξηγηματικών μεταβλητών.

Άλλο ένα γνωστό πρόβλημα με την παλινδρόμηση Poisson είναι η παρουσία μεγάλου αριθμού μηδενικών. Αν εκτελούνται δυο διεργασίες, μια να αποφασίζει αν υπάρχουν μηδέν γεγονότα ή αν υπάρχουν, και μια άλλη Poisson διεργασία να υπολογίζει τον αριθμό των γεγονότων, τότε θα υπάρχουν παραπάνω μηδενικά από όσα μπορεί να προβλέψει μια παλινδρόμηση Poisson. Ένα παράδειγμα για αυτό θα ήταν η κατανομή των τσιγάρων που καπνίζονται μέσα σε μια ώρα από τα μέλη μιας ομάδας ατόμων, όπου στην ομάδα αυτή υπάρχουν και μη-καπνιστές.

Αυτή η μέθοδος δουλεύει με την υπόθεση ότι οι παρατηρούμενες συχνότητες των απλότυπων αντιπροσωπεύουν μια τυχαία μεταβλητή η οποία ακολουθεί την κατανομή Poisson. Έτσι μπορούμε να εκτελέσουμε αυτό το μοντέλο, βάζοντας σαν εξαρτημένη μεταβλητή τις συχνότητες, και σαν παράγοντα επιρροής (ανεξάρτητες μεταβλητές) τους απλότυπους, αλλά και την κατάσταση (case/control).

Ιστορικά αυτά τα λεγόμενα log-linear μοντέλα έχουν χρησιμοποιηθεί αρκετά νωρίς για ανάλυση απλότυπων, για παράδειγμα στον υπολογισμό του συνελεστή συσχέτισης (LD), (Weir and Wilson 1986) αλλά και σε μελέτες γενετικής συσχέτισης ασθενειών με απλότυπους (Tiret, Amouyel et al. 1991). Το μοντέλο αυτό μπορεί να διατυπωθεί με τον εξής τύπο:

$$\log(n_j) = \mu_0 + \beta_0 y_j + \sum_{j=2}^r a_j z_j + \sum_{j=2}^r \beta_j z_j y_j \quad (26)$$

Το οποίο είναι ένα μοντέλο για την περιγραφή των 2 x r πινάκων συνάφειας απλότυπων σε σχέση με την ασθένεια. Τα  $\beta_j$  είναι οι συντελεστές που αντιστοιχούν στην αλληλεπίδραση (interaction) απλότυπου-ασθένειας και είναι όμοια με αυτά της

λογιστικής και της πολυωνυμικής παλινδρόμησης. Εδώ βγαίνει εύκολα το συμπέρασμα πως οι συντελεστές  $a_j$  και  $\beta_j$  είναι όμοιοι και στα τρία μοντέλα.

Για να εκτελέσουμε αυτόν τον έλεγχο στο Stata, αφού εισάγουμε τα δεδομένα, πρέπει να πληκτρολογήσουμε την παρακάτω εντολή:

```
xi: poisson frequency i.case*i.haplotype_nr
```

Με το xi στην αρχή της εντολής δηλώνουμε ότι θέλουμε να χρησιμοποιήσουμε μια κατηγορική μεταβλητή η οποία θα έχει έναν δείκτη. Σε αυτή την περίπτωση θα είναι το γινόμενο των μεταβλητών case και haplotype\_nr όπου βάζοντας μπροστά τους το i. δηλώνουμε στο Stata ότι είναι οι μεταβλητές μας. Πέρα από αυτά με την εντολή poisson δηλώνουμε πως θέλουμε να κάνουμε παλινδρόμηση Poisson με εξαρτημένη μεταβλητή την frequency και ανεξάρτητη την i.case\*i.haplotype\_nr.

Και τα αποτελέσματα εμφανίζονται στην εικόνα 30:

```
i.case          _Icase_0-1          (naturally coded; _Icase_0 omitted)
i.haplotype_nr  _Ihaplotype_1-7        (naturally coded; _Ihaplotype_1 omitted)
i.case*i.hapl~r  _Icasxhap_#_#          (coded as above)
```

```
Iteration 0:  log likelihood = -6133.2732
Iteration 1:  log likelihood = -5823.1098
Iteration 2:  log likelihood = -5820.7762
Iteration 3:  log likelihood = -5820.7737
Iteration 4:  log likelihood = -5820.7737
```

```
Poisson regression          Number of obs   =          84
                           LR chi2(13)           =       16592.09
                           Prob > chi2          =          0.0000
                           Pseudo R2           =          0.5877
```

```
Log likelihood = -5820.7737
```

frequency	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
_Icase_1	-.3083014	.0268178	-11.50	0.000	-.3608632 - .2557395
_Ihaplotyp~2	-.2244389	.0261886	-8.57	0.000	-.2757676 - .1731102
_Ihaplotyp~3	-1.583283	.042288	-37.44	0.000	-1.666166 -1.5004
_Ihaplotyp~4	-1.873937	.0478406	-39.17	0.000	-1.967703 -1.780171
_Ihaplotyp~5	-3.015109	.0807204	-37.35	0.000	-3.173318 -2.8569
_Ihaplotyp~6	-4.225312	.1453889	-29.06	0.000	-4.510269 -3.940355
_Ihaplotyp~7	-3.053088	.0821962	-37.14	0.000	-3.214189 -2.891986
_Icasxhap~2	-.1698943	.0414179	-4.10	0.000	-.2510719 -.0887167
_Icasxhap~3	-.0890366	.0664107	-1.34	0.180	-.2191992 .0411259
_Icasxhap~4	-.3048031	.0797892	-3.82	0.000	-.4611872 -.1484191
_Icasxhap~5	-.8412774	.1629055	-5.16	0.000	-1.160566 -.5219884
_Icasxhap~6	-.0663921	.2277182	-0.29	0.771	-.5127115 .3799273
_Icasxhap~7	-.0811634	.1292076	-0.63	0.530	-.3344057 .1720789
_cons	6.304753	.0174528	361.25	0.000	6.270547 6.33896

Εικόνα 30: Αποτελέσματα της παλινδρόμησης poisson

Τα τελικά αποτελέσματα που λαμβάνουμε υπόψη είναι οι μεταβλητές `_IcasXhap_~i`,  $i=2\dots 7$  οι οποίες φαίνονται στις γραμμές 8-13, στην στήλη frequency του πίνακα της εικόνας 30.

Όπως και στις προηγούμενες μεθόδους, μπορούμε να εμφανίσουμε τα αποτελέσματα με την στήλη των λόγων των πιθανοτήτων κάθε απλότυπου:

```
xi: poisson frequency i.case*i.haplotype_nr, irr
```

Όπου εδώ το `irr` σημαίνει incidence-rate ratio (σχετικός λόγος πιθανοτήτων), και έχουμε το αποτέλεσμα της εικόνας 31.

```
i.case          _Icase_0-1          (naturally coded; _Icase_0 omitted)
i.haplotype_nr  _Ihaplotype_1-7        (naturally coded; _Ihaplotype_1 omitted)
i.case*i.hapl~r  _IcasXhap_#_#    (coded as above)
```

```
Iteration 0:    log likelihood = -6133.2732
Iteration 1:    log likelihood = -5823.1098
Iteration 2:    log likelihood = -5820.7762
Iteration 3:    log likelihood = -5820.7737
Iteration 4:    log likelihood = -5820.7737
```

```
Poisson regression                                Number of obs   =           84
                                                    LR chi2(13)     =       16592.09
                                                    Prob > chi2     =           0.0000
                                                    Pseudo R2      =           0.5877
```

```
Log likelihood = -5820.7737
```

frequency	IRR	Std. Err.	z	P> z	[95% Conf. Interval]
_Icase_1	.7346939	.0197029	-11.50	0.000	.6970743 .7743437
_Ihaplotyp~2	.7989644	.0209238	-8.57	0.000	.7589893 .8410449
_Ihaplotyp~3	.2053	.0086817	-37.44	0.000	.1889702 .223041
_Ihaplotyp~4	.1535181	.0073444	-39.17	0.000	.1397776 .1686094
_Ihaplotyp~5	.0490405	.0039586	-37.35	0.000	.0418645 .0574466
_Ihaplotyp~6	.0146208	.0021257	-29.06	0.000	.0109955 .0194413
_Ihaplotyp~7	.0472129	.0038807	-37.14	0.000	.0401879 .0554659
_IcasXhap~2	.843754	.0349465	-4.10	0.000	.7779664 .9151047
_IcasXhap~3	.9148121	.0607533	-1.34	0.180	.8031617 1.041983
_IcasXhap~4	.7372685	.0588261	-3.82	0.000	.6305347 .8620698
_IcasXhap~5	.4311594	.0702383	-5.16	0.000	.3133087 .5933396
_IcasXhap~6	.9357639	.2130904	-0.29	0.771	.5988696 1.462178
_IcasXhap~7	.922043	.119135	-0.63	0.530	.7157633 1.187772

**Εικόνα 31: Αποτελέσματα της παλινδρόμησης poisson με incidence-rate ratio**

Για την εκτέλεση του ελέγχου ισότητας των συντελεστών ισχύει η εντολή που χρησιμοποιούμε στην λογιστική παλινδρόμηση, απλά την εφαρμόζουμε στην μεταβλητή `_IcasXhap_1_i`,  $i=2\dots 7$ .

```
testparm _IcasXhap_1_*
```

Όπου το σύμβολο `'*'` δηλώνει πως θέλουμε όλα τα περιεχόμενα της μεταβλητής `_casXhap_1_`.



Και το αποτέλεσμα είναι το αναμενόμενο (εικόνα 32):

```
( 1) [frequency]_IcasXhap_1_2 = 0
( 2) [frequency]_IcasXhap_1_3 = 0
( 3) [frequency]_IcasXhap_1_4 = 0
( 4) [frequency]_IcasXhap_1_5 = 0
( 5) [frequency]_IcasXhap_1_6 = 0
( 6) [frequency]_IcasXhap_1_7 = 0

      chi2( 6) =    47.57
      Prob > chi2 =    0.0000
```

**Εικόνα 32: Αποτελέσματα testparm για παλινδρόμηση poisson**

Όπως και με τις άλλες δυο μεθόδους, μπορούμε να κάνουμε αυτόν τον έλεγχο για διάφορους συνδιασμούς των συντελεστών, για παράδειγμα:

```
test _IcasXhap_1_3 _IcasXhap_1_7
```

με αποτέλεσμα την εικόνα 33.

```
( 1) [frequency]_IcasXhap_1_3 = 0
( 2) [frequency]_IcasXhap_1_7 = 0

      chi2( 2) =    2.07
      Prob > chi2 =    0.3560
```

**Εικόνα 33: Αποτέλεσμα του test για δυο μεταβλητές**

Φυσικά μπορούμε να κάνουμε αυτόν τον έλεγχο με όσες από τις μεταβλητές που έχουν αντίστοιχο συντελεστή παλινδρόμησης επιθυμούμε, πχ:

```
test _IcasXhap_1_3 _IcasXhap_1_6 _IcasXhap_1_7
test _IcasXhap_1_2 _IcasXhap_1_4 _IcasXhap_1_5 _IcasXhap_1_7
```

Άσχετα με το ποια είναι η στρατηγική δειγματοληψίας που δημιούργησε τα δεδομένα μας, είναι ευρέως γνωστό ότι τα αποτελέσματα της προσαρμογής αυτών των μοντέλων στα δεδομένα είναι σχεδόν όμοια (Agresti 2002). Επομένως η επιλογή του μοντέλου πρέπει να βασίζεται σε παράγοντες που έχουν να κάνουν με την αληθοφάνεια και την ερμηνευτικότητα των αποτελεσμάτων. Να μπορούν δηλαδή να αναλυθούν εύκολα, και να παράγουν χρήσιμα και ξεκάθαρα αποτελέσματα. Επίσης ένας παράγοντας που πρέπει να λάβουμε υπόψη είναι το πόσο εύκαμπτο μπορεί να

είναι το μοντέλο και να επιτρέπει αλλαγές στα δεδομένα ή τις παραμέτρους του. Έχει αποδειχτεί ότι οι υπολογισμοί μέγιστης πιθανοφάνειας που αποκτούμε από τα μοντέλα με την μέθοδο της αναμενόμενης πιθανοφάνειας (prospective likelihood) είναι ίδιοι με αυτούς που αποκτώνται με την μέθοδο της αναδρομικής πιθανοφάνειας (retrospective likelihood) (Prentice and Pyke 1979; Chen 2003). Η ισότητα της λογιστικής παλινδρόμησης και των μοντέλων Poisson έχουν επίσης χρησιμοποιηθεί στο παρελθόν για την δημιουργία μεθόδων για τον εντοπισμό της αλληλεπίδρασης γονιδίων σε βιολογικό επίπεδο, (Umbach and Weinberg 1997).

Οι παραπάνω μέθοδοι είναι απλές προεκτάσεις των γενικευμένων γραμμικών μοντέλων και υποθέτουν πως α) Η απλοτυπική πιθανότητα ακολουθεί ένα αθροιστικό μοντέλο, β) Ότι η φάση του απλότυπου είναι γνωστή (βλέπε εισαγωγή – το πρόβλημα εύρεσης του απλότυπου (Lin, Cutler et al. 2002)) και γ) ότι ο πληθυσμός είναι σε ισορροπία Hardy-Weinberg. Το γενετικό μοντέλο της κληρονομικότητας μπορεί να λειτουργήσει με την χρήση των «απλο-γονότυπων» αντί για τους γονότυπους στην ανάλυση. Αυτό πραγματοποιείται εύκολα με τις μεθόδους που παρουσιάστηκαν χρησιμοποιώντας όλους τους συνδυασμούς των απλότυπων ( $h_1h_1, h_1h_2$ , κλπ). Ωστόσο στις μελέτες συσχέτισης με ασθενείς και μάρτυρες που χρησιμοποιούμε άτομα χωρίς συγγένεια, για να αντιμετωπίσουμε το πρόβλημα της εύρεσης του απλότυπου, χρησιμοποιούμε στατιστικές μεθόδους, συνήθως με χρήση του αλγορίθμου EM, ή κάποιον παρεμφερή (Niu 2004; Xu, Wu et al. 2004; Marchini, Cutler et al. 2006). Το να χρησιμοποιούμε τους απλότυπους σαν γνωστούς ίσως αποβεί προβληματικό (Lin and Huang 2007), επειδή όταν χρησιμοποιούμε πιθανοτικές μεθόδους, ο πιο πιθανός απλότυπος, δεν είναι απαραίτητα και ο σωστός.

Για την αντιμετώπιση αυτών των προβλημάτων έχουν αναπτυχθεί διάφορες μέθοδοι κατά καιρούς, για παράδειγμα η πρόσθεση κάποιου βάρους στους απλότυπους σύμφωνα με την πιθανότητα εμφάνισής τους (Zaykin, Westfall et al. 2002; French, Lumley et al. 2006). Η διαδικασία της παλινδρόμησης δεν αλλάζει, απλά προσθέτουμε επιπλέον μια στήλη με όνομα prob στον πίνακα 4, η οποία περιέχει τις αντίστοιχες πιθανότητες για κάθε απλότυπο.

