

Nested Hypotheses: An example in Genetics

Carlos Alberto de Bragança Pereira

Departamento de Estatística

Inst Mat & Estat: U São Paulo

Nested Hypotheses

- Most clinical trials today have the objective of deciding if a polymorphism is associated with a disease: That is, one is question if an allele a could be the cause of a the disease.
- Two groups, called case and control are observed to obtain the genotypic frequencies.

| Genotype | AA | Aa | aa | S. Size |
|----------|----------|----------|----------|---------|
| Case | x_{AA} | x_{Aa} | x_{aa} | n |
| Control | y_{AA} | y_{Aa} | y_{aa} | m |

Nested Hypotheses

- The theoretical population frequencies are:

| | | | |
|----------|---------------|---------------|---------------|
| Genotype | AA | Aa | aa |
| Case | γ_{AA} | γ_{Aa} | γ_{aa} |
| Control | π_{AA} | π_{Aa} | π_{aa} |

- Test of genotypic homogeneity: $\pi = \gamma$
- Test of allelic homogeneity:

$$Q(A) = \gamma_{AA} + .5\gamma_{Aa} = \pi_{AA} + .5\pi_{Aa} = P(A)$$

| Group | Observed frequency | | | |
|---------|--------------------|----|----|--------|
| | AA | AB | BB | sample |
| Case | 55 | 83 | 50 | 188 |
| Control | 24 | 42 | 39 | 105 |

| Group | ML estimates: genotypic homogeneity | | | |
|---------|-------------------------------------|------|------|-----|
| | AA | AB | BB | sum |
| Case | 0,27 | 0,43 | 0,30 | 1 |
| Control | 0,27 | 0,43 | 0,30 | 1 |

| Group | Expected frequency: genotypic homogeneity | | | |
|---------|---|-------|-------|--------|
| | AA | AB | BB | sample |
| Case | 50,69 | 80,20 | 57,11 | 188 |
| Control | 28,31 | 44,80 | 31,89 | 105 |

| | | | |
|-------------|------|------|---------------|
| 0,37 | 0,10 | 0,88 | 1,35 |
| 0,66 | 0,17 | 1,58 | 2,41 |
| Chi-squared | | | 3,76 |
| p-value | | | 15,24% |

| Group | Observed frequency | | | sample |
|---------|--------------------|----|----|--------|
| | AA | AB | BB | |
| Case | 55 | 83 | 50 | 188 |
| Control | 24 | 42 | 39 | 105 |

| Group | Allelic Frequency | | |
|---------|-------------------|-----|-----|
| | A | B | sum |
| Case | 193 | 183 | 376 |
| Control | 90 | 120 | 210 |

| Group | ML allelic estimates | | |
|---------|----------------------|------|-----|
| | A | B | sum |
| Case | 0,51 | 0,49 | 1 |
| Control | 0,43 | 0,57 | 1 |

| | | |
|------|------|------|
| 0,72 | 0,67 | 1,39 |
| 1,29 | 1,20 | 2,49 |

| Group | Allelic Exp. Frequency | | |
|---------|------------------------|--------|-----|
| | A | B | Sum |
| Case | 181,58 | 194,42 | 376 |
| Control | 101,42 | 108,58 | 210 |

Chi-square 3,87

p-value 4,91%

| Group | Frecuencia observada | | | |
|---------|----------------------|----|----|--------|
| | AA | AB | BB | sample |
| Case | 55 | 83 | 50 | 188 |
| Control | 24 | 42 | 39 | 105 |

| ML estimates; allelic homogeneity | | | | |
|-----------------------------------|------|------|------|-----|
| | AA | AB | BB | sum |
| Case | 0,26 | 0,44 | 0,30 | 1 |
| Control | 0,28 | 0,40 | 0,31 | 1 |

| Group | Expected frequency: allelic homogeneity | | | |
|---------|---|-------|-------|--------|
| | AA | AB | BB | sample |
| Case | 49,66 | 82,73 | 55,61 | 188 |
| Control | 29,72 | 42,25 | 33,04 | 105 |

| | | | |
|------|-------|------|------|
| 0.57 | 0.001 | 0.57 | 1.14 |
| 1.10 | 0.001 | 1.08 | 2.18 |