Για παράδειγμα μετά την εισαγωγή των δεδομένων στο Stata εκτελούμε τις παρακάτω εντολές:

```
expand frequency
```



Αυτή η εντολή επεκτείνει τα δεδομένα έτσι ώστε να μπορέσουμε να εφαρμόσουμε στην συνέχεια την λογιστική παλινδρόμηση αλλά με πιθανοτικά βάρη. Αυτό γίνεται με την εντολή:

```
xi: logit case i.haplotype_nr [pweight=prob]
```

Με το [pweight=prob] να δηλώνει πως θα χρησιμοποιήσουμε πιθανοτικά βάρη, θα λάβουμε δηλαδή υπόψιν την αλληλεπίδραση των απλότυπων με βάση την πιθανότητα που έχουν να είναι αληθείς. Οι τιμές των συντελεστών παλινδρόμησης βλέπουμε πως έχουν επηρεαστεί, αν συγκρίνουμε τις εικόνες 18 και 34.

```
i.haplotype_nr    _Ihaplotype_1-7    (naturally coded; _Ihaplotype_1 omitted)
(sum of wgt is    9.3215e+03)
Iteration 0:      log pseudolikelihood = -8371.4204
Iteration 1:      log pseudolikelihood = -8354.2923
Iteration 2:      log pseudolikelihood = -8354.2368
Iteration 3:      log pseudolikelihood = -8354.2367

Logistic regression              Number of obs    =    12406
                                Wald chi2(6)     =     47.56
                                Prob > chi2        =     0.0000
Log pseudolikelihood = -8354.2367 Pseudo R2       =     0.0021
```

case	Coef.	Robust Std. Err.	z	P> z	[95% Conf. Interval]	
_Ihaplotype~2	-.1698943	.0414196	-4.10	0.000	-.2510752	-.0887135
_Ihaplotype~3	-.0890366	.0664134	-1.34	0.180	-.2192044	.0411312
_Ihaplotype~4	-.3048031	.0797925	-3.82	0.000	-.4611935	-.1484128
_Ihaplotype~5	-.8412774	.1629121	-5.16	0.000	-1.160579	-.5219756
_Ihaplotype~6	-.0773611	.2239085	-0.35	0.730	-.5162136	.3614914
_Ihaplotype~7	-.0781156	.1287041	-0.61	0.544	-.330371	.1741399
_cons	-.3083014	.0268189	-11.50	0.000	-.3608654	-.2557374

**Εικόνα 34: Αποτελέσματα της λογιστικής παλινδρόμησης με πιθανοτικά βάρη**

Ένας άλλος τρόπος θα ήταν η περίπτωση της εντολής `hapirpf` η οποία δέχεται γονοτυπικά δεδομένα και υπολογίζει αυτόματα την απλοτυπική φάση (Mander 2001).

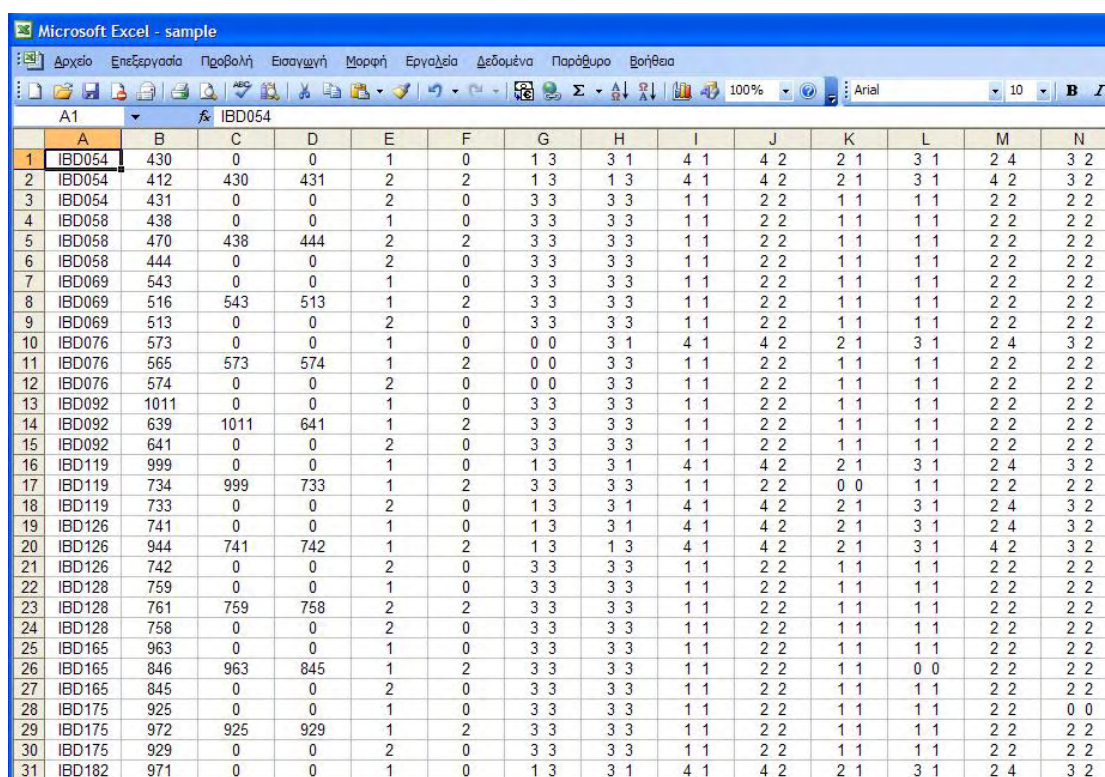
Επίσης έχουν αναπτυχθεί μέθοδοι τύπου score για αναμενόμενη πιθανότητα (Schaid, Rowland et al. 2002) και για αναδρομική πιθανότητα (Epstein and Satten 2003), καθώς και μεθόδους που αφορούν την αλληλεπίδραση των γονιδίων (Chen and Li 2008). Μια σύγκριση των μεθόδων έδειξε ότι τα αποτελέσματα είναι δύσκολο να συγκριθούν όταν η επιρροή του απλότυπου στην πιθανότητα της ασθένειας ακολουθεί ένα πολυπαραγοντικό μοντέλο. Ωστόσο για κυριαρχικά και υποχωρητικά μοντέλα, η

μέθοδος αναδρομικής πιθανότητας έχει αυξημένη αποτελεσματικότητα, χωρίς να παραμελούμε την αναμενόμενη πιθανότητα (Satten and Epstein 2004).

Επίσης γραφικά μοντέλα έχουν προταθεί από τον Thomas (Thomas 2005), τα οποία επιδιώκουν την ανάλυση γονοτυπικών δεδομένων αντί για απλοειδή, και λογαριθμικά-γραμμικά μοντέλα από τον Baker (Baker 2005) ο οποίος συμπεριλαμβάνει κάποιους περιβαλλοντικούς συντελεστές οι οποίοι συμβάλουν στο αποτέλεσμα και λαμβάνει επίσης υπόψη τυχόν αποκλίσεις από την ισορροπία Hardy-Weinberg. Ο Lin και συνεργάτες έχουν επεκτείνει τις μεθόδους που παρουσιάστηκαν παραπάνω, με το να συμπεριλάβουν διάφορα είδη δειγματοληψίας σε ένα ενοποιημένο πλαίσιο εργασίας (Lin, Zeng et al. 2005). Εφάρμοσαν τις μεθόδους τους σε μια μελέτη για τον καρκίνο του μαστού η οποία έδειξε κάποιες σημαντικές επιπλοκές του καπνίσματος σε απλότυπους σχετικά με τον καρκίνο αυτού του τύπου.

## 2.4.6 Η εφαρμογή Harpview

Η εφαρμογή Harpview είναι ένα συχνά χρησιμοποιούμενο εργαλείο σε θέματα βιοπληροφορικής, το οποίο έχει δημιουργηθεί για την ανάλυση και την οπτικοποίηση ακολουθιών, αλλά και υπολογισμό της ανισοροπίας σύνδεσης (linkage disequilibrium) σε γενετικά δεδομένα. Επίσης είναι κατάλληλο για μελέτες γενετικής συσχέτισης, για την διαλογή κάποιων πολυμορφισμών και τέλος για τον υπολογισμό της συχνότητας των απλότυπων. Η εφαρμογή δημιουργήθηκε και εξελίσσεται από το εργαστήριο του Dr.Mark Daly στο ινστιτούτο του MIT/Harvard Broad Institute (<http://www.broadinstitute.org>).



The screenshot shows a Microsoft Excel spreadsheet with a table of linkage disequilibrium data. The table has 14 columns labeled A through N and 31 rows labeled 1 through 31. Each row represents a SNP, and the columns represent the pairwise LD values between that SNP and the others. The diagonal elements are all 1.0. The table is as follows:

	A	B	C	D	E	F	G	H	I	J	K	L	M	N
1	IBD054	430	0	0	1	0	1 3	3 1	4 1	4 2	2 1	3 1	2 4	3 2
2	IBD054	412	430	431	2	2	1 3	1 3	4 1	4 2	2 1	3 1	4 2	3 2
3	IBD054	431	0	0	2	0	3 3	3 3	1 1	2 2	1 1	1 1	2 2	2 2
4	IBD058	438	0	0	1	0	3 3	3 3	1 1	2 2	1 1	1 1	2 2	2 2
5	IBD058	470	438	444	2	2	3 3	3 3	1 1	2 2	1 1	1 1	2 2	2 2
6	IBD058	444	0	0	2	0	3 3	3 3	1 1	2 2	1 1	1 1	2 2	2 2
7	IBD069	543	0	0	1	0	3 3	3 3	1 1	2 2	1 1	1 1	2 2	2 2
8	IBD069	516	543	513	1	2	3 3	3 3	1 1	2 2	1 1	1 1	2 2	2 2
9	IBD069	513	0	0	2	0	3 3	3 3	1 1	2 2	1 1	1 1	2 2	2 2
10	IBD076	573	0	0	1	0	0 0	3 1	4 1	4 2	2 1	3 1	2 4	3 2
11	IBD076	565	573	574	1	2	0 0	3 3	1 1	2 2	1 1	1 1	2 2	2 2
12	IBD076	574	0	0	2	0	0 0	3 3	1 1	2 2	1 1	1 1	2 2	2 2
13	IBD092	1011	0	0	1	0	3 3	3 3	1 1	2 2	1 1	1 1	2 2	2 2
14	IBD092	639	1011	641	1	2	3 3	3 3	1 1	2 2	1 1	1 1	2 2	2 2
15	IBD092	641	0	0	2	0	3 3	3 3	1 1	2 2	1 1	1 1	2 2	2 2
16	IBD119	999	0	0	1	0	1 3	3 1	4 1	4 2	2 1	3 1	2 4	3 2
17	IBD119	734	999	733	1	2	3 3	3 3	1 1	2 2	0 0	1 1	2 2	2 2
18	IBD119	733	0	0	2	0	1 3	3 1	4 1	4 2	2 1	3 1	2 4	3 2
19	IBD126	741	0	0	1	0	1 3	3 1	4 1	4 2	2 1	3 1	2 4	3 2
20	IBD126	944	741	742	1	2	1 3	1 3	4 1	4 2	2 1	3 1	4 2	3 2
21	IBD126	742	0	0	2	0	3 3	3 3	1 1	2 2	1 1	1 1	2 2	2 2
22	IBD128	759	0	0	1	0	3 3	3 3	1 1	2 2	1 1	1 1	2 2	2 2
23	IBD128	761	759	758	2	2	3 3	3 3	1 1	2 2	1 1	1 1	2 2	2 2
24	IBD128	758	0	0	2	0	3 3	3 3	1 1	2 2	1 1	1 1	2 2	2 2
25	IBD165	963	0	0	1	0	3 3	3 3	1 1	2 2	1 1	1 1	2 2	2 2
26	IBD165	846	963	845	1	2	3 3	3 3	1 1	2 2	1 1	0 0	2 2	2 2
27	IBD165	845	0	0	2	0	3 3	3 3	1 1	2 2	1 1	1 1	2 2	2 2
28	IBD175	925	0	0	1	0	3 3	3 3	1 1	2 2	1 1	1 1	2 2	0 0
29	IBD175	972	925	929	1	2	3 3	3 3	1 1	2 2	1 1	1 1	2 2	2 2
30	IBD175	929	0	0	2	0	3 3	3 3	1 1	2 2	1 1	1 1	2 2	2 2
31	IBD182	971	0	0	1	0	1 3	3 1	4 1	4 2	2 1	3 1	2 4	3 2

Εικόνα 35: Τα δεδομένα του linkage format

Το Harpview μέχρι τώρα υποστηρίζει τις παρακάτω διαδικασίες: Μπορεί να υπολογίσει απλότυπους και ανισοροπία σύνδεσης με γονοτυπικά δεδομένα. Μπορεί να υπολογίσει την συχνότητα των απλότυπων στον πληθυσμό. Είναι δυνατή η εκτέλεση μελέτης γενετικής συσχέτισης με απλότυπους ή πολυμορφισμούς. Περιέχει υλοποίηση του αλγορίθμου επιλογής πολυμορφισμών ως δεικτών από τον Paul de Bakker. Υπάρχει η δυνατότητα λήψης δεδομένων από το HarMap, και τέλος μπορεί να απεικονίσει γραφικά τα αποτελέσματα όπως για παράδειγμα τα LD plots (Barrett 2009).

Στο πλαίσιο αυτής της εργασίας, έγινε εισαγωγή δεδομένων τα οποία παρέχονται ήδη από την εφαρμογή, και ακολουθεί μια αναφορά των αποτελεσμάτων. Τα δεδομένα γίνονται δεκτά σε πέντε συνολικά μορφές (linkage, Haps, HarMap, PHASE, PLINK) αλλά εμάς μας ενδιέφερε περισσότερο η πρώτη (linkage format) η οποία μπορούσε να χρησιμοποιηθεί και για μελέτες ασθενών - μαρτύρων. Τα αρχεία αυτών των δεδομένων είναι αρχεία .ped, και απαιτείται η χρήση ενός προγράμματος που ονομάζεται Makedped, για να αποκτήσουν την σωστή μορφή, και να είναι αποδεκτά από το Harlview. Επίσης υπάρχει ένα συνοδευτικό αρχείο το οποίο περιέχει πληροφορίες για την θέση κάθε πολυμορφισμού πάνω στο εκάστοτε χρωμόσωμα.

	A	B	C	D	E	F
1	IGR1118a_1	274044				
2	IGR1119a_1	274541				
3	IGR1143a_1	286593				
4	IGR1144a_1	287261				
5	IGR1169a_2	299755				
6	IGR1218a_2	324341				
7	IGR1219a_2	324379				
8	IGR1286a_1	358048				
9	TSC0101718	366811				
10	IGR1373a_1	395079				
11	IGR1371a_1	396353				
12	IGR1369a_2	397334				
13	IGR1369a_1	397381				
14	IGR1367a_1	398352				
15	IGR2008a_2	411823				
16	IGR2008a_1	411873				
17	IGR2010a_3	412456				
18	IGR2011b_1	413233				
19	IGR2016a_1	415579				
20	IGR2020a_1	417617				
21						
22						

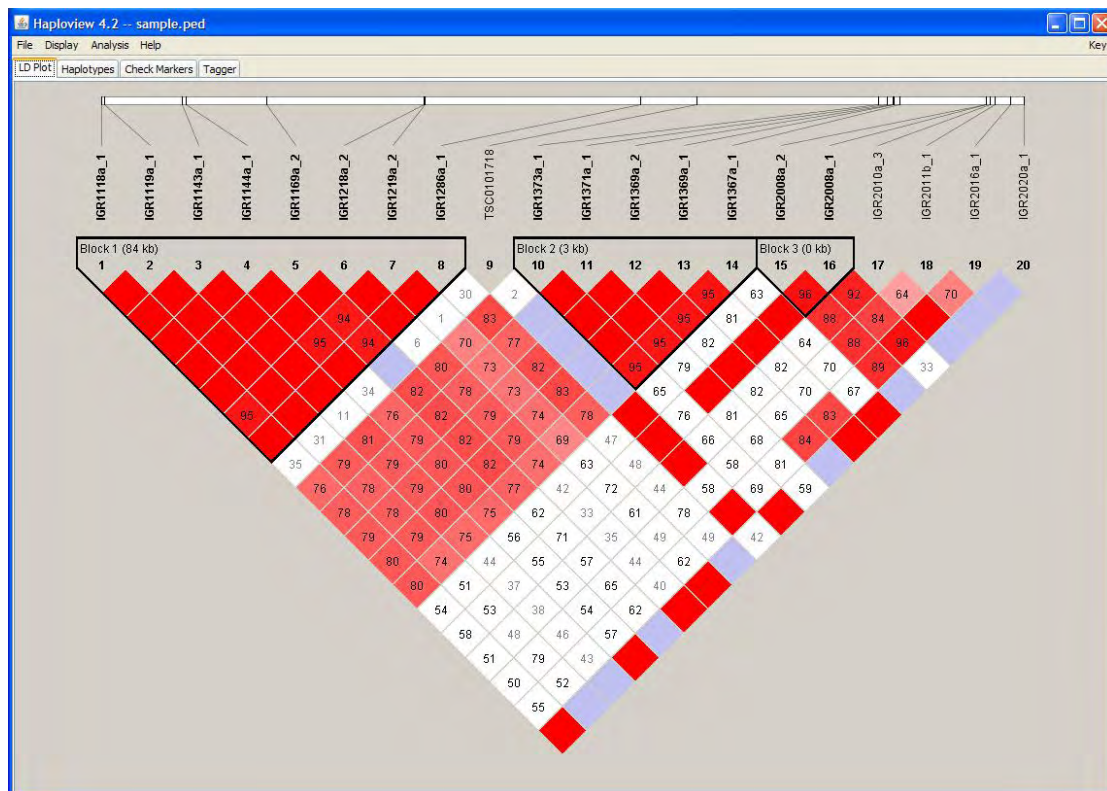
**Εικόνα 36: Οι πολυμορφισμοί και η θέση τους**

Στην εικόνα 35 βλέπουμε τα περιεχόμενα του linkage format. Στις πρώτες έξι στήλες καταγράφουμε κάποια δεδομένα που αφορούν τα άτομα. Στην πρώτη στήλη έχουμε ένα μοναδικό αλφαριθμητικό κωδικό που προσδιορίζει την οικογένεια που ανήκει το άτομο. Στην δεύτερη στήλη έχουμε έναν μοναδικό κωδικό που προσδιορίζει ένα άτομο. Στην τρίτη και τέταρτη στήλη έχουμε τους κωδικούς του πατέρα και της μητέρας του ατόμου αντίστοιχα. Στην πέμπτη στήλη έχουμε το φύλο

(1=αρσενικό, 2=θυλικό), ενώ στην έκτη στήλη έχουμε την κατάσταση επίδρασης. Η κατάσταση επίδρασης δηλώνει αν το άτομο αυτό επηρεάζεται από κάποιον παράγοντα που μελετάμε (0=άγνωστο, 1=δεν επηρεάζεται, 2=επηρεάζεται).

Από την έβδομη στήλη και μετά γίνεται η καταγραφή των γονότυπων. Στην εικόνα 35 φαίνονται με αριθμούς (1-4), και η αντιστοιχία είναι 1=A, 2=C, 3=G και 4=T. Ο χρήστης έχει τη δυνατότητα να εισάγει τα δεδομένα και με τους χαρακτήρες των βάσεων, δηλαδή A T G και C.

Το αρχείο με τα δεδομένα αυτά, πρέπει να συνοδεύεται και από ένα αρχείο info, το οποίο περιέχει μια λίστα με τους πολυμορφισμούς που ερευνούμε (εικόνα 36), καθώς και την θέση τους στο χρωμόσωμα. Η θέση των πολυμορφισμών είναι ο αριθμός των βάσεων μετρώντας από την αρχή του χρωμοσώματος.



**Εικόνα 37: LD plot**

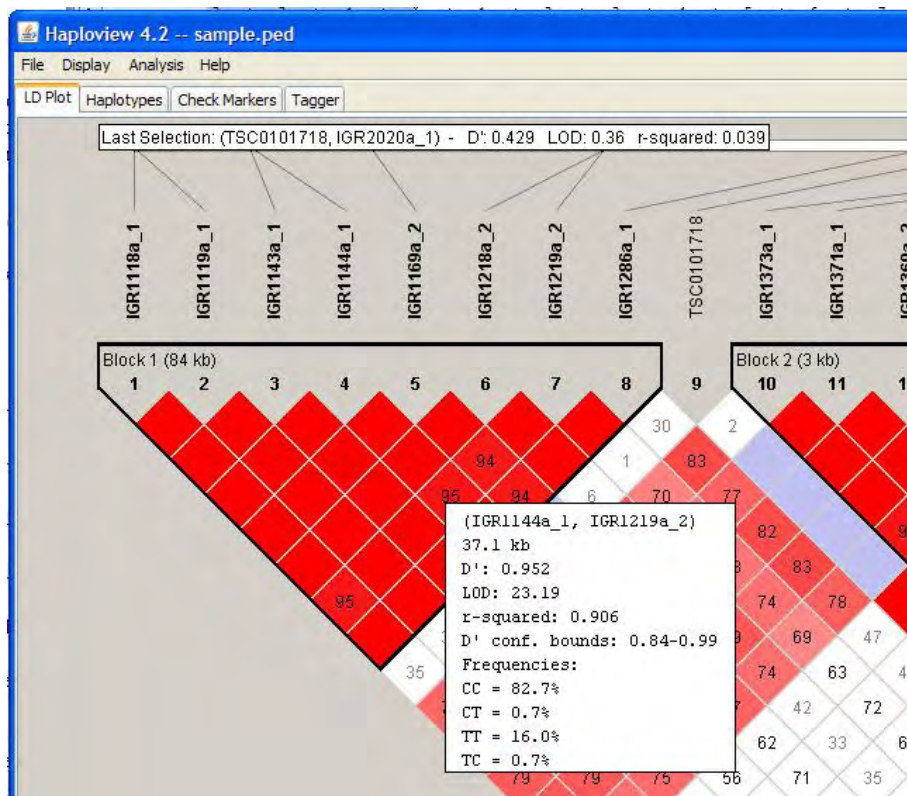
Αφού εισάγουμε τα δεδομένα, η εφαρμογή βγάζει κάποια αποτελέσματα. Το πρώτο είναι ένα διάγραμμα που δείχνει την συσχέτιση των πολυμορφισμών, και το πόσο ισχυρή είναι (εικόνα 37). Οι πιο ισχυρά σχετιζόμενοι πολυμορφισμοί ομαδοποιούνται και σχηματίζουν απλότυπους. Ο χρήστης έχει την δυνατότητα να επιλέξει έναν από τους διάφορους τρόπους που υπάρχουν για την επιλογή της ομάδας

των πολυμορφισμών που θα αποτελέσουν έναν απλότυπο, δεδομένου της ισχυρής μεταξύ τους συσχέτισης (Gabriel, Schaffner et al. 2002; Wang, Akey et al. 2002).

Επίσης ο χρήστης μπορεί να επιλέξει χειροκίνητα άλλες περιοχές για ανάλυση και όχι απαραίτητα αυτές που επιλέγει αυτόματα το Harlowiew. Μετά την επιλογή των περιοχών συσχέτισης, η εφαρμογή δημιουργεί τους απλότυπους και υπολογίζει τις συχνότητες τους στον πληθυσμό. Αυτή η απεικόνιση φαίνεται στην εικόνα 39. Οι γραμμές απεικονίζουν την μετάβαση από μια περιοχή ισχυρής συσχέτισης στην επόμενη, και το πάχος της γραμμής είναι ανάλογο με την συχνότητα μετάβασης σε αυτήν την περιοχή. Η γραμμή απεικονίζει επίσης τον βαθμό συσχέτισης (LD)  $D'$  μεταξύ των δυο περιοχών, χρησιμοποιώντας τον απλότυπο σαν το αλληλόμορφο για αυτήν την περιοχή. Διάφορες παράμετροι είναι στην επιλογή του χρήστη, όπως τα αλληλόμορφα να απεικονίζονται με γράμματα ή αριθμούς ή χρωματιστά κουτάκια, αλλά και η απεικόνιση μόνο των απλότυπων που ξεπερνούν ένα ρυθμιζόμενο κατώφλι (ή όριο) του πληθυσμού.

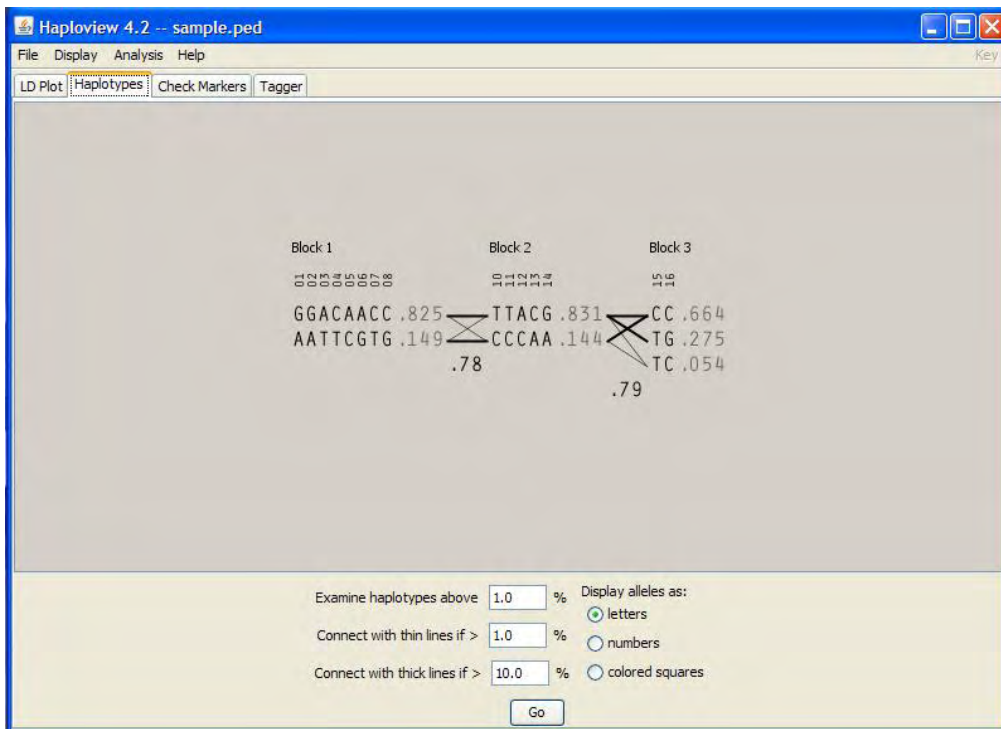
Κάθε τετραγωνάκι στο σχήμα απεικονίζει την συσχέτιση μεταξύ δύο κάθε φορά πολυμορφισμών. Το χρώμα υποδεικνύει το πόσο μεγάλη είναι η συσχέτιση, για παράδειγμα το κόκκινο χρώμα σημαίνει πολύ ισχυρή συσχέτιση. Η εφαρμογή υπολογίζει τους συντελεστές συσχέτισης  $D'$  και  $r^2$  καθώς και τις συχνότητες κάθε περίπτωσης βάσεων όπως φαίνεται στην εικόνα 38.





Εικόνα 38: Τα στοιχεία που αντιστοιχούν σε ένα ζευγάρι πολυμορφισμών

Το τελευταίο κομμάτι των αποτελεσμάτων αποτελείται από μια λίστα ζευγαριών πολυμορφισμών που εξετάσαμε. Κάθε πολυμορφισμός εξετάστηκε ως προς την συσχέτισή του με τους υπόλοιπους, με αποτέλεσμα μια λίστα που φαίνεται στην εικόνα 40. Για κάθε ζευγάρι υπολογίστηκε ο συντελεστής συσχέτισης  $r^2$ , και η εφαρμογή μπορεί να εμφανίσει όλους τους πολυμορφισμούς οι οποίοι έχουν συντελεστή συσχέτισης πάνω από κάποιο όριο, το οποίο επίσης μπορεί να επιλεχθεί από τον χρήστη. Συνοπτικά έχουμε μια αναφορά που περιέχει  $n$  πολυμορφισμούς, και  $n \times n$  υπολογισμούς συσχέτισης (γίνεται έλεγχος και με ζευγάρι δυο ίδιων πολυμορφισμών από το Haploview).



Εικόνα 39: Αποτελέσματα απλότυπων

The screenshot shows the Haploview 4.2 interface with the 'Tests' tab selected. The 'Tests' section lists various tests and their r<sup>2</sup> values. The 'Alleles captured by Current Selection' section lists the alleles captured by the current selection.

Allele	Test	r <sup>2</sup>
IGR1118a_1	IGR1219a_2	0.903
IGR1119a_1	IGR1219a_2	0.949
IGR1143a_1	IGR1219a_2	0.911
IGR1144a_1	IGR1219a_2	0.906
IGR1169a_2	IGR1219a_2	0.836
IGR1218a_2	IGR1219a_2	0.908
IGR1219a_2	IGR1219a_2	1.0
IGR1286a_1	IGR1219a_2	0.946
TSC0101718	TSC0101718	1.0
IGR1373a_1	IGR1371a_1	0.908
IGR1371a_1	IGR1371a_1	1.0
IGR1369a_2	IGR1371a_1	0.951
IGR1369a_1	IGR1371a_1	0.952
IGR1367a_1	IGR1371a_1	0.859
IGR2008a_2	IGR2008a_2	1.0
IGR2008a_1	IGR2010a_3	0.814
IGR2010a_3	IGR2010a_3	1.0
IGR2011b_1	IGR2011b_1	1.0
IGR2016a_1	IGR2010a_3	0.968
IGR2020a_1	IGR2020a_1	1.0

**Alleles captured by Current Selection**

IGR1118a\_1  
IGR1143a\_1  
IGR1119a\_1  
IGR1219a\_2  
IGR1286a\_1  
IGR1144a\_1  
IGR1169a\_2

7 SNPs in 7 tests captured 20 of 20 (100%) alleles at r<sup>2</sup> >= 0.8  
Mean max r<sup>2</sup> is 0.941

Buttons: Dump Tests File, Dump Tags File

Εικόνα 40: Η λίστα με τις συσχετίσεις των πολυμορφισμών



## 2.5 Μελέτες με μεταβλητή συνεχούς τιμής

Εκτός από τις μελέτες με ασθενείς και μάρτυρες, ασχοληθήκαμε και με ένα άλλο είδος μελετών που πραγματοποιούν απλοτυπική ανάλυση. Σε αυτή την περίπτωση έχουμε έναν παράγοντα κινδύνου, και ερευνούμε την σχέση του με τον εκάστοτε απλότυπο. Ψάχνουμε δηλαδή μια συσχέτιση μεταξύ απλότυπου και παράγοντα κινδύνου. Αυτό βέβαια δεν έχει μεγάλη διαφορά από τις μελέτες ασθενών και μαρτύρων, η μόνη διαφορά είναι ότι ο παράγοντας είναι μεταβλητή που δεν παίρνει δυο τιμές για κάθε άτομο. Είναι δηλαδή αριθμητική μεταβλητή και όχι κατηγορική. Χαρακτηριστικό παράδειγμα τέτοιου παράγοντα είναι οι τιμές της χοληστερόλης στο αίμα (Lu, Inazu et al. 2003). Συνεπώς θα έχουμε μια λίστα ανθρώπων οι οποίοι θα ανήκουν σε κάποια απλοτυπική ομάδα, και θα έχουν ο καθένας την δικιά του τιμή για την χοληστερόλη.

### 2.5.1 Απλή γραμμική παλινδρόμηση

Η απλή γραμμική παλινδρόμηση είναι ένας στατιστικός έλεγχος που εξετάζει την σχέση μεταξύ δυο μεταβλητών  $X$  και  $Y$ . Η κύρια διαδικασία είναι ο υπολογισμός των ελάχιστων τετραγώνων των μεταβλητών, από μια ευθεία γραμμή η οποία περνάει μέσα από τις παρατηρήσεις στην γραφική παράσταση. Αυτή η ευθεία είναι τοποθετημένη με τέτοιο τρόπο, ώστε να ελαχιστοποιεί το άθροισμα των κάθετων αποστάσεων των στοιχείων από αυτήν (Eberly 2007). Το επίθετο «απλή» παλινδρόμηση αναφέρεται στο γεγονός ότι αυτή η παλινδρόμηση είναι η πιο απλή στην στατιστική. Όπως είδαμε υπάρχουν και άλλες μέθοδοι παλινδρόμησης.

Η ευθεία που αναφέρεται παραπάνω έχει κλίση ίση με την συσχέτιση των δυο μεταβλητών  $X$  και  $Y$  διορθωμένη από τον λόγο των τυπικών αποκλίσεων των μεταβλητών αυτών.

Αν υποθέσουμε ότι τα δεδομένα μας έχουν μέγεθος  $n$  με  $\{Y_i, X_i\}$  όπου  $i=1,2,\dots,n$ , τότε ο σκοπός μας είναι να βρούμε την εξίσωση της ευθείας:

$$y = \alpha + \beta x \tag{27}$$

η οποία θα ταιριάζει καλύτερα στα δεδομένα. Όταν λέμε ότι θα ταιριάζει καλύτερα αναφερόμαστε στην προσέγγιση των ελάχιστων τετραγώνων, έτσι ώστε να έχουμε το

μικρότερο άθροισμα (εικόνα 41). Αν συμβαίνει αυτό, τότε οι αριθμοί  $\alpha$  και  $\beta$  λύνουν το ακόλουθο πρόβλημα ελαχιστοποίησης:

$$\min_{\alpha, \beta} Q(\alpha, \beta), \text{ όπου } Q(\alpha, \beta) = \sum_{i=1}^n \hat{\varepsilon}_i^2 = \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2 \quad (28)$$

Με την χρήση μαθηματικής ανάλυσης ή γεωμετρίας αποδεικνύεται ότι οι τιμές των  $\alpha$  και  $\beta$  που ελαχιστοποιούν την συνάρτηση  $Q$  είναι οι εξής:

$$\begin{aligned} \hat{\beta} &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n x_i y_i - \sum_{i=1}^n x_i \sum_{j=1}^n y_j / n}{\sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2 / n} \\ &= \frac{\bar{xy} - \bar{x} \bar{y}}{x^2 - \bar{x}^2} = \frac{Cov[x, y]}{Var[x]} = r_{xy} \frac{s_y}{s_x} \end{aligned} \quad (29)$$

$$\hat{\alpha} = \bar{y} - \hat{\beta} \bar{x} \quad (30)$$

όπου  $r_{xy}$  είναι ο συντελεστής συσχέτισης του δείγματος μεταξύ του  $X$  και  $Y$ ,  $s_x$  η τυπική απόκλιση του  $X$  και  $s_y$  αντίστοιχα η τυπική απόκλιση του  $Y$ .