Chi-squar 3.32

p-value 6.85%

| Group | Observed frequency | | | |
|---------|--------------------|----|----|--------|
| | AA | AB | BB | sample |
| Case | 55 | 83 | 50 | 188 |
| Control | 24 | 42 | 39 | 105 |

| Group | Expected under HWE | | | |
|---------|--------------------|--------|--------|--------|
| | AA | AB | BB | sample |
| Case | 49.533 | 93.934 | 44.533 | 188 |
| Control | 19.286 | 51.429 | 34.286 | 105 |

| HWE | Alele Freq. | |
|---------|-------------|--------|
| | p(A) | p(B) |
| case | 0.5133 | 0.4867 |
| control | 0.4286 | 0.5714 |

Chi **p-value**
2.5470 **11.05%**
3.5292 **6.03%**

P_1 = population frequency of AA

P_2 = population frequency of Aa

P_3 = population frequency of aa

Next generation

Q_1, Q_2, Q_3 .

| Type of coupling | Type of offspring | | |
|------------------------|-------------------------|-------------------------|-------------------------|
| | AA | Aa | aa |
| AA x AA | $(P_1)(P_1)$ | 0 | 0 |
| AA x Aa | $(P_1)(P_2)/2$ | $(P_1)(P_2)/2$ | 0 |
| AA x aa | 0 | $(P_1)(P_3)$ | 0 |
| Aa x AA | $(P_1)(P_2)/2$ | $(P_1)(P_2)/2$ | 0 |
| Aa x Aa | $(P_2)(P_2)/4$ | $(P_2)(P_2)/2$ | $(P_2)(P_2)/4$ |
| Aa x aa | 0 | $(P_3)(P_2)/2$ | $(P_3)(P_2)/2$ |
| aa x AA | 0 | $(P_1)(P_3)$ | 0 |
| aa x Aa | 0 | $(P_3)(P_2)/2$ | $(P_3)(P_2)/2$ |
| aa x aa | 0 | 0 | $(P_3)(P_3)$ |
| Next Generation | Q_1 | Q_2 | Q_3 |

$$Q_1 = P_1^2 + 2 \frac{P_1 P_2}{2} + \frac{P_2^2}{4} = \left(P_1 + \frac{P_2}{2} \right)^2; \quad Q_3 = \left(P_3 + \frac{P_2}{2} \right)^2$$

$$Q_2 = 2 \frac{P_1 P_2}{2} + 2 \frac{P_3 P_2}{2} + 2 P_1 P_3 + 2 \frac{P_2^2}{2} = 2 \left(P_1 + \frac{P_2}{2} \right) \left(P_3 + \frac{P_2}{2} \right)$$

$$Q_1 = P_1^2 + 2\frac{P_1P_2}{2} + \frac{P_2^2}{4} = \left(P_1 + \frac{P_2}{2}\right)^2; \quad Q_3 = \left(P_3 + \frac{P_2}{2}\right)^2$$

$$Q_2 = 2\frac{P_1P_2}{2} + 2\frac{P_3P_2}{2} + 2P_1P_3 + 2\frac{P_2^2}{2} = 2\left(P_1 + \frac{P_2}{2}\right)\left(P_3 + \frac{P_2}{2}\right)$$