Αν αντικαταστήσουμε τα  $\hat{\alpha}$  και  $\hat{\beta}$  στην εξίσωση 30 προκύπτει η ακόλουθη μορφή:

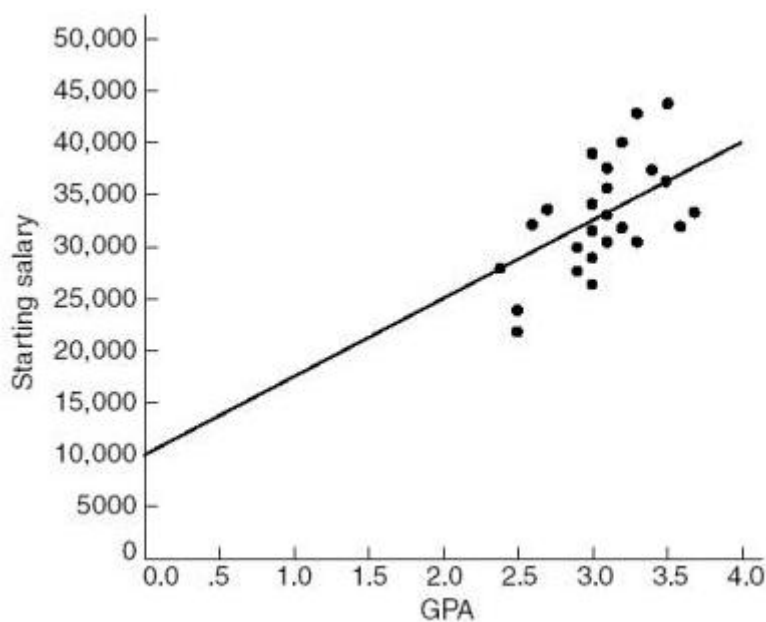
$$\frac{y - \bar{y}}{s_y} = r_{xy} \frac{x - \bar{x}}{s_x} \quad (31)$$

Αυτό δείχνει το ρόλο που παίζει ο συντελεστής συσχέτισης στην ευθεία της παλινδρόμησης.

- Η ευθεία της παλινδρόμησης περνάει από το κέντρο μάζας  $(\bar{x}, \bar{y})$ .
- Το σύνολο των στοιχείων είναι μηδέν, εάν το μοντέλο περιλαμβάνει σφάλμα:

$$\sum_{i=1}^n \hat{\varepsilon}_i = 0$$

- Τα  $\hat{\alpha}$  και  $\hat{\beta}$  είναι αμερόληπτα, αυτό σημαίνει ότι ερμηνεύουμε το μοντέλο στοχαστικά, πράγμα που σημαίνει ότι πρέπει να υποθέσουμε πως για κάθε τιμή του  $x$ , η αντίστοιχη τιμή του  $y$  παράγεται σαν ένα μέσο αποτέλεσμα  $\alpha + \beta x$  με μια επιπλέον μεταβλητή σφάλματος  $\varepsilon$ . Αυτό το σφάλμα πρέπει να είναι κατά μέσο όρο μηδέν, για κάθε τιμή του  $x$ . Σύμφωνα με αυτήν την ερμηνεία οι εκτιμήτριες ελάχιστων τετραγώνων  $\hat{\alpha}$  και  $\hat{\beta}$  θα είναι τυχαίες μεταβλητές, και θα εκτιμούν αμερόληπτα τις πραγματικές τιμές των  $\alpha$  και  $\beta$ .



**Εικόνα 41: Απλή παλινδρόμηση που δείχνει την συσχέτιση μεταξύ πρώτου μισθού και μέσου όρου βαθμολογίας φοιτητών**

Στην περίπτωση των απλοτύπων χρησιμοποιούμε την πολλαπλή παλινδρόμηση με τον ακόλουθο μαθηματικό τύπο να προσαρμόζεται στα δεδομένα:

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

(32)

για  $k$  απλότυπους.

Τα δεδομένα που παίρνουμε για την εκτέλεση της παλινδρόμησης έχουν γενικά την μορφή του πίνακα 6.

Απλότυπος	Συχνότητα	Συνεχής μεταβλητή (πχ τιμή χοληστερόλης στο αίμα)	Διακύμανση
απλότυπ_1	96	50.1	1.4
απλότυπ_2	50	55.1	3.2
...	43	53.2	3.4
...	37	50.3	3.6
απλότυπ_n	19	50.1	5.4

Πίνακας 6: Γενική μορφή των δεδομένων

Τις περισσότερες φορές, η διακύμανση δεν μας δίνεται από τις μελέτες και συνεπώς πρέπει να την υπολογίσουμε.

Για τον υπολογισμό αυτό χρησιμοποιούμε τον τύπο:

$$\text{var} = sd^2 / n \quad (33)$$

Αλλά επίσης χρήσιμος είναι και ο τύπος της τυπικής απόκλισης:

$$se^2 = sd^2 / n \quad (34)$$

και έτσι προκύπτουν άλλες δυο στήλες στον πίνακα 6 μια για την τυπική απόκλιση ( $sd$ ) και μια για το  $sd^2$  όπως φαίνεται στην εικόνα 42.

	A	B	C	D	E	F	G
1	haplotype	haplotype_nr	frequency	hdl	sd	sd_squared	variance
2	GCLCCCTGCTC B1	1	96	50.1	11.6	134.6	1.4
3	GASTTCTACCA B2	2	50	55.1	12.7	161.3	3.2
4	TCLTTCCGTCA B2	3	43	53.2	12.1	146.4	3.4
5	GCLTTCTGCCC B1	4	37	50.3	11.6	134.6	3.6
6	GCLCCTTGCCA B1	5	19	50.1	10.1	102	5.4
7	TCLTTCTACCA B1	6	16	47.8	10.6	112.4	7
8	GCLCCCTGCTA B1	7	3	45.6	9	81	27
9	TCLCTCCGTCA B2	8	3	48.6	8.2	67.2	22.4
10							

**Εικόνα 42: Τα δεδομένα στο excel για την απλή παλινδρόμηση**

Το Stata έχει την εντολή `regress` για την απλή παλινδρόμηση, αλλά σε αυτήν την εργασία χρησιμοποιούμε την διακύμανση σαν επιπλέον παράγοντα που επηρεάζει το αποτέλεσμα. Γι' αυτό το λόγο χρησιμοποιείται η εντολή `vwls` (variance-weighted least-squares regression). Για να εκτελέσουμε αυτόν τον έλεγχο πληκτρολογούμε τις παρακάτω εντολές:

```
gen se=sqrt(variance)
```

Με την εντολή `gen` δημιουργούμε μια μεταβλητή `se` (standard error) η οποία θα χρειαστεί αργότερα και είναι η τετραγωνική ρίζα της διακύμανσης, όπως προκύπτει από τις σχέσεις 32 και 33.

```
xi: vwls hdl i.haplotype_nr, sd(se)
```

Με το `xi` στην αρχή της εντολής δηλώνουμε ότι θέλουμε να χρησιμοποιήσουμε μια κατηγορική μεταβλητή η οποία θα έχει έναν δείκτη. Με την εντολή `vwls` δηλώνουμε πως θέλουμε να κάνουμε `vwls` παλινδρόμηση με εξαρτημένη μεταβλητή την `hdl` και ανεξάρτητη την `i.haplotype_nr`.

Χρησιμοποιούμε έμμεσα την διακύμανση βάζοντας σαν επιλογή μετά την εντολή το `sd(se)`, αντί για την χρήση του `fweight` ή του `aweight` που είδαμε στις παλινδρομήσεις των μελετών ασθενών και μαρτύρων. Μετά την εκτέλεση της εντολής, έχουμε τα αποτελέσματα της εικόνας 43.

```

i.haplotype_nr   _Ihaplotype_1-8   (naturally coded; _Ihaplotype_1 omitted)
Variance-weighted least-squares regression
Goodness-of-fit chi2(0) = .
Prob > chi2 = .
Number of obs = 8
Model chi2(7) = 10.33
Prob > chi2 = 0.1707
-----+-----
      hd1 |      Coef.   Std. Err.      z    P>|z|     [95% Conf. Interval]
-----+-----
_Ihaplotyp~2 |          5   2.144761    2.33   0.020   .7963456   9.203654
_Ihaplotyp~3 |    3.100002   2.19089    1.41   0.157  -1.194064   7.394068
_Ihaplotyp~4 |    .2000008   2.236068    0.09   0.929  -4.182612   4.582613
_Ihaplotyp~5 |    7.76e-16   2.607681    0.00   1.000  -5.110961   5.110961
_Ihaplotyp~6 |   -2.299999   2.898275   -0.79   0.427  -7.980514   3.380516
_Ihaplotyp~7 |     -4.5     5.329165   -0.84   0.398  -14.94497   5.944971
_Ihaplotyp~8 |     -1.5     4.878524   -0.31   0.758  -11.06173   8.061732
   _cons |     50.1     1.183216   42.34   0.000   47.78094   52.41906
-----+-----

```

**Εικόνα 43: Αποτελέσματα της vnlw παλινδρόμησης**

Όσον αφορά την ερμηνεία των αποτελεσμάτων, ισχύουν και εδώ οι βασικές αρχές της παλινδρόμησης. Κάθε απλότυπος έχει έναν συντελεστή που δηλώνει την ένταση της επιρροής στο συνολικό μοντέλο. Με το p-value βλέπουμε αν ισχύει η μηδενική υπόθεση ή όχι. Σε αυτήν την περίπτωση η εξαρτημένη μεταβλητή συνεχούς τιμής είναι τα επίπεδα της χοληστερόλης στο αίμα. Οι ανεξάρτητες μεταβλητές που επηρεάζουν τις τιμές της χοληστερόλης είναι οι διάφοροι απλότυποι, που σε αυτήν την έρευνα είναι οκτώ διαφορετικοί.

Η μέθοδος αυτή της vnlw παλινδρόμησης μπορεί να μας πει, ότι οι τιμές της εξαρτημένης μεταβλητής επηρεάζονται από κάποιον απλότυπο. Όταν έχουμε θετικό συντελεστή παλινδρόμησης αυτό σημαίνει πως ο απλότυπος αυτός συμβάλει στην αύξηση των τιμών HDL-c ή CETP αντίστοιχα. Ενώ αντίθετα ένα αρνητικό πρόσημο σημαίνει πως ο απλότυπος αυτός ρίχνει τις αντίστοιχες τιμές. Η τελευταία τιμή από τους συντελεστές παλινδρόμησης την εικόνα 43 μας δείχνει την τιμή της HDL που αντιστοιχεί στον πιο κοινό απλότυπο του δείγματος, δηλαδή τον πρώτο. Οι συντελεστές των υπόλοιπων απλότυπων δηλώνουν την απόκλιση από αυτή την τιμή. Για παράδειγμα η τιμή της HDL που αντιστοιχεί στον δεύτερο απλότυπο είναι 50.1, αν προσθέσουμε τον αντίστοιχο συντελεστή 5 θα έχουμε ως αποτέλεσμα την τιμή που αντιστοιχεί στα δεδομένα, δηλαδή 55.1 (πίνακας 7).

# ΚΕΦΑΛΑΙΟ 3 – ΑΠΟΤΕΛΕΣΜΑΤΑ ΚΑΙ ΣΥΖΗΤΗΣΗ

## 3.1 ΑΠΟΤΕΛΕΣΜΑΤΑ

### 3.1.1 Μελέτες απλότυπων με τις τιμές της HDL χοληστερίνης (συνεχής μεταβλητή)

Συνολικά έγιναν τρεις παλινδρομήσεις, μια για κάθε μελέτη. Οι δυο από αυτές τις παλινδρομήσεις εμφάνισαν στατιστικά σημαντικά αποτελέσματα σύμφωνα με τον  $X^2$  του Pearson που είδαμε παραπάνω (2.4.1), με p-value 0.0011 και 0.0032. Η παλινδρόμηση με το μικρότερο p-value είχε και τον υψηλότερο  $X^2=25.84$ . Στο σύνολο μελετήθηκαν 22 απλότυποι σε σχέση με τις τιμές της HDL-χοληστερόλης αλλά και 5 απλότυποι σε σχέση με την συγκέντρωση της CETP (cholesterol ester transfer protein), μιας πρωτεΐνης η οποία μεταφέρει τα μόρια της χοληστερόλης στο αίμα.

Η πρώτη μελέτη περιέχει απλοτυπική ανάλυση της περιοχής του υποκινητή του γονιδίου της CET πρωτεΐνης, και έχει τα δεδομένα που ψάχνουμε (Lu, Inazu et al. 2003). Αυτά αποτελούνται από οκτώ απλότυπους (εικόνα 44), τρεις από αυτούς επιδρούν αρνητικά στις τιμές της HDL-χοληστερόλης, ενώ από τους υπόλοιπους οι δυο έχουν αρκετά υψηλό συντελεστή παλινδρόμησης δείχνοντας ότι παίζουν σημαντικό ρόλο. Για τον τέταρτο απλότυπο ο συντελεστής είναι κοντά στο μηδέν, πράγμα που σημαίνει πως δεν επηρεάζει τόσο την HDL.

Στατιστικά σημαντική επιρροή θεωρούμε όμως μόνο αυτήν του δεύτερου απλότυπου, καθώς μόνο εκεί βλέπουμε p-value μικρότερο από την τιμή 0.05 και μπορούμε να απορρίψουμε την υπόθεση πως δεν υπάρχει σχέση με την ασθένεια. Γενικά βλέπουμε επίσης ότι για τον  $X^2$  στην μελέτη αυτή, το p-value είναι επίσης μεγαλύτερο από την προαναφερθείσα τιμή και συνεπώς είναι ένα στοιχείο που προδιαθέτει αρνητικά για την αξιοπιστία των αποτελεσμάτων. Τα δεδομένα που πήραμε από αυτήν την μελέτη φαίνονται στον πίνακα 7.

haplotype	haplotype_nr	frequency	hdl	sd	sd_squared	variance
GCLCCCTGCTC B1	1	96	50.1	11.6	134.6	1.4
GASTTCTACCA B2	2	50	55.1	12.7	161.3	3.2
TCLTTCCGTCA B2	3	43	53.2	12.1	146.4	3.4
GCLTTCTGCCC B1	4	37	50.3	11.6	134.6	3.6
GCLCCTTGCCA B1	5	19	50.1	10.1	102	5.4
TCLTTCTACCA B1	6	16	47.8	10.6	112.4	7
GCLCCCTGCTA B1	7	3	45.6	9	81	27
TCLCTCCGTCA B2	8	3	48.6	8.2	67.2	22.4

Πίνακας 7: Απλοτυπικά δεδομένα για τις τιμές της hdl (mg/dl)

```

i.haplotype_nr    _Ihaplotype_1-8    (naturally coded; _Ihaplotype_1 omitted)
Variance-weighted least-squares regression
Goodness-of-fit chi2(0) = .
Prob > chi2 = .
Number of obs = 8
Model chi2(7) = 10.33
Prob > chi2 = 0.1707
-----+-----
          hdl |          Coef.   Std. Err.      z    P>|z|    [95% Conf. Interval]
-----+-----
_Ihaplotyp~2 |             5    2.144761    2.33   0.020    .7963456    9.203654
_Ihaplotyp~3 |    3.100002    2.19089    1.41   0.157   -1.194064    7.394068
_Ihaplotyp~4 |    .2000008    2.236068    0.09   0.929   -4.182612    4.582613
_Ihaplotyp~5 |    7.76e-16    2.607681    0.00   1.000   -5.110961    5.110961
_Ihaplotyp~6 |   -2.299999    2.898275   -0.79   0.427   -7.980514    3.380516
_Ihaplotyp~7 |    -4.5        5.329165   -0.84   0.398   -14.94497    5.944971
_Ihaplotyp~8 |    -1.5        4.878524   -0.31   0.758   -11.06173    8.061732
   _cons |    50.1        1.183216   42.34   0.000    47.78094    52.41906
-----+-----

```

Εικόνα 44: Η πρώτη και λιγότερο σημαντική παλινδρόμηση (Lu, Inazu et al. 2003)

Η δεύτερη μελέτη η οποία αφορά την απλοτυπική ανάλυση του γονιδίου της CETP, είναι η πιο αξιόπιστη από πλευράς αποτελεσμάτων, με p-value < 0.05 και  $X^2=25.84$ . Βλέπουμε αρκετούς απλότυπους να συμβάλουν σημαντικά στις τιμές της HDL, και επιπλέον με θετικούς συντελεστές παλινδρόμησης (εικόνα 45). Οι πρώτοι τρεις απλότυποι συμβάλουν με υψηλούς συντελεστές, ο απλότυπος 2 με συντελεστή 7.8, ο απλότυπος 3 με 11.3 και ο τέταρτος με συντελεστή 10.1. Τα δεδομένα αυτής της μελέτης φαίνονται στον πίνακα 8.



haplotype	haplotype_nr	frequency	hdl	sd	sd_squared	variance
a1	1	80	47.2	14	196	2.5
c1	2	45	55	14	196	4.4
c2	3	35	58.5	14	196	5.6
a2	4	13	57.3	14	196	15
b2	5	9	55.7	14	196	22
d1	6	6	56.5	14	196	32.6
e1	7	6	64.3	14	196	32.6
d2	8	4	56.1	14	196	49
b1	9	2	54.6	14	196	98

Πίνακας 8: Απλοτυπικά δεδομένα για τις τιμές της hdl (mg/dl)

```

i.haplotype_nr      _Ihaplotype_1-9      (naturally coded; _Ihaplotype_1 omitted)
Variance-weighted least-squares regression
Goodness-of-fit chi2(0) = .
Prob > chi2 = .
Number of obs = 9
Model chi2(8) = 25.84
Prob > chi2 = 0.0011
-----
          hdl |          Coef.   Std. Err.      z    P>|z|      [95% Conf. Interval]
-----+-----
_Ihaplotype~2 |    7.799999    2.626785    2.97  0.003    2.651595    12.9484
_Ihaplotype~3 |     11.3      2.84605    3.97  0.000    5.721844    16.87815
_Ihaplotype~4 |     10.1      4.1833    2.41  0.016    1.900881    18.29912
_Ihaplotype~5 |      8.5      4.949748    1.72  0.086   -1.201327    18.20133
_Ihaplotype~6 |    9.299999    5.924525    1.57  0.116   -2.311857    20.91186
_Ihaplotype~7 |     17.1      5.924525    2.89  0.004    5.488146    28.71186
_Ihaplotype~8 |    8.899998    7.17635    1.24  0.215   -5.16539    22.96539
_Ihaplotype~9 |    7.399998   10.02497    0.74  0.460   -12.24858    27.04858
   _cons      |    47.2      1.581139   29.85  0.000    44.10103    50.29898
-----

```

Εικόνα 45: Η παλινδρόμηση της ισχυρότερης σε στατιστική σημαντικότητα μελέτης (Bauerfeind, Knoblauch et al. 2002)

Επιπλέον την μεγαλύτερη επιρροή φαίνεται να έχει ο απλότυπος 7, ο οποίος με  $p\text{-value} = 0.004$  έχει τον υψηλότερο συντελεστή παλινδρόμησης (17.1) από όλους και θεωρούμε ότι έχει αρκετά μεγάλη επιρροή στις τιμές της HDL. Όσον αφορά τους υπόλοιπους απλότυπους, παρουσιάζουν επίσης υψηλούς συντελεστές παλινδρόμησης, αλλά δεν είναι στατιστικά σημαντικοί καθώς το  $p\text{-value}$  είναι πάνω από το όριο που έχουμε θέσει στα πλαίσια του διαστήματος εμπιστοσύνης 95%.

Στην Τρίτη και τελευταία μελέτη η οποία αφορά επίσης απλοτυπική ανάλυση στο γονίδιο της CETP, έχουμε πέντε απλότυπους και ερευνούμε την συσχέτισή τους με τις τιμές της HDL (εικόνα 46), αλλά και την συγκέντρωση της πρωτεΐνης CETP (εικόνα 47). Όσον αφορά την HDL, βλέπουμε γενικά πολύ μικρούς συντελεστές παλινδρόμησης, από τους οποίους δυο είναι στατιστικά σημαντικοί. Οστόσο η επιρροή τους είναι ουδέτερη, με συντελεστή παλινδρόμησης κοντά στο μηδέν δεν έχουμε ούτε θετική ούτε αρνητική επιρροή. Τα δεδομένα όσον αφορά τις τιμές της HDL φαίνονται στον πίνακα 9.

haplotype	haplotype_nr	frequency	hdl	se	variance
1	1	290	0.9	0.007	0.000049
2	2	87	0.9	0.011	0.000121
3	3	32	0.92	0.013	0.000169
4	4	50	0.93	0.012	0.000144
5	5	197	0.94	0.009	0.000081

**Πίνακας 9: Απλοτυπικά δεδομένα για τις τιμές της hdl (mmol/l)**

```

i.haplotype_nr    _Ihaplotype_1-5    (naturally coded; _Ihaplotype_1 omitted)
Variance-weighted least-squares regression
Goodness-of-fit chi2(0) = .
Prob > chi2 = .
Number of obs = 5
Model chi2(4) = 15.88
Prob > chi2 = 0.0032
-----+-----
      hdl |          Coef.   Std. Err.      z    P>|z|     [95% Conf. Interval]
-----+-----
_Ihaplotyp~2 | -1.66e-18   .0130384    -0.00   1.000   -.0255548   .0255548
_Ihaplotyp~3 |      .02   .0147648     1.35   0.176   -.0089385   .0489386
_Ihaplotyp~4 |      .03   .0138924     2.16   0.031   .0027713   .0572287
_Ihaplotyp~5 |      .04   .0114018     3.51   0.000   .017653    .062347
   _cons |      .9     .007    128.57   0.000   .8862802   .9137197
-----+-----

```

**Εικόνα 46: Η παλινδρόμηση της μελέτης με το μεγαλύτερο δείγμα (HDL-c) (Klerkx, Tanck et al. 2003)**

Επίσης βλέπουμε στον δεύτερο απλότυπο έναν πολύ μεγάλο συντελεστή (-1.66e-18) με το e-18 να σημαίνει μετακίνηση της υποδιαστολής κατά 18 θέσεις προς τα δεξιά) και p-value πολύ μεγαλύτερο από το αποδεκτό όριο. Αυτό ίσως οφείλεται στην μεγάλη διαφορά συχνοτήτων μεταξύ του απλότυπου 2 και του 1.

Η αντίστοιχη παλινδρόμηση για τις τιμές της CETP έδειξε υψηλότερο  $X^2$  από αυτήν της HDL. Οι συντελεστές παλινδρόμησης είναι όλοι κοντά στο μηδέν, και από αυτούς μόνο η επιρροή που αντιστοιχεί στον τέταρτο και πέμπτο απλότυπο είναι στατιστικά σημαντική. Η επιρροή αυτή είναι της τάξεως του -0.07 για τον απλότυπο 4 και -0.15 περίπου για τον απλότυπο 5. Σε συνδυασμό με την τιμή της CETP η οποία είναι χαμηλότερη για τον απλότυπο αυτό, μπορούμε να συμπεράνουμε μια μικρή επιρροή, όχι όμως και για τις αντίστοιχες τιμές της HDL, όπου δεν υπάρχει σχεδόν καθόλου. Τα δεδομένα όσον αφορά τις τιμές της CETP φαίνονται στον πίνακα 10.

haplotype	haplotype_nr	frequency	CETP	se	variance
1	1	290	1.96	0.018	0.000324
2	2	87	2	0.030	0.0009
3	3	32	1.92	0.035	0.001225
4	4	50	1.89	0.036	0.001296
5	5	197	1.81	0.023	0.000529

Πίνακας 10: Απλοτυπικά δεδομένα για τις τιμές της CETP (mg/ml)

```

i.haplotype_nr    _Ihaplotype_1-5    (naturally coded; _Ihaplotype_1 omitted)
Variance-weighted least-squares regression
Goodness-of-fit chi2(0) = .
Prob > chi2 = .
Number of obs = 5
Model chi2(4) = 35.57
Prob > chi2 = 0.0000
-----
          cetp |          Coef.   Std. Err.      z    P>|z|     [95% Conf. Interval]
-----+-----
_Ihaplotyp~2 |           .04    .0349857     1.14   0.253    - .0285708   .1085707
_Ihaplotyp~3 |    -.0400001    .0393573    -1.02   0.309    - .117139   .0371389
_Ihaplotyp~4 |    -.0700001    .0402492    -1.74   0.082    - .1488871   .008887
_Ihaplotyp~5 |   -.1500001    .0292062    -5.14   0.000    - .2072431  -.0927571
   _cons      |           1.96     .018     108.89  0.000     1.924721   1.995279
-----

```

Εικόνα 47: Η παλινδρόμηση της μελέτης με το μεγαλύτερο δείγμα (CETP) (Klerkx, Tanck et al. 2003)

### 3.1.2 Μελέτες με ασθενείς και μάρτυρες (case-control)

Συνολικά πήραμε δεδομένα απλότυπων από 19 σχετικές μελέτες ασθενών-μαρτύρων οι οποίες αναφέρονται στην βιβλιογραφία. Για κάθε μελέτη εκτελέσαμε και τις τρεις μεθόδους που περιγράψαμε παραπάνω (λογιστική, πολυωνυμική και Poisson παλινδρόμηση), αν και σχεδόν ισοδύναμες (Agresti 2002) διαφέρουν αρκετά στην μέθοδο υλοποίησης και τον υπολογισμό των συντελεστών. Σε κάποιες μελέτες υπήρχαν δεδομένα απλότυπων για περισσότερα από ένα γονίδια, ή δυο και παραπάνω σύνολα απλότυπων, συνεπώς ο αριθμός των παλινδρομήσεων είναι μεγαλύτερος από τον αναμενόμενο αριθμό.

Παρατηρήσαμε 8 μελέτες που είχαν  $p$ -value < 0.05 στον έλεγχο  $X^2$  πράγμα που κάνει τα αποτελέσματά τους στατιστικώς σημαντικά (Παράρτημα 1). Οι αντίστοιχοι  $X^2$  είχαν αρκετά μεγάλο εύρος τιμών το οποίο βρίσκεται μεταξύ 76 και 9. Επίσης υπήρχε μεγάλη ποικιλία περιπτώσεων στα αποτελέσματα των παλινδρομήσεων. Για παράδειγμα είδαμε μελέτη με στατιστική σημαντικότητα, αλλά κανέναν απλότυπο να διαφοροποιείται από το σύνολο (Zee, Cook et al. 2005). Οι περισσότερες μελέτες που είχαν  $p$ -value > 0.05 στον έλεγχο  $X^2$  δεν είχαν ιδιαίτερα υψηλούς συντελεστές παλινδρόμησης και κατ' επέκταση δεν φανέρωσαν κάποια συσχέτιση με την ασθένεια. Μια μελέτη η οποία είχε  $p$ -value λίγο πάνω από την περιοχή απόρριψης της μηδενικής υπόθεσης, δηλαδή λίγο παραπάνω από 0.05



εμφάνισε κάποια συσχέτιση με ασθένεια σε δυο απλότυπους (Kankona, Stejskalova et al. 2007). Τέλος υπήρξαν και μελέτες οι οποίες έδειξαν την συσχέτιση ενός και μόνο απλότυπου με την ασθένεια, καθώς μόνο για αυτόν τον απλότυπο η συσχέτιση ήταν στατιστικά σημαντική.

Ακολουθεί η περαιτέρω ανάλυση κάποιων χαρακτηριστικών αποτελεσμάτων.

Logistic regression Number of obs = 12401  
LR chi2(6) = 50.22  
Prob > chi2 = 0.0000  
Pseudo R2 = 0.0030  
Log likelihood = -8317.8958

case	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
_Ihaplotyp~2	-.1698943	.0414179	-4.10	0.000	-.2510719 -.0887167
_Ihaplotyp~3	-.0890366	.0664107	-1.34	0.180	-.2191992 .0411259
_Ihaplotyp~4	-.3048031	.0797892	-3.82	0.000	-.4611872 -.1484191
_Ihaplotyp~5	-.8412774	.1629055	-5.16	0.000	-1.160566 -.5219884
_Ihaplotyp~6	-.0663921	.2277182	-0.29	0.771	-.5127115 .3799273
_Ihaplotyp~7	-.0811634	.1292076	-0.63	0.530	-.3344057 .1720789
_cons	-.3083014	.0268178	-11.50	0.000	-.3608632 -.2557395

Multinomial logistic regression Number of obs = 12401  
LR chi2(6) = 50.22  
Prob > chi2 = 0.0000  
Pseudo R2 = 0.0030  
Log likelihood = -8317.8958

case	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
1					
_Ihaplotyp~2	-.1698943	.0414179	-4.10	0.000	-.2510719 -.0887167
_Ihaplotyp~3	-.0890366	.0664107	-1.34	0.180	-.2191992 .0411259
_Ihaplotyp~4	-.3048031	.0797892	-3.82	0.000	-.4611872 -.1484191
_Ihaplotyp~5	-.8412774	.1629055	-5.16	0.000	-1.160566 -.5219884
_Ihaplotyp~6	-.0663921	.2277182	-0.29	0.771	-.5127115 .3799273
_Ihaplotyp~7	-.0811634	.1292076	-0.63	0.530	-.3344057 .1720789
_cons	-.3083014	.0268178	-11.50	0.000	-.3608632 -.2557395

(case==0 is the base outcome)

Poisson regression Number of obs = 84  
LR chi2(13) = 16592.09  
Prob > chi2 = 0.0000  
Pseudo R2 = 0.5877  
Log likelihood = -5820.7737

frequency	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
_Icase_1	-.3083014	.0268178	-11.50	0.000	-.3608632 -.2557395
_Ihaplotyp~2	-.2244389	.0261886	-8.57	0.000	-.2757676 -.1731102
_Ihaplotyp~3	-1.583283	.042288	-37.44	0.000	-1.666166 -1.5004
_Ihaplotyp~4	-1.873937	.0478406	-39.17	0.000	-1.967703 -1.780171
_Ihaplotyp~5	-3.015109	.0807204	-37.35	0.000	-3.173318 -2.8569
_Ihaplotyp~6	-4.225312	.1453889	-29.06	0.000	-4.510269 -3.940355
_Ihaplotyp~7	-3.053088	.0821962	-37.14	0.000	-3.214189 -2.891986
_IcasXhap~2	-.1698943	.0414179	-4.10	0.000	-.2510719 -.0887167
_IcasXhap~3	-.0890366	.0664107	-1.34	0.180	-.2191992 .0411259
_IcasXhap~4	-.3048031	.0797892	-3.82	0.000	-.4611872 -.1484191
_IcasXhap~5	-.8412774	.1629055	-5.16	0.000	-1.160566 -.5219884
_IcasXhap~6	-.0663921	.2277182	-0.29	0.771	-.5127115 .3799273
_IcasXhap~7	-.0811634	.1292076	-0.63	0.530	-.3344057 .1720789
_cons	6.304753	.0174528	361.25	0.000	6.270547 6.33896

**Εικόνα 48: Λογιστική, πολυωνμική και Poisson παλινδρόμηση στα δεδομένα της μελέτης “Haplotypes extending across ACE are associated with Alzheimer’s disease” (Kehoe, Katzov et al. 2003)**

Η πρώτη μελέτη που αναλύθηκε αφορούσε τη σχέση των απλότυπων του γονιδίου ACE με την ασθένεια του Alzheimer. Στα αποτελέσματα που βλέπουμε στην εικόνα 48, έχουμε μια καλή περίπτωση συσχέτισης. Παρατηρούμε υψηλό  $X^2$  με  $p\text{-value} < 0.05$ , πράγμα που σημαίνει αρκετά αξιόπιστα αποτελέσματα. Μπορούμε να διακρίνουμε στατιστική σημαντικότητα για τρεις απλότυπους συγκεκριμένα, τους 2, 4 και 5. Παρατηρούμε πως οι συντελεστές τους έχουν αρνητικό πρόσημο. Αυτό σημαίνει πως έχουν αρνητική επίδραση στην ασθένεια, μειώνουν δηλαδή τις πιθανότητες να εμφανιστεί αυτή η ασθένεια στο μέλλον. Ένα τμήμα των δεδομένων που χρησιμοποιήσαμε φαίνεται στον πίνακα 11.

haplotype	frequency	haplotype_nr	case	frequency_percent
AAT	202	1	1	0,5
GTC	136	2	1	0,33
GAT	37	3	1	0,09
GAC	16	4	1	0,04
ATC	4	5	1	0,01
GTT	2	6	1	0,01
AAC	11	7	1	0,03
AAT	150	1	0	0,45
GTC	130	2	0	0,39
GAT	26	3	0	0,08
GAC	19	4	0	0,06
ATC	5	5	0	0,01
GTT	0	6	0	0
AAC	6	7	0	0,02

**Πίνακας 11: Απλοτυπικά δεδομένα ασθενών-μαρτύρων**

Για να δούμε καλύτερα τι συμβαίνει επαναλαμβάνουμε την πολυωνμική παλινδρόμηση, και ζητάμε τους συντελεστές με την μορφή του σχετικού λόγου πιθανοτήτων (εικόνα 49).