$$P_1 = p^2; \quad P_2 = 2p(1-p) \quad P_3 = (1-p)^2$$

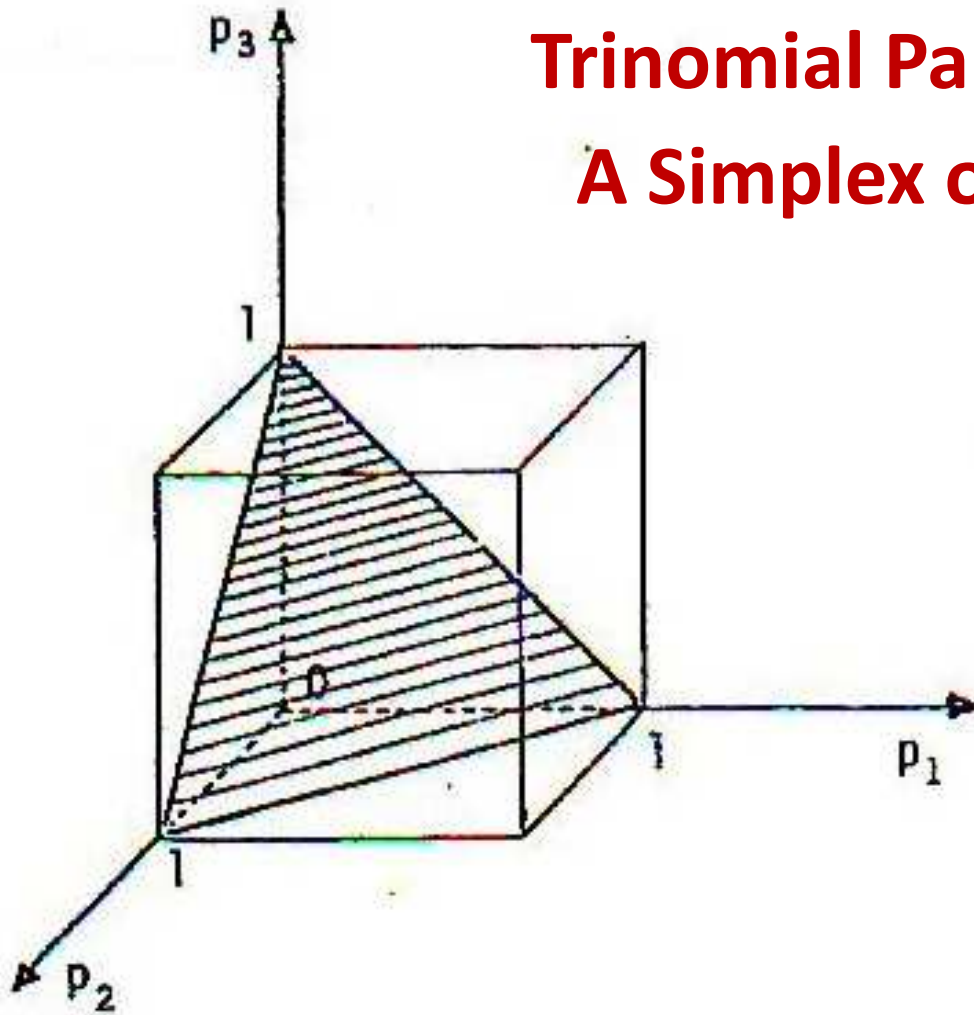
$$Q_1 = p^4 + 2p^3(1-p) + p^2(1-p)^2 = p^2; \quad Q_3 = (1-p)^2$$

$$Q_2 = 2p^3(1-p) + 4p^2(1-p)^2 + 2p(1-p)^3 = 2p(1-p)$$

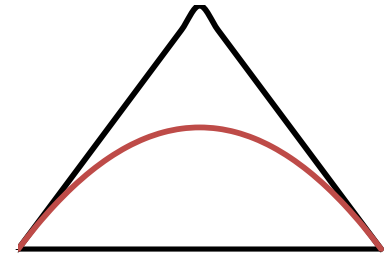
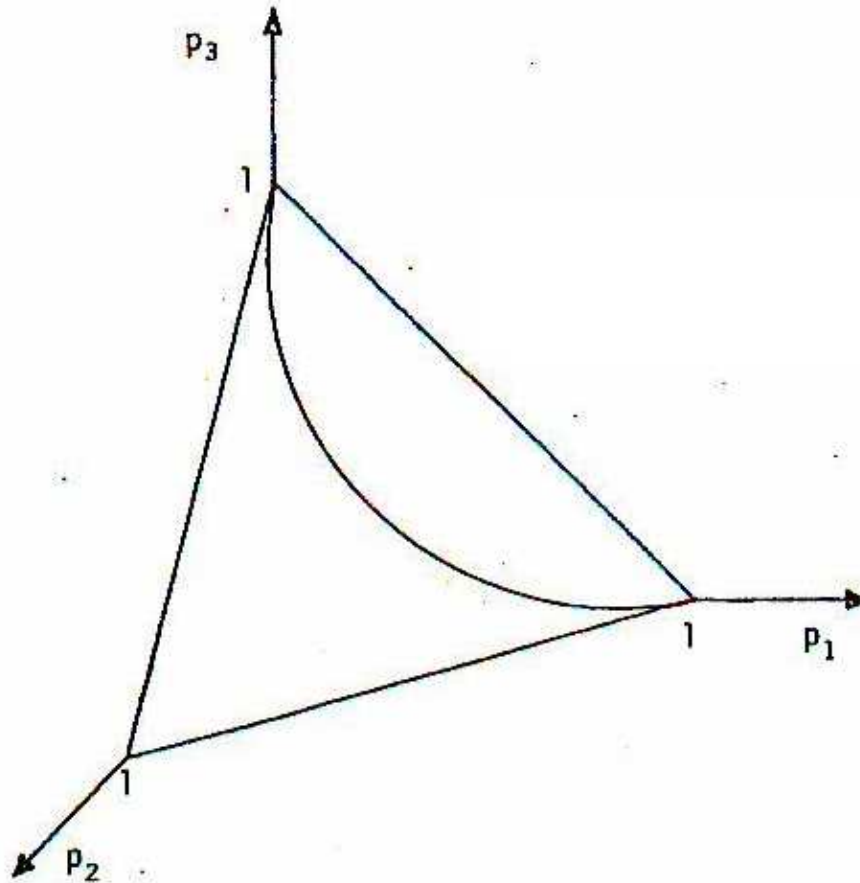
Note that $p = P_1 + \frac{1}{2}P_2 = 1 - \left(P_3 + \frac{1}{2}P_2\right)$