τον σχετικό λόγο πιθανοτήτων με μηδενική συχνότητα σε εκείνη τη θέση. Τα απλοτυπικά δεδομένα που επεξεργαστήκαμε φαίνονται στον πίνακα 12. Παρατηρήστε τις συχνότητες (στήλη frequency) των απλότυπων 14 και 15 (στήλη haplotype\_nr) για ασθενείς (case=1) και μάρτυρες (case=0).

haplotype	haplotype_nr	frequency	case
GCGAGCCCCA	1	190	1
GCGCGTCCCCA	2	92	1
ACGAGCCCCA	3	73	1
GCGCGTCCGA	4	61	1
ATAAACACGA	5	49	1
ACGCATCCCCA	6	24	1
GCGCGCCCCA	7	12	1
GCGAGTCCCCA	8	18	1
ATGCACATCG	9	6	1
ATGCGTCCCCA	10	6	1
ACGCGTCCGA	11	3	1
GCGCGCCCGT	12	6	1
GCAAGCACCA	13	6	1
ATGAACCCGA	14	6	1
ACGAGTCCCCA	15	6	1
ACGCACCCCCA	16	49	1
GCGAGCCCCA	1	161	0
GCGCGTCCCCA	2	80	0
ACGAGCCCCA	3	47	0
GCGCGTCCGA	4	47	0
ATAAACACGA	5	43	0
ACGCATCCCCA	6	24	0
GCGCGCCCCA	7	9	0
GCGAGTCCCCA	8	4	0
ATGCACATCG	9	9	0
ATGCGTCCCCA	10	4	0
ACGCGTCCGA	11	4	0
GCGCGCCCGT	12	2	0
GCAAGCACCA	13	1	0
ATGAACCCGA	14	0	0
ACGAGTCCCCA	15	0	0
ACGCACCCCCA	16	33	0

**Πίνακας 12: Απλοτυπικά δεδομένα ασθενών-μαρτύρων**









Multinomial logistic regression  
 Log likelihood = -721.04525

Number of obs = 1075  
 LR chi2(15) = 30.15  
 Prob > chi2 = 0.0114  
 Pseudo R2 = 0.0205

case	RRR	Std. Err.	z	P> z	[95% Conf. Interval]
1					
_Ihaploty~2	.9744737	.1819003	-0.14	0.890	.6758968 1.404947
_Ihaploty~3	1.316125	.2836535	1.27	0.202	.8626705 2.007935
_Ihaploty~4	1.099776	.2438042	0.43	0.668	.7122056 1.698256
_Ihaploty~5	.9656059	.2267392	-0.15	0.882	.6094311 1.529943
_Ihaploty~6	.8473684	.260912	-0.54	0.591	.4634256 1.549403
_Ihaplotyp~7	1.129825	.5126949	0.27	0.788	.4642484 2.749613
_Ihaplotyp~8	3.813158	2.147014	2.38	0.017	1.264772 11.49628
_Ihaplotyp~9	.5649123	.303822	-1.06	0.288	.1968731 1.620972
_Ihaploty~10	1.271053	.8316813	0.37	0.714	.3525342 4.582747
_Ihaploty~11	.6355263	.4901419	-0.59	0.557	.1401699 2.881458
_Ihaploty~12	2.542105	2.093406	1.13	0.257	.5060983 12.76886
_Ihaploty~13	5.084211	5.518514	1.50	0.134	.605774 42.67136
_Ihaploty~14	4.65e+15	1.27e+23	0.00	1.000	0 .
_Ihaploty~15	4.65e+15	1.27e+23	0.00	1.000	0 .
_Ihaploty~16	1.258214	.3137611	0.92	0.357	.7717738 2.051251

(case==0 is the base outcome)

**Εικόνα 52: Απεικόνιση των συντελεστών με Odds ratio – επέκταση της εικόνας 51**

Άλλη περίπτωση αποτελεσμάτων ήταν αυτή της απλοτυπικής ανάλυσης του γονιδίου του υποδοχέα της αδρεναλίνης (Zee, Cook et al. 2005). Εδώ παρατηρήσαμε στατιστική σημαντικότητα, αλλά υπήρξαν διαφορετικοί  $X^2$  μεταξύ λογιστικής και πολυωνμικής παλινδρόμησης (εικόνα 53). Κάτι τέτοιο είναι αναμενόμενο καθώς και σε αυτήν την περίπτωση υπάρχουν απλότυποι με μηδενική συχνότητα. Ωστόσο κανένας από τους απλότυπους δεν έδειξε σημαντική συσχέτιση και τα p-value ήταν όλα παραπάνω από την περιοχή απόρριψης. Θα μπορούσαμε να πούμε για τον απλότυπο 2 ότι έχει μια μικρή τάση καθώς είναι μόλις 0.0061 μονάδες εκτός της περιοχής απόρριψης. Ίσως αν είχαμε μεγαλύτερο δείγμα να μπορούσαμε να βγάλουμε ένα καλύτερο συμπέρασμα για αυτόν τον απλότυπο. Τα δεδομένα που επεξεργαστήκαμε φαίνονται στον πίνακα 13.

haplotype	haplotype_nr	frequency	case
G16-E27-T164	1	223	1
G16-Q27-I164	2	1	1
G16-Q27-T164	3	89	1
G16-E27-I164	4	0	1
R16-Q27-T164	5	209	1
R16-Q27-I164	6	0	1
R16-E27-T164	7	1	1
R16-E27-I164	8	0	1
G16-E27-T164	1	858	0
G16-Q27-I164	2	27	0
G16-Q27-T164	3	414	0
G16-E27-I164	4	0	0
R16-Q27-T164	5	791	0
R16-Q27-I164	6	0	0
R16-E27-T164	7	0	0
R16-E27-I164	8	0	0

**Πίνακας 13: Απλοτυπικά δεδομένα ασθενών-μαρτύρων**

Ακόμη μια οριακή περίπτωση αποτελεσμάτων είχαμε και σε ακόμη μια μελέτη η οποία αφορούσε εκτίμηση γενετικού κινδύνου για την διαβητική νεφροπάθεια (Kankona, Stejskalova et al. 2007). Αυτή τη φορά όμως είχαμε οριακό αριθμό p-value στον έλεγχο  $X^2$ . Είχαμε υψηλότερο  $X^2$  από την προηγούμενη μελέτη (Zee, Cook et al. 2005), αλλά το p-value είναι πάνω από την περιοχή απόρριψης κατά 0.0004 μονάδες. Συνεπώς θα πρέπει να αναφερθούμε σε αυτήν την μελέτη. Βλέπουμε στην εικόνα 55 ότι δυο απλότυποι εμφανίζουν σημαντική συσχέτιση, οι 4 και 5 έχουν p-value<0.05. Επαναλαμβάνουμε την λογιστική παλινδρόμηση και εκτυπώνουμε τα αποτελέσματα με Odds ratio στην εικόνα 54.

Βλέπουμε λοιπόν πως οι απλότυποι 4 και 5 έχουν τον ίδιο σχετικό κίνδυνο για την ασθένεια, και αυτός είναι 2.21 φορές παραπάνω από τον απλότυπο με την μεγαλύτερη συχνότητα, δηλαδή τον πρώτο. Τα δεδομένα που επεξεργαστήκαμε φαίνονται παρακάτω στον πίνακα 14.

haplotype	haplotype_nr	frequency	case
1111	1	90	1
1211	2	43	1
2111	3	30	1
1122	4	25	1
2211	5	19	1
1222	6	12	1
2122	7	8	1
2222	8	4	1
rare	9	5	1
1111	1	199	0
1211	2	68	0
2111	3	52	0
1122	4	25	0
2211	5	19	0
1222	6	18	0
2122	7	7	0
2222	8	5	0
rare	9	19	0

**Πίνακας 14: Απλοτυπικά δεδομένα ασθενών-μαρτύρων**







Logistic regression Number of obs = 648  
 LR chi2(8) = 15.49  
 Prob > chi2 = 0.0504  
 Pseudo R2 = 0.0182

Log likelihood = -417.21246

case	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
__Ihaplotyp~2	.3351876	.2325896	1.44	0.150	-.1206797 .7910549
__Ihaplotyp~3	.2434488	.2621075	0.93	0.353	-.2702724 .7571701
__Ihaplotyp~4	.7934952	.3100584	2.56	0.010	.1857918 1.401199
__Ihaplotyp~5	.7934952	.3484242	2.28	0.023	.1105964 1.476394
__Ihaplotyp~6	.38803	.3937323	0.99	0.324	-.3836711 1.159731
__Ihaplotyp~7	.9270265	.5329103	1.74	0.082	-.1174584 1.971512
__Ihaplotyp~8	.5703516	.6827417	0.84	0.404	-.7677976 1.908501
__Ihaplotyp~9	-.5415059	.5184278	-1.04	0.296	-1.557606 .474594
__cons	-.7934952	.1270285	-6.25	0.000	-1.042466 -.5445239

Multinomial logistic regression Number of obs = 648  
 LR chi2(8) = 15.49  
 Prob > chi2 = 0.0504  
 Pseudo R2 = 0.0182

Log likelihood = -417.21246

case	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
1					
__Ihaplotyp~2	.3351876	.2325896	1.44	0.150	-.1206797 .7910549
__Ihaplotyp~3	.2434488	.2621075	0.93	0.353	-.2702724 .7571701
__Ihaplotyp~4	.7934952	.3100584	2.56	0.010	.1857918 1.401199
__Ihaplotyp~5	.7934952	.3484242	2.28	0.023	.1105964 1.476394
__Ihaplotyp~6	.38803	.3937323	0.99	0.324	-.3836711 1.159731
__Ihaplotyp~7	.9270265	.5329103	1.74	0.082	-.1174584 1.971512
__Ihaplotyp~8	.5703516	.6827417	0.84	0.404	-.7677976 1.908501
__Ihaplotyp~9	-.5415059	.5184278	-1.04	0.296	-1.557606 .474594
__cons	-.7934952	.1270285	-6.25	0.000	-1.042466 -.5445239

(case==0 is the base outcome)

Poisson regression Number of obs = 18  
 LR chi2(17) = 709.82  
 Prob > chi2 = 0.0000  
 Pseudo R2 = 0.8900

Log likelihood = -43.859131

frequency	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]
__Icase_1	-.7934952	.1270285	-6.25	0.000	-1.042466 -.5445239
__Ihaplotyp~2	-1.073797	.1404671	-7.64	0.000	-1.349108 -.7984866
__Ihaplotyp~3	-1.342061	.155743	-8.62	0.000	-1.647312 -1.03681
__Ihaplotyp~4	-2.074429	.2121912	-9.78	0.000	-2.490316 -1.658542
__Ihaplotyp~5	-2.348866	.2401181	-9.78	0.000	-2.819489 -1.878243
__Ihaplotyp~6	-2.402933	.2461314	-9.76	0.000	-2.885342 -1.920524
__Ihaplotyp~7	-3.347395	.3845546	-8.70	0.000	-4.101108 -2.593681
__Ihaplotyp~8	-3.683867	.452797	-8.14	0.000	-4.571333 -2.796401
__Ihaplotyp~9	-2.348866	.2401181	-9.78	0.000	-2.819489 -1.878243
__IcasXhap~2	.3351876	.2325896	1.44	0.150	-.1206797 .7910549
__IcasXhap~3	.2434488	.2621075	0.93	0.353	-.2702724 .7571701
__IcasXhap~4	.7934952	.3100584	2.56	0.010	.1857918 1.401199
__IcasXhap~5	.7934952	.3484242	2.28	0.023	.1105964 1.476394
__IcasXhap~6	.38803	.3937323	0.99	0.324	-.3836711 1.159731
__IcasXhap~7	.9270265	.5329103	1.74	0.082	-.1174584 1.971512
__IcasXhap~8	.5703516	.6827417	0.84	0.404	-.7677976 1.908501
__IcasXhap~9	-.5415059	.5184282	-1.04	0.296	-1.557607 .4745947
__cons	5.293305	.0708881	74.67	0.000	5.154367 5.432243

**Εικόνα 55: Λογιστική, πολυωνμική και Poisson παλινδρόμηση στα δεδομένα της μελέτης “Genetic risk factors for diabetic nephropathy on chromosomes 6p and 7q identified by the set-association approach” (Kankova, Stejskalova et al. 2007)**

## 3.2 ΣΥΖΗΤΗΣΗ

Τα αποτελέσματα αυτής της εργασίας συμφώνησαν με αυτά των αντίστοιχων μελετών. Ο αριθμός των μελετών που κατέληξαν σε κάποιο τρανταχτό συμπέρασμα ήταν μικρός, αυτό οφείλεται στο ότι οι μελέτες αυτού του τύπου είναι κυρίως ερευνητικές και έτσι είναι νωρίς να γίνεται λόγος για σίγουρα αποτελέσματα, εκτός και αν τα στοιχεία είναι πολύ ενθαρρυντικά. Στα σημεία που υπήρξαν στατιστικώς σημαντικά αποτελέσματα, ήρθαν σε συμφωνία με τα δικά μας, και αυτό απέδειξε πως η μέθοδος της παλινδρόμησης αποτέλεσε ένα καλό και αξιόπιστο εργαλείο για την ανάλυση των δεδομένων.

Ένα βασικό συμπέρασμα είναι ότι πρέπει κανείς να προσέχει την οργάνωση των δεδομένων στους πίνακες. Η κωδικοποίηση που γίνεται στο Stata έχει την αρχή να δεσμεύει τον πρώτο απλότυπο και να τον χρησιμοποιεί ως αυτόν με τη μεγαλύτερη συχνότητα. Σε αρκετές μελέτες όμως τα δεδομένα δεν ήταν ταξινομημένα κατά αύξουσα συχνότητα, και ως συνέπεια έπρεπε να αλλάξουν κατά την μεταφορά. Επίσης προσοχή χρειάζεται στην ομοιογένεια των δεδομένων. Ο τρόπος απεικόνισης μπορεί να διαφέρει από μελέτη σε μελέτη, καθώς και οι μονάδες μέτρησης.

Σε πολλές περιπτώσεις χρειάστηκαν αλλαγές κυρίως στην μετατροπή των συχνοτήτων, αλλά και τις τιμές της χοληστερίνης στο αίμα. Για την ανάλυση των μελετών με μεταβλητή συνεχούς τιμής είχαμε σαν κατάλληλο εργαλείο την απλή γραμμική παλινδρόμηση, ωστόσο τα δεδομένα ήταν λίγα και αυτή η μέθοδος δεν παρουσίαζε αποτελέσματα. Έτσι χρησιμοποιήσαμε έναν άλλο τύπο παλινδρόμησης που παρέχει το Stata, την λεγόμενη *vwls* (*variance weighted least squares*) παλινδρόμηση, η οποία ταίριαξε στα δεδομένα καθώς θέλαμε να χρησιμοποιήσουμε ως παράγοντα και την διακύμανση των τιμών της χοληστερίνης.