is the allelic frequency

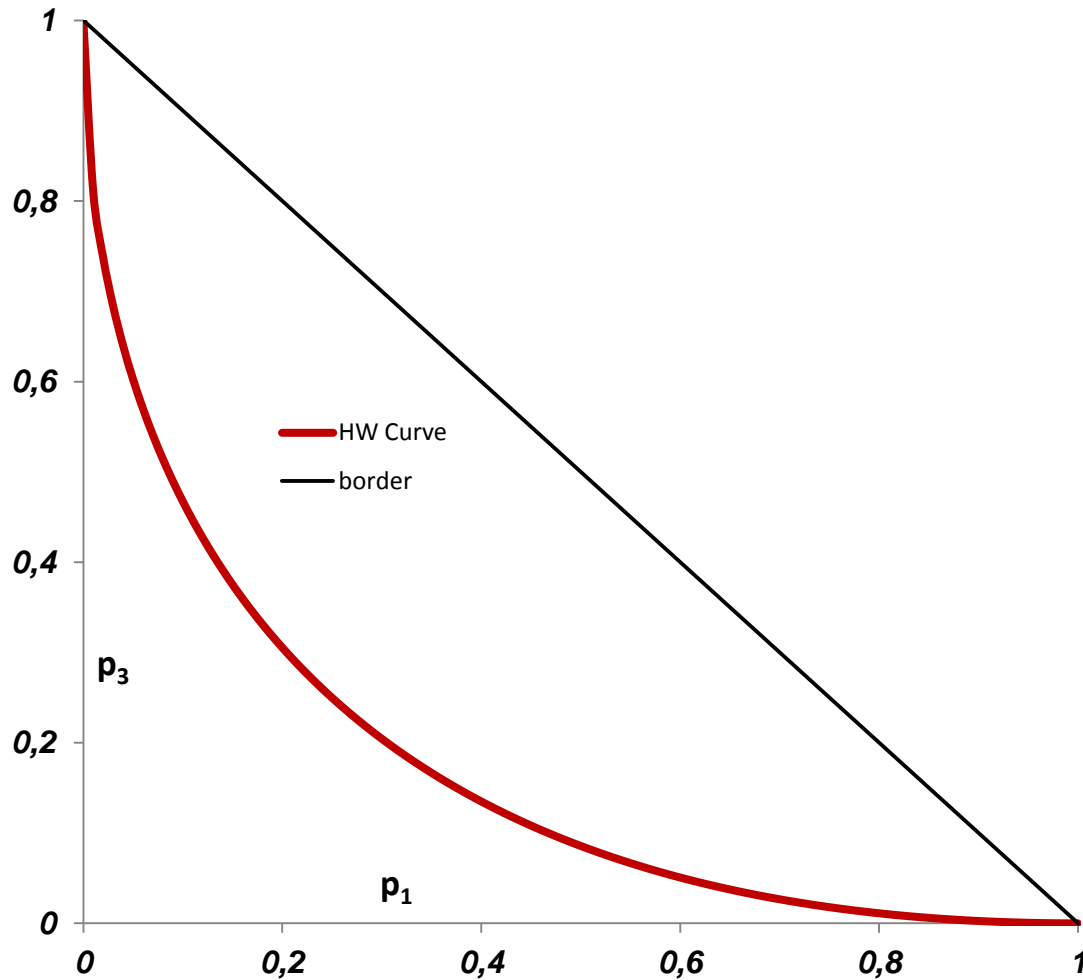
Trinomial Parameter Space: A Simplex of Dim 2 in \mathcal{R}^3