Όσον αφορά τις μελέτες με ασθενείς και μάρτυρες τα αποτελέσματα ήταν ενθαρρυντικά και έδειξαν πως με αυτές τις απλές μεθόδους μπορούμε να κάνουμε ανάλυση απλότυπων, αλλά και παρόμοιων δεδομένων όπως είναι οι γονότυποι αλλά και μεμονωμένοι πολυμορφισμοί. Η πιο κατάλληλη για την παρουσίαση των αποτελεσμάτων αποδείχθηκε η πολυωνμική παλινδρόμηση λόγω της δομής της, ωστόσο η διαφορά της με την λογιστική δεν είναι μεγάλη οπτικά. Από την υπολογιστική όμως πλευρά η πολυωνμική παλινδρόμηση έδειξε ένα βασικό πλεονέκτημα, το οποίο ήταν ο υπολογισμός των συντελεστών παλινδρόμησης με μεγαλύτερη ακρίβεια. Σε δυο περιπτώσεις η λογιστική παλινδρόμηση εξείρεσε

κάποιους απλότυπους από την διαδικασία επειδή είχαν μηδενική συχνότητα, ενώ η πολυωνυμική όχι. Η Poisson αρκετές φορές δεν συμφωνούσε στους υπολογισμούς του p-value για τον  $X^2$ , αλλά ούτε για τον ίδιο. Ωστόσο είχε το ίδιο αποτέλεσμα με τις υπόλοιπες στους συντελεστές παλινδρόμησης και στα αντίστοιχα p-value. Αυτό ίσως μας παροτρύνει να την χρησιμοποιούμε συνδιαστικά με τις άλλες, αλλά ίσως και καθόλου, ανάλογα κάθε περίπτωση.

Πολύ σημαντικός παράγοντας που επηρεάζει την αξιοπιστία των αποτελεσμάτων είναι το μέγεθος του δείγματος. Ένας ικανοποιητικός αριθμός ατόμων μπορεί να συμβάλει στην εξαγωγή αξιόπιστων αποτελεσμάτων. Στις μελέτες που είδαμε ο μέσος αριθμός των ατόμων ήταν 1000-1200. Υπήρξαν βέβαια περιπτώσεις πολύ κάτω από αυτό το πλαίσιο, αλλά και πολύ παραπάνω. Σε όσες περιπτώσεις είχαμε πολλά παραπάνω άτομα, παρατηρήσαμε πως αποφεύγαμε τις οριακές καταστάσεις των p-value και το αποτέλεσμα ήταν πιο ξεκάθαρο.

Αυτές οι παλινδρομήσεις είναι η βάση στην οποία μπορεί να στηριχθεί κανείς στην ανάλυση τέτοιων δεδομένων, και αποτελούν εξέλιξη του υπολογισμού του σχετικού λόγου πιθανοτήτων (Odds ratio). Για την ακρίβεια στηρίζονται πάνω στο μαθηματικό αυτό εργαλείο. Αν λάβουμε υπόψη τα πρόσθετα που μπορούμε να βάλουμε σε αυτές τις μεθόδους, όπως τα πιθανοτικά βάρη, θα έχουμε στα χέρια μας απλά και χρήσιμα εργαλεία για απλοτυπική ανάλυση.

Επιπλέον, σημαντικό είναι να ξέρουμε πως δεν μπορούμε να βγάλουμε εύκολα συμπέρασμα για μια ασθένεια, εάν οφείλεται ή όχι σε έναν απλότυπο, πολύ απλά γιατί υπάρχουν και άλλοι παράγοντες που επηρεάζουν την πορεία της. Οι απλότυποι αποτελούν ένα τμήμα της ανάλυσης του ανθρώπινου γονιδιώματος. Για την εξαγωγή κάποιου συμπεράσματος χρειάζεται αρκετός χρόνος έρευνας και στατιστικών μελετών. Το συμπέρασμα αυτής της εργασίας είναι πως θα πρέπει να συμπεριλαμβάνουμε στην μελέτη γενετικής συσχέτισης την ανάλυση των απλότυπων ως ένα παραπάνω στοιχείο που θα οδηγήσει σε κάποιο συμπέρασμα. Η βελτιστοποίηση αυτής της ανάλυσης ελπίζουμε πως θα οδηγήσει σε πιο αξιόπιστα και αληθή αποτελέσματα στο μέλλον.

Η σημαντικότητα της απλοτυπικής ανάλυσης φαίνεται καθώς πολλές μελέτες πλέον την συμπεριλαμβάνουν στην δουλειά τους. Από την πρώτη εμφάνιση των απλοτυπικών αναλύσεων την δεκαετία του '90 και μέχρι σήμερα παρατηρούμε έναν αυξανόμενο αριθμό μελετών που περιέχουν τέτοια δεδομένα. Το μέλλον είναι αρκετά



υποσχόμενο καθώς ήδη αποτυπικά δεδομένα είναι διαθέσιμα και στο διαδίκτυο μέσω του HarMap project, πράγμα που ενθαρρύνει την ενασχόληση φοιτητών/ερευνητών με το θέμα.

Ακολουθούν όλες οι υπόλοιπες μελέτες που δεν αναλύθηκαν. Τις χωρίσαμε σε 2 κατηγορίες. Η πρώτη κατηγορία είναι μελέτες με απόρριψη της μηδενικής υπόθεσης ( $p\text{-value}<0.05$ ) του  $X^2$  και η δεύτερη με αποδοχή της μηδενικής υπόθεσης ( $p\text{-value}>0.05$ ) του  $X^2$ .



## **Παράρτημα 1 – απόρριψη της μηδενικής υπόθεσης**





haplotype	haplotype_nr	frequency	percent	case
GAGCCA	1	7094	0,31	1
GAGCCG	2	5950	0,26	1
GGCACCG	3	4577	0,2	1
GGCACTG	4	2060	0,09	1
TGCACCG	5	1831	0,08	1
GAGCCA	1	3880	0,3	1
GAGCCG	2	3363	0,26	1
GGCACCG	3	2587	0,2	1
GGCACTG	4	1035	0,08	1
TGCACCG	5	1164	0,09	1
GAGCCA	1	9720	0,32	0
GAGCCG	2	7594	0,25	0
GGCACCG	3	6379	0,21	0
GGCACTG	4	3037	0,1	0
TGCACCG	5	2430	0,08	0
GAGCCA	1	3768	0,28	0
GAGCCG	2	3499	0,26	0
GGCACCG	3	2692	0,2	0
GGCACTG	4	1077	0,08	0
TGCACCG	5	1211	0,09	0

**Πίνακας 16: Δεδομένα της μελέτης “Genetic variation in IL6 gene and type 2 diabetes: tagging-SNP haplotype analysis in large-scale case-control study and meta-analysis” (Qi, van Dam et al. 2006)**



haplotype	haplotype_nr	frequency	case
CGGAT	1	355	1
CGGTT	2	163	1
CAGTT	3	45	1
AAGTT	4	35	1
AAATC	5	129	1
Others	6	45	1
CGGAT	1	286	0
CGGTT	2	185	0
CAGTT	3	63	0
AAGTT	4	21	0
AAATC	5	68	0
Others	6	129	0

**Πίνακας 17: Δεδομένα της μελέτης “Tagging SNPs in non-homologous end-joining pathway genes and risk of glioma” (Liu, Zhang et al. 2007) – για το γονίδιο XRCC5**











haplotype	haplotype_nr	frequency	case
TTGGAG	1	135	1
CTGGGG	2	86	1
TAGGAG	3	81	1
TAGGAA	4	54	1
TTGTAG	5	25	1
TTSGAG	6	5	1
TTGGAA	7	4	1
CTGGAG	8	0	1
TTGGGG	9	3	1
Others	10	3	1
TTGGAG	1	185	0
CTGGGG	2	63	0
TAGGAG	3	109	0
TAGGAA	4	47	0
TTGTAG	5	21	0
TTSGAG	6	5	0
TTGGAA	7	2	0
CTGGAG	8	15	0
TTGGGG	9	6	0
Others	10	3	0

**Πίνακας 20: Δεδομένα της μελέτης “Haplotype analysis of the RAGE gene: identification of a haplotype marker for diabetic nephropathy in type 2 diabetes mellitus” (Kankova, Stejskalova et al. 2005)**





## **Παράρτημα 2 – αποδοχή της μηδενικής υπόθεσης**





haplotype	haplotype_nr	case	frequency
GTG	1	1	177
TCGA	2	1	128
CCCC	3	1	26
TCGC	4	1	5
CGCC	5	1	167
AGA	6	1	126
GGG	7	1	23
GTA	8	1	9
GTG	1	0	256
TCGA	2	0	168
CCCC	3	0	36
TCGC	4	0	9
CGCC	5	0	219
AGA	6	0	156
GGG	7	0	23
GTA	8	0	11

**Πίνακας 21: Δεδομένα της μελέτης “Analysis of Candidate Genes on Chromosomes 5q and 19p  
in  
Celiac Disease” (Latiano, Mora et al. 2007)**





haplotype	haplotype_nr	frequency	case
CGC	1	28	1
TTT	2	24	1
CGT	3	2	1
TGC	4	6	1
CTT	5	6	1
TGT	6	2	1
TTC	7	1	1
CTC	8	1	1
CGC	1	23	0
TTT	2	26	0
CGT	3	14	0
TGC	4	5	0
CTT	5	5	0
TGT	6	3	0
TTC	7	2	0
CTC	8	1	0

**Πίνακας 23: Δεδομένα της μελέτης “Common ABCB1 polymorphisms are not associated with multidrug resistance in epilepsy using a gene-wide tagging approach” (Leschziner, Andrew et al. 2007)**

















haplotype	haplotype_nr	frequency	case
ATAA	1	786	1
GGGG	2	646	1
GTGA	3	172	1
GGGA	4	171	1
GTAA	5	46	1
AGGG	6	36	1
ATGA	7	21	1
ATAA	1	756	0
GGGG	2	661	0
GTGA	3	170	0
GGGA	4	158	0
GTAA	5	48	0
AGGG	6	52	0
ATGA	7	31	0

**Πίνακας 28: Δεδομένα της μελέτης “RGS4 is not a susceptibility gene for schizophrenia in Japanese: Association study in a large case-control population” (Ishiguro, Horiuchi et al. 2007)**





haplotype	haplotype_nr	frequency	case
GTGTGCCGGAGCCCATATA	1	42	1
GCGTGCCGGAGCCGGCA	2	39	1
TTACTTGGAAACCGCATA	3	21	1
GTGTGCCGGAGCCACACA	4	14	1
TTACTTGGAAACCGCATG	5	12	1
GCGTGCCGGAGTCCGGCA	6	5	1
RARE	7	18	1
GTGTGCCGGAGCCCATATA	1	50	0
GCGTGCCGGAGCCGGCA	2	34	0
TTACTTGGAAACCGCATA	3	33	0
GTGTGCCGGAGCCACACA	4	17	0
TTACTTGGAAACCGCATG	5	9	0
GCGTGCCGGAGTCCGGCA	6	14	0
RARE	7	17	0

**Πίνακας 29: Δεδομένα της μελέτης “Polymorphism discovery in 62 DNA repair genes and haplotype associations with risks for lung and head and neck cancers” (Michiels, Danoy et al. 2007) – Για το γονίδιο MSH3**





haplotype	haplotype_nr	frequency	case
A T C G A	1	82	1
A C T A C	2	27	1
A C C G A	3	32	1
C C C G A	4	11	1
A T C G A	1	89	0
A C T A C	2	43	0
A C C G A	3	33	0
C C C G A	4	5	0

**Πίνακας 30: Δεδομένα της μελέτης “Polymorphism discovery in 62 DNA repair genes and haplotype associations with risks for lung and head and neck cancers” (Michiels, Danoy et al. 2007) – Για το γονίδιο ERCC5**

```

Logistic regression
Log likelihood = -220.04458
Number of obs = 322
LR chi2(3) = 5.29
Prob > chi2 = 0.1517
Pseudo R2 = 0.0119

-----
case | Coef. Std. Err. z P>|z| [95% Conf. Interval]
-----+-----
_Ihaplotyp~2 | -.3834461 .2893509 -1.33 0.185 [-.9505634 .1836712]
_Ihaplotyp~3 | .0511455 .2915203 0.18 0.861 [-.5202239 .6225148]
_Ihaplotyp~4 | .8703745 .5606605 1.55 0.121 [-.2284999 1.969249]
_cons | -.0819171 .1530721 -0.54 0.593 [-.381933 .2180987]
-----

Multinomial logistic regression
Log likelihood = -220.04458
Number of obs = 322
LR chi2(3) = 5.29
Prob > chi2 = 0.1517
Pseudo R2 = 0.0119

-----
case | Coef. Std. Err. z P>|z| [95% Conf. Interval]
-----+-----
1
_Ihaplotyp~2 | -.3834461 .2893509 -1.33 0.185 [-.9505634 .1836712]
_Ihaplotyp~3 | .0511455 .2915203 0.18 0.861 [-.5202239 .6225148]
_Ihaplotyp~4 | .8703745 .5606605 1.55 0.121 [-.2284999 1.969249]
_cons | -.0819171 .1530721 -0.54 0.593 [-.381933 .2180987]
-----
(case==0 is the base outcome)

Poisson regression
Log likelihood = -20.848817
Number of obs = 8
LR chi2(7) = 164.89
Prob > chi2 = 0.0000
Pseudo R2 = 0.7982

-----
frequency | Coef. Std. Err. z P>|z| [95% Conf. Interval]
-----+-----
_Icase_1 | -.0819171 .1530721 -0.54 0.593 [-.381933 .2180987]
_Ihaplotyp~2 | -.7274363 .1857196 -3.92 0.000 [-1.09144 -.3634325]
_Ihaplotyp~3 | -.9921288 .2038112 -4.87 0.000 [-1.391591 -.5926663]
_Ihaplotyp~4 | -2.879198 .4596041 -6.26 0.000 [-3.780006 -1.978391]
_IcasXhap~2 | -.3834461 .2893509 -1.33 0.185 [-.9505634 .1836712]
_IcasXhap~3 | .0511455 .2915203 0.18 0.861 [-.5202239 .6225148]
_IcasXhap~4 | .8703745 .5606605 1.55 0.121 [-.2284999 1.969249]
_cons | 4.488636 .1059998 42.35 0.000 [4.280881 4.696392]
-----

```

**Εικόνα 71: Λογιστική, πολυωνμική και Poisson παλινδρόμηση στα δεδομένα της μελέτης “Polymorphism discovery in 62 DNA repair genes and haplotype associations with risks for lung and head and neck cancers” (Michiels, Danoy et al. 2007) – Για το γονίδιο ERCC5**





haplotype	haplotype_nr	frequency	case
TTGTCC	1	674	1
CGAATT	2	111	1
TTGACC	3	30	1
TGAATT	4	33	1
CGATTT	5	12	1
CGAATC	6	9	1
23 others	7	44	1
TTGTCC	1	681	0
CGAATT	2	115	0
TTGACC	3	29	0
TGAATT	4	26	0
CGATTT	5	12	0
CGAATC	6	6	0
23 others	7	43	0

**Πίνακας 32: Δεδομένα της μελέτης “Haplotypic variation in MRE11, RAD50 and NBS1 and risk of non-Hodgkin’s lymphoma” (Rollinson, Kesby et al. 2006)**



haplotype	haplotype_nr	frequency	case
AACGCCT	1	295	1
AACGCTC	2	161	1
AGCACCC	3	105	1
AACGCC	4	35	1
TACGCCT	5	42	1
AACGGTC	6	35	1
Rare	7	21	1
AACGCCT	1	640	0
AACGCTC	2	359	0
AGCACCC	3	219	0
AACGCC	4	109	0
TACGCCT	5	94	0
AACGGTC	6	62	0
Rare	7	62	0

**Πίνακας 33: Δεδομένα της μελέτης “Effect of ATM, CHEK2 and ERBB2 TAGSNPs and haplotypes on endometrial cancer risk” (Einarsdottir, Humphreys et al. 2007) – Για το γονίδιο ATM**



Logistic regression  
 Log likelihood = -1383.0763  
 Number of obs = 2239  
 LR chi2(6) = 6.04  
 Prob > chi2 = 0.4192  
 Pseudo R2 = 0.0022

case	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
_Ihaplotyp~2	-.0274252	.118106	-0.23	0.816	-.2589087	.2040583
_Ihaplotyp~3	.0393814	.137994	0.29	0.775	-.2310819	.3098448
_Ihaplotyp~4	-.361507	.2066351	-1.75	0.080	-.7665044	.0434904
_Ihaplotyp~5	-.0311323	.1984947	-0.16	0.875	-.4201748	.3579102
_Ihaplotyp~6	.2027065	.2228291	0.91	0.363	-.2340304	.6394434
_Ihaplotyp~7	-.3081191	.2621076	-1.18	0.240	-.8218406	.2056024
_cons	-.7744928	.0703728	-11.01	0.000	-.912421	-.6365647