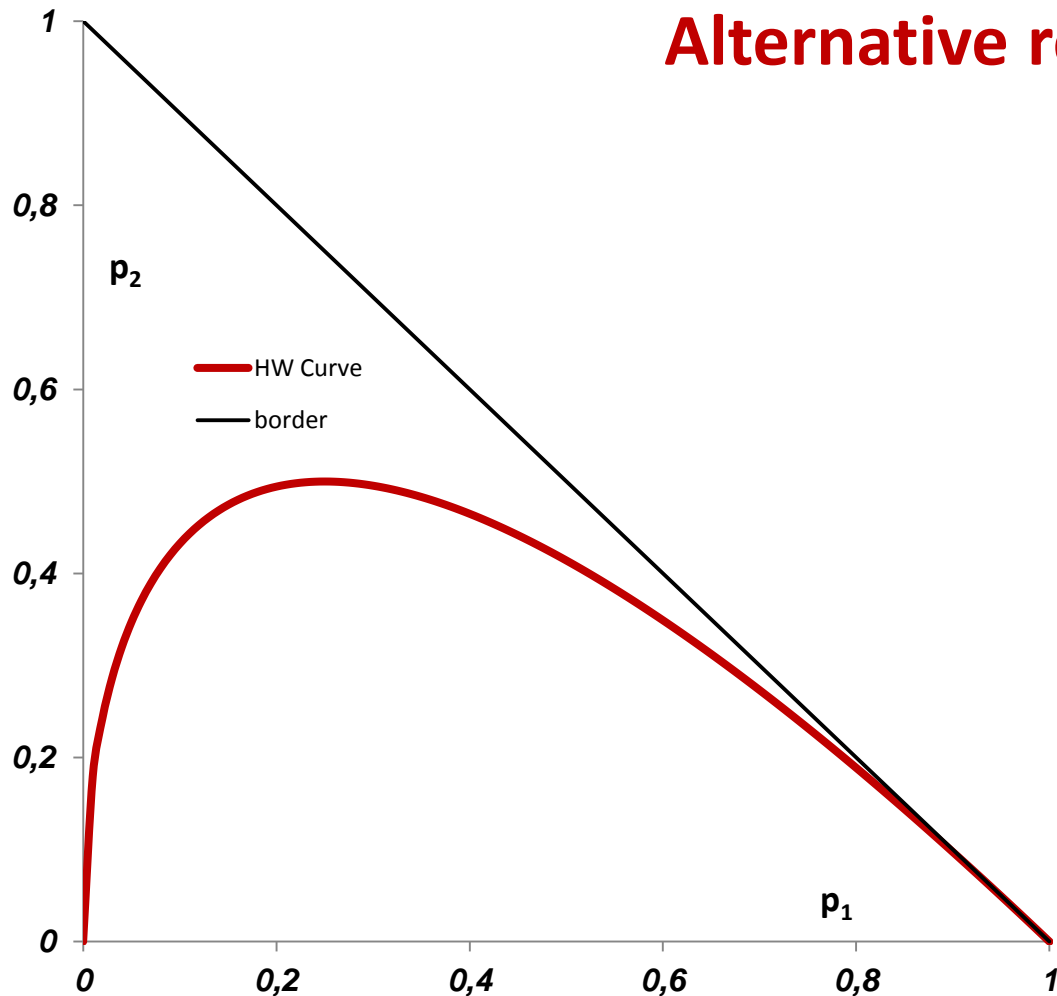
HW Equilibrium hypothesis subspace of Dim 1 in \mathcal{R}^3



HW Equilibrium hypothesis subspace of Dim 1 in \mathcal{R}^2



HW Equilibrium hypothesis subspace of Dim 1 in \mathfrak{R}^2 : Alternative representation



| aleles | A | B |
|---------------|--------------|--------------|
| A | P(AA) | P(AB) |
| B | P(BA) | P(BB) |

| aleles | A | B |
|---------------|-------------|-------------|
| A | p1 | .5p2 |
| B | .5p2 | p3 |

Association index for binary data: Disequilibrium Coefficient

Brentani H; Nakano EY; Martins CB; Izbicki R; Pereira CAB (2011), Disequilibrium coefficient: Bayesian perspective, *Statistical Applications in Genetics and Molecular Biology* 10(1): Article 22

Carlos Alberto de Bragança Pereira

Department of Statistic

Inst of Mat & Stat of U São Paulo

Hardy-Weimberg Law

- **Hardy-Weinberg law shows that in a Mendelian population, under certain restrictions, allele frequencies will be constant through generations; alternatively, been in HWE means that the relation of the proportions of the genotypes in a specific locus exhibit a special association between the alleles transmitted from parents.**

MATERIAL & METHODS

- In order to see if a population is in HWE one takes a sample of size $n=n_1+n_2+n_3$ and observe n_1, n_2, n_3 , the absolute frequencies of the genotypes **AA, Aa, aa**, respectively. Let π_1, π_2, π_3 be the population frequencies of these genotypes, with $\pi_1+\pi_2+\pi_3=1$ and $\pi_i \geq 0, i=1,2,3$. If we consider genotypes from different individuals to be statistically independent, likelihood function can be expressed as

Likelihood

$$L(\pi_1, \pi_2, \pi_3 | \mathbf{n}_1, \mathbf{n}_2, \mathbf{n}_3) \propto \pi_1^{\mathbf{n}_1} \pi_2^{\mathbf{n}_2} \pi_3^{\mathbf{n}_3}$$

with \propto denoting proportionality. Note that the parametric space is

$$\Theta = \{(\mathbf{p}_1, \mathbf{p}_2, \mathbf{p}_3) : \mathbf{p}_1 \geq \mathbf{0} \wedge \mathbf{p}_2 \geq \mathbf{0} \wedge \mathbf{p}_3 \geq \mathbf{0} \wedge \mathbf{p}_1 + \mathbf{p}_2 + \mathbf{p}_3 = \mathbf{1}\}$$

- HWE holds if, and only if, there exists $0 < \pi < 1$ such that $\pi_1 = \pi^2$, $\pi_2 = 2\pi(1-\pi)$, $\pi_3 = (1-\pi)^2$