Multinomial logistic regression  
 Log likelihood = -1383.0763  
 Number of obs = 2239  
 LR chi2(6) = 6.04  
 Prob > chi2 = 0.4192  
 Pseudo R2 = 0.0022

case	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
1						
_Ihaplotyp~2	-.0274252	.118106	-0.23	0.816	-.2589087	.2040583
_Ihaplotyp~3	.0393814	.137994	0.29	0.775	-.2310819	.3098448
_Ihaplotyp~4	-.361507	.2066351	-1.75	0.080	-.7665044	.0434904
_Ihaplotyp~5	-.0311323	.1984947	-0.16	0.875	-.4201748	.3579102
_Ihaplotyp~6	.2027065	.2228291	0.91	0.363	-.2340304	.6394434
_Ihaplotyp~7	-.3081191	.2621076	-1.18	0.240	-.8218406	.2056024
_cons	-.7744928	.0703728	-11.01	0.000	-.912421	-.6365647

(case==0 is the base outcome)

Poisson regression  
 Log likelihood = -45.15153  
 Number of obs = 14  
 LR chi2(13) = 1939.51  
 Prob > chi2 = 0.0000  
 Pseudo R2 = 0.9555

frequency	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
_Icase_1	-.7744928	.0703728	-11.01	0.000	-.912421	-.6365647
_Ihaplotyp~2	-.5781458	.0659395	-8.77	0.000	-.7073848	-.4489068
_Ihaplotyp~3	-1.072396	.0782861	-13.70	0.000	-1.225834	-.9189586
_Ihaplotyp~4	-1.77012	.1036186	-17.08	0.000	-1.973209	-1.567032
_Ihaplotyp~5	-1.918173	.1104572	-17.37	0.000	-2.134666	-1.701681
_Ihaplotyp~6	-2.334334	.1330095	-17.55	0.000	-2.595028	-2.07364
_Ihaplotyp~7	-2.334334	.1330095	-17.55	0.000	-2.595028	-2.07364
_IcasXhap~2	-.0274252	.118106	-0.23	0.816	-.2589087	.2040583
_IcasXhap~3	.0393814	.137994	0.29	0.775	-.2310819	.3098448
_IcasXhap~4	-.361507	.2066351	-1.75	0.080	-.7665044	.0434904
_IcasXhap~5	-.0311323	.1984947	-0.16	0.875	-.4201748	.3579102
_IcasXhap~6	.2027065	.2228291	0.91	0.363	-.2340304	.6394434
_IcasXhap~7	-.3081191	.2621076	-1.18	0.240	-.8218406	.2056024
_cons	6.461468	.0395285	163.46	0.000	6.383994	6.538943

**Εικόνα 74: Λογιστική, πολωνομική και Poisson παλινδρόμηση στα δεδομένα της μελέτης “Effect of ATM, CHEK2 and ERBB2 TAGSNPs and haplotypes on endometrial cancer risk” (Einarsdottir, Humphreys et al. 2007) – Για το γονίδιο ATM**

haplotype	haplotype_nr	frequency	case
GCCCCC	1	154	1
GGCTGC	2	161	1
GCCCCG	3	91	1
ACCCGC	4	91	1
GCTCGG	5	56	1
GGCCGC	6	35	1
Rare	7	105	1
GCCCCC	1	375	0
GGCTGC	2	359	0
GCCCCG	3	203	0
ACCCGC	4	156	0
GCTCGG	5	125	0
GGCCGC	6	94	0
Rare	7	234	0

**Πίνακας 34: Δεδομένα της μελέτης “Effect of ATM, CHEK2 and ERBB2 TAGSNPs and haplotypes on endometrial cancer risk” (Einarsdottir, Humphreys et al. 2007) – Για το γονίδιο CHEK2**











haplotype	haplotype_nr	frequency	case
G G C G A	1	163	1
G G A A A	2	11	1
G G C G G	3	52	1
G G A A G	4	49	1
G A A A G	5	10	1
G A A A A	6	25	1
G A A G G	7	23	1
G A A G A	8	12	1
G A C G A	9	57	1
A A A A G	10	46	1
A A A A A	11	98	1
G G C G A	1	180	0
G G A A A	2	17	0
G G C G G	3	47	0
G G A A G	4	51	0
G A A A G	5	21	0
G A A A A	6	27	0
G A A G G	7	15	0
G A A G A	8	16	0
G A C G A	9	64	0
A A A A G	10	75	0
A A A A A	11	116	0

**Πίνακας 37: Δεδομένα της μελέτης “Cholesteryl Ester Transfer Protein (CETP) Genetic Variation and Early Onset of Non-Fatal Myocardial Infarction” (Meiner, Friedlander et al. 2008)**







## BIBΛΙΟΓΡΑΦΙΑ

- Agresti, A. (2002). "Categorical Data Analysis, 2nd edn.: John Wiley & Sons."
- Baker, S. G. (2005). "A simple loglinear model for haplotype effects in a case-control study involving two unphased genotypes." Stat Appl Genet Mol Biol **4**: Article14.
- Barrett, J. C. (2009). "Haploview: Visualization and analysis of SNP genotype data." CSH Protoc **2009**(10): pdb ip71.
- Bauerfeind, A., H. Knohlauch, et al. (2002). "Single nucleotide polymorphism haplotypes in the cholesteryl-ester transfer protein (CETP) gene and lipid phenotypes." Hum Hered **54**(4): 166-73.
- Berk, R. and J. MacDonald (2007). "Overdispersion and Poisson Regression." Journal of Quantitative Criminology **24**(3): 269-284.
- Bernig, T., B. J. Boersma, et al. (2007). "The mannose-binding lectin (MBL2) haplotype and breast cancer: an association study in African-American and Caucasian women." Carcinogenesis **28**(4): 828-36.
- Chen, H. Y. (2003). "A note on the prospective analysis of outcome-dependent samples." Journal of the Royal Statistical Society **65**(2): 575-584.
- Chen, X. and Z. Li (2008). "Inference of haplotype effects in case-control studies using unphased genotype and environmental data." Biom J **50**(2): 270-82.
- Chen, Y. H. and J. T. Kao (2006). "Multinomial logistic regression approach to haplotype association analysis in population-based case-control studies." BMC Genet **7**: 43.
- Clark, A. G. (1990). "Inference of haplotypes from PCR-amplified samples of diploid populations." Mol Biol Evol **7**(2): 111-22.
- Clegg, M. T., J. F. Kidwell, et al. (1976). "Dynamics of correlated genetic systems. I. Selection in the region of the Glued locus of *Drosophila melanogaster*." Genetics **83**(4): 793-810.
- Consonni, D., S. De Matteis, et al. (2010). "Lung cancer and occupation in a population-based case-control study." Am J Epidemiol **171**(3): 323-33.
- Deloukas, P. and D. Bentley (2004). "The HapMap project and its application to genetic studies of drug response." Pharmacogenomics J **4**(2): 88-90.
- Devlin, B. and N. Risch (1995). "A comparison of linkage disequilibrium measures for fine-scale mapping." Genomics **29**(2): 311-22.
- Eberly, L. E. (2007). "Correlation and simple linear regression." Methods Mol Biol **404**: 143-64.
- Einarsdottir, K., K. Humphreys, et al. (2007). "Effect of ATM, CHEK2 and ERBB2 TAGSNPs and haplotypes on endometrial cancer risk." Hum Mol Genet **16**(2): 154-64.
- Epstein, M. P. and G. A. Satten (2003). "Inference on haplotype effects in case-control studies using unphased genotype data." Am J Hum Genet **73**(6): 1316-29.
- Fallin, D., A. Cohen, et al. (2001). "Genetic analysis of case/control data using estimated haplotype frequencies: application to APOE locus variation and Alzheimer's disease." Genome Res **11**(1): 143-51.
- French, B., T. Lumley, et al. (2006). "Simple estimates of haplotype relative risks in case-control data." Genet Epidemiol **30**(6): 485-94.
- Gabriel, S. B., S. F. Schaffner, et al. (2002). "The structure of haplotype blocks in the human genome." Science **296**(5576): 2225-9.

- Gusfield, D. (2000). "A practical algorithm for optimal inference of haplotypes from diploid populations." Proc Int Conf Intell Syst Mol Biol **8**: 183-9.
- Gusfield, D. (2001). "Inference of haplotypes from samples of diploid populations: complexity and algorithms." J Comput Biol **8**(3): 305-23.
- HapMap (2003). "The International HapMap Project." Nature **426**(6968): 789-96.
- HapMap (2004). "Integrating ethics and science in the International HapMap Project." Nat Rev Genet **5**(6): 467-75.
- HapMap (2005). "A haplotype map of the human genome." Nature **437**(7063): 1299-320.
- Hill, W. G. and B. S. Weir (1994). "Maximum-likelihood estimation of gene location by linkage disequilibrium." Am J Hum Genet **54**(4): 705-14.
- Ishiguro, H., Y. Horiuchi, et al. (2007). "RGS4 is not a susceptibility gene for schizophrenia in Japanese: association study in a large case-control population." Schizophr Res **89**(1-3): 161-4.
- Johnatty, S. E., A. B. Spurdle, et al. (2008). "Progesterone receptor polymorphisms and risk of breast cancer: results from two Australian breast cancer studies." Breast Cancer Res Treat **109**(1): 91-9.
- Kankova, K., A. Stejskalova, et al. (2005). "Haplotype analysis of the RAGE gene: identification of a haplotype marker for diabetic nephropathy in type 2 diabetes mellitus." Nephrol Dial Transplant **20**(6): 1093-102.
- Kankova, K., A. Stejskalova, et al. (2007). "Genetic risk factors for diabetic nephropathy on chromosomes 6p and 7q identified by the set-association approach." Diabetologia **50**(5): 990-9.
- Kedda, M. A., D. L. Duffy, et al. (2006). "ADAM33 haplotypes are associated with asthma in a large Australian population." Eur J Hum Genet **14**(9): 1027-36.
- Kehoe, P. G., H. Katzov, et al. (2003). "Haplotypes extending across ACE are associated with Alzheimer's disease." Hum Mol Genet **12**(8): 859-67.
- Klerkx, A. H., M. W. Tanck, et al. (2003). "Haplotype analysis of the CETP gene: not TaqIB, but the closely linked -629C-->A polymorphism and a novel promoter variant are independently associated with CETP concentration." Hum Mol Genet **12**(2): 111-23.
- Latiano, A., B. Mora, et al. (2007). "Analysis of candidate genes on chromosomes 5q and 19p in celiac disease." J Pediatr Gastroenterol Nutr **45**(2): 180-6.
- Leschziner, G. D., T. Andrew, et al. (2007). "Common ABCB1 polymorphisms are not associated with multidrug resistance in epilepsy using a gene-wide tagging approach." Pharmacogenet Genomics **17**(3): 217-20.
- Lin, C. P. and C. S. Fann (2009). "A novel tool for individual haplotype inference using mixed data." J Biomed Sci **16**: 52.
- Lin, D. Y. and B. E. Huang (2007). "The Use of Inferred Haplotypes in Downstream Analyses." The American Journal of Human Genetics **80**(3): 577-579.
- Lin, D. Y., D. Zeng, et al. (2005). "Maximum likelihood estimation of haplotype effects and haplotype-environment interactions in association studies." Genet Epidemiol **29**(4): 299-312.
- Lin, S., D. J. Cutler, et al. (2002). "Haplotype inference in random population samples." Am J Hum Genet **71**(5): 1129-37.
- Liu, Y., H. Zhang, et al. (2007). "Tagging SNPs in non-homologous end-joining pathway genes and risk of glioma." Carcinogenesis **28**(9): 1906-13.
- Lu, H., A. Inazu, et al. (2003). "Haplotype analyses of cholesteryl ester transfer protein gene promoter: a clue to an unsolved mystery of TaqIB polymorphism." J Mol Med **81**(4): 246-55.

- Mander, A. P. (2001). "Haplotype analysis in population-based association studies." The Stata Journal **1**(1): 58–75.
- Marchini, J., D. Cutler, et al. (2006). "A comparison of phasing algorithms for trios and unrelated individuals." Am J Hum Genet **78**(3): 437-50.
- McCullagh, P. N., J. (1999). Generalized linear models.
- Meiner, V., Y. Friedlander, et al. (2008). "Cholesteryl ester transfer protein (CETP) genetic variation and early onset of non-fatal myocardial infarction." Ann Hum Genet **72**(Pt 6): 732-41.
- Michiels, S., P. Danoy, et al. (2007). "Polymorphism discovery in 62 DNA repair genes and haplotype associations with risks for lung and head and neck cancers." Carcinogenesis **28**(8): 1731-9.
- Nei, M. and W. H. Li (1980). "Non-random association between electromorphs and inversion chromosomes in finite populations." Genet Res **35**(1): 65-83.
- Niu, T. (2004). "Algorithms for inferring haplotypes." Genet Epidemiol **27**(4): 334-47.
- Olson, J. M. and E. M. Wijsman (1994). "Design and sample-size considerations in the detection of linkage disequilibrium with a disease locus." Am J Hum Genet **55**(3): 574-80.
- Prentice, R. L. and R. Pyke (1979). "Logistic disease incidence models and case-control studies." Biometrika **66**(3): 403-411.
- Qi, L., R. M. van Dam, et al. (2006). "Genetic variation in IL6 gene and type 2 diabetes: tagging-SNP haplotype analysis in large-scale case-control study and meta-analysis." Hum Mol Genet **15**(11): 1914-20.
- Qian, L., J. Zhao, et al. (2007). "Brain-derived neurotrophic factor and risk of schizophrenia: an association study and meta-analysis." Biochem Biophys Res Commun **353**(3): 738-43.
- Rollinson, S., H. Kesby, et al. (2006). "Haplotypic variation in MRE11, RAD50 and NBS1 and risk of non-Hodgkin's lymphoma." Leuk Lymphoma **47**(12): 2567-83.
- Satten, G. A. and M. P. Epstein (2004). "Comparison of prospective and retrospective methods for haplotype inference in case-control studies." Genet Epidemiol **27**(3): 192-201.
- Schaid, D. J., C. M. Rowland, et al. (2002). "Score tests for association between traits and haplotypes when linkage phase is ambiguous." Am J Hum Genet **70**(2): 425-34.
- Stephens, M. and P. Donnelly (2003). "A comparison of bayesian methods for haplotype reconstruction from population genotype data." Am J Hum Genet **73**(5): 1162-9.
- Sun, S., C. M. Greenwood, et al. (2007). "Haplotype inference using a Bayesian Hidden Markov model." Genet Epidemiol **31**(8): 937-48.
- Thomas, A. (2005). "Characterizing allelic associations from unphased diploid data by graphical modeling." Genet Epidemiol **29**(1): 23-35.
- Tiret, L., P. Amouyel, et al. (1991). "Testing for association between disease and linked marker loci: a log-linear-model analysis." Am J Hum Genet **48**(5): 926-34.
- Umbach, D. M. and C. R. Weinberg (1997). "Designing and analysing case-control studies to exploit independence of genotype and exposure." Stat Med **16**(15): 1731-43.

- Valdes, A. M., J. Loughlin, et al. (2007). "Sex and ethnic differences in the association of ASPN, CALM1, COL2A1, COMP, and FRZB with genetic susceptibility to osteoarthritis of the knee." *Arthritis Rheum* **56**(1): 137-46.
- Wallenstein, S., S. E. Hodge, et al. (1998). "Logistic regression model for analyzing extended haplotype data." *Genet Epidemiol* **15**(2): 173-81.
- Wang, N., J. M. Akey, et al. (2002). "Distribution of recombination crossovers and the origin of haplotype blocks: the interplay of population history, recombination, and mutation." *Am J Hum Genet* **71**(5): 1227-34.
- Weir, B. S. and S. R. Wilson (1986). "Log-linear models for linked loci." *Biometrics* **42**(3): 665-70.
- Xing, E. P., M. I. Jordan, et al. (2007). "Bayesian haplotype inference via the Dirichlet process." *J Comput Biol* **14**(3): 267-84.
- Xu, H., X. Wu, et al. (2004). "Comparison of haplotype inference methods using genotypic data from unrelated individuals." *Hum Hered* **58**(2): 63-8.
- Zaykin, D. V., P. H. Westfall, et al. (2002). "Testing association of statistically inferred haplotypes with discrete and continuous traits in samples of unrelated individuals." *Hum Hered* **53**(2): 79-91.
- Zee, R. Y., N. R. Cook, et al. (2005). "Haplotype analysis of the beta2 adrenergic receptor gene and risk of myocardial infarction in humans." *Genetics* **169**(3): 1583-7.