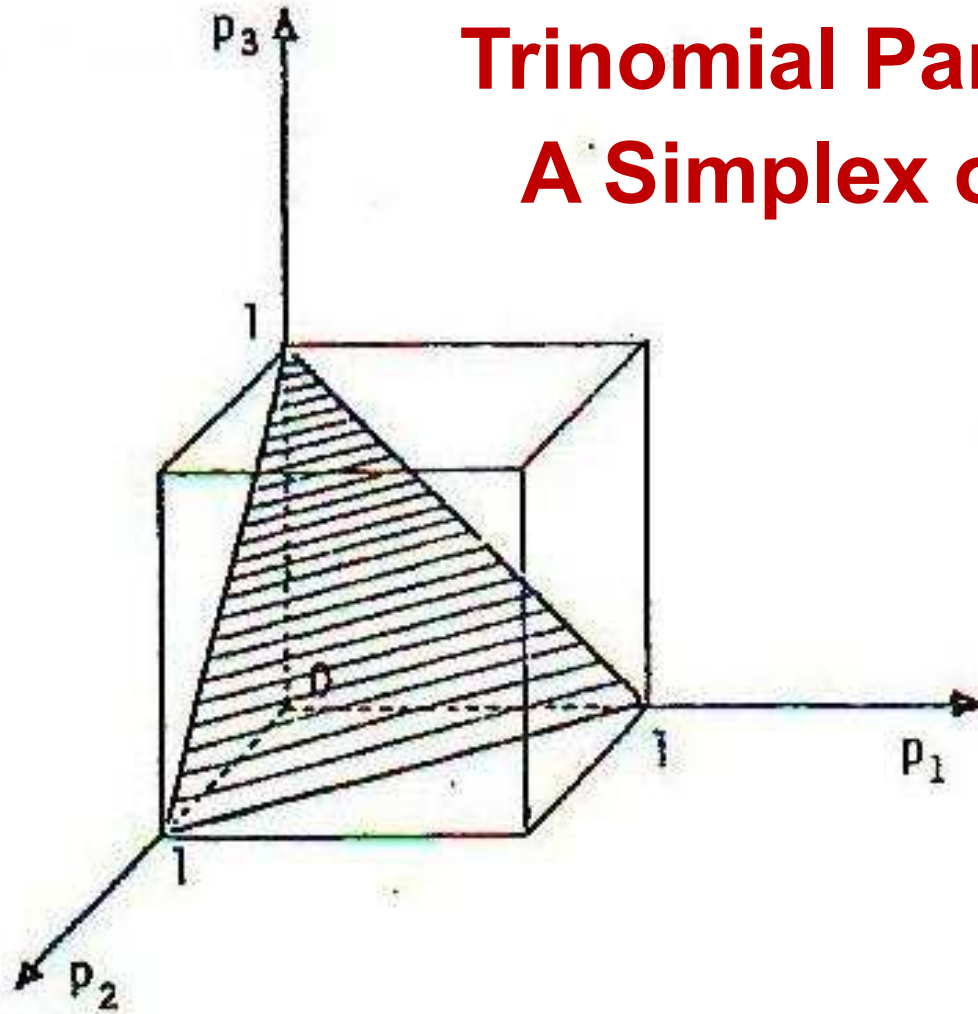
- **The HWE null hypothesis is:**

$H: \theta \in \Theta_H$ with $\theta = (\pi_1; \pi_2; \pi_3)$ &

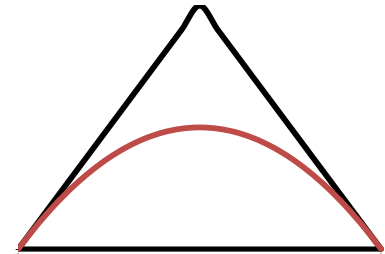
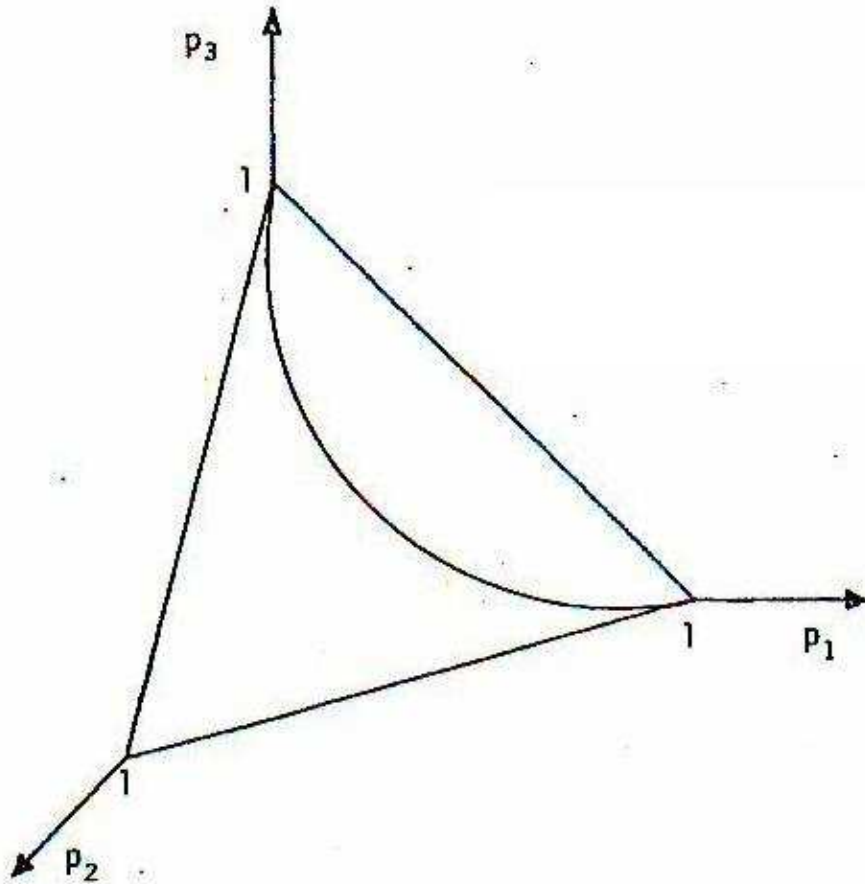
$$\Theta_H = \left\{ (\mathbf{p}_1, \mathbf{p}_2, \mathbf{p}_3) : \exists \mathbf{p} \in [0;1] : \left(\mathbf{p}_1 = \mathbf{p}^2 \right) \wedge \left(\mathbf{p}_3 = (1 - \mathbf{p})^2 \right) \right\} \subset \Theta.$$

$$\Theta = \left\{ (\mathbf{p}_1, \mathbf{p}_2, \mathbf{p}_3) : \forall i; \mathbf{p}_i \geq \mathbf{0} \wedge \mathbf{p}_1 + \mathbf{p}_2 + \mathbf{p}_3 = \mathbf{1} \right\}$$

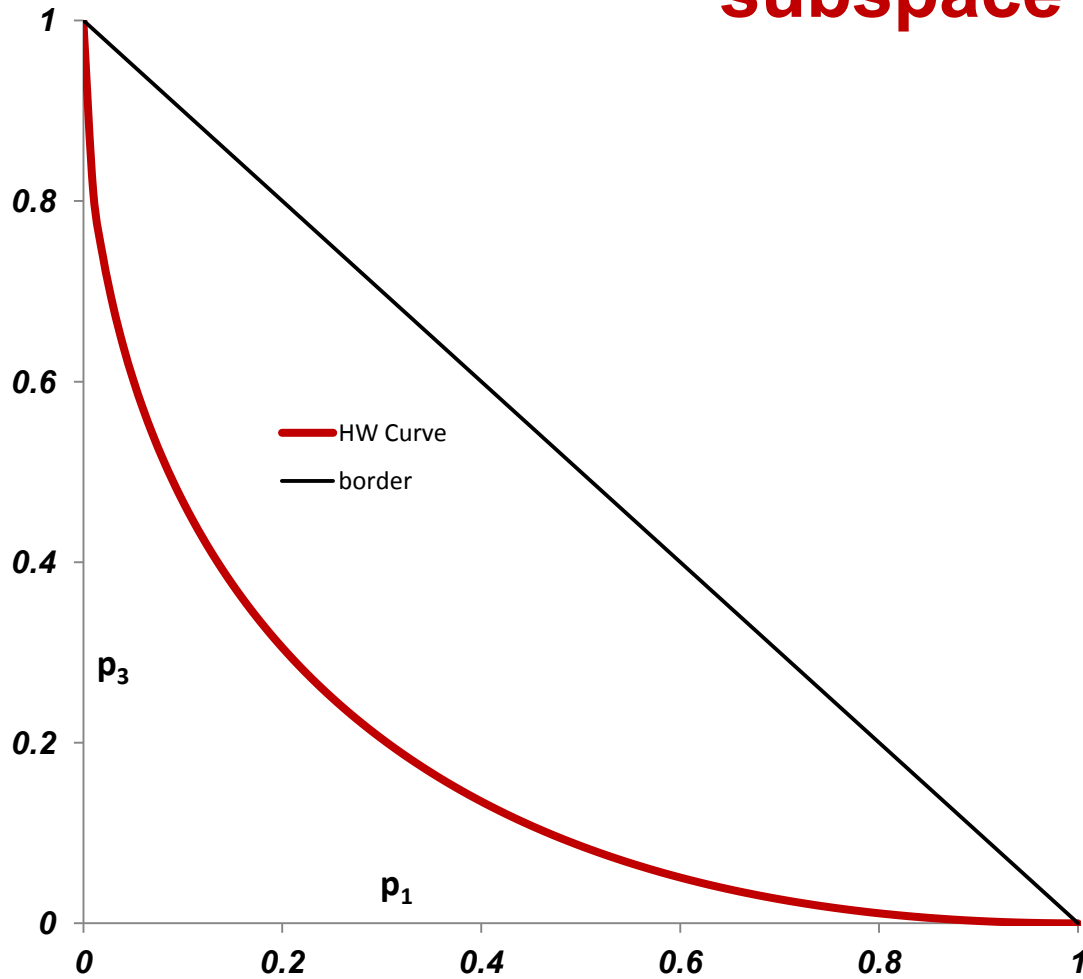
Trinomial Parameter Space: A Simplex of Dim 2 in \mathfrak{R}^3



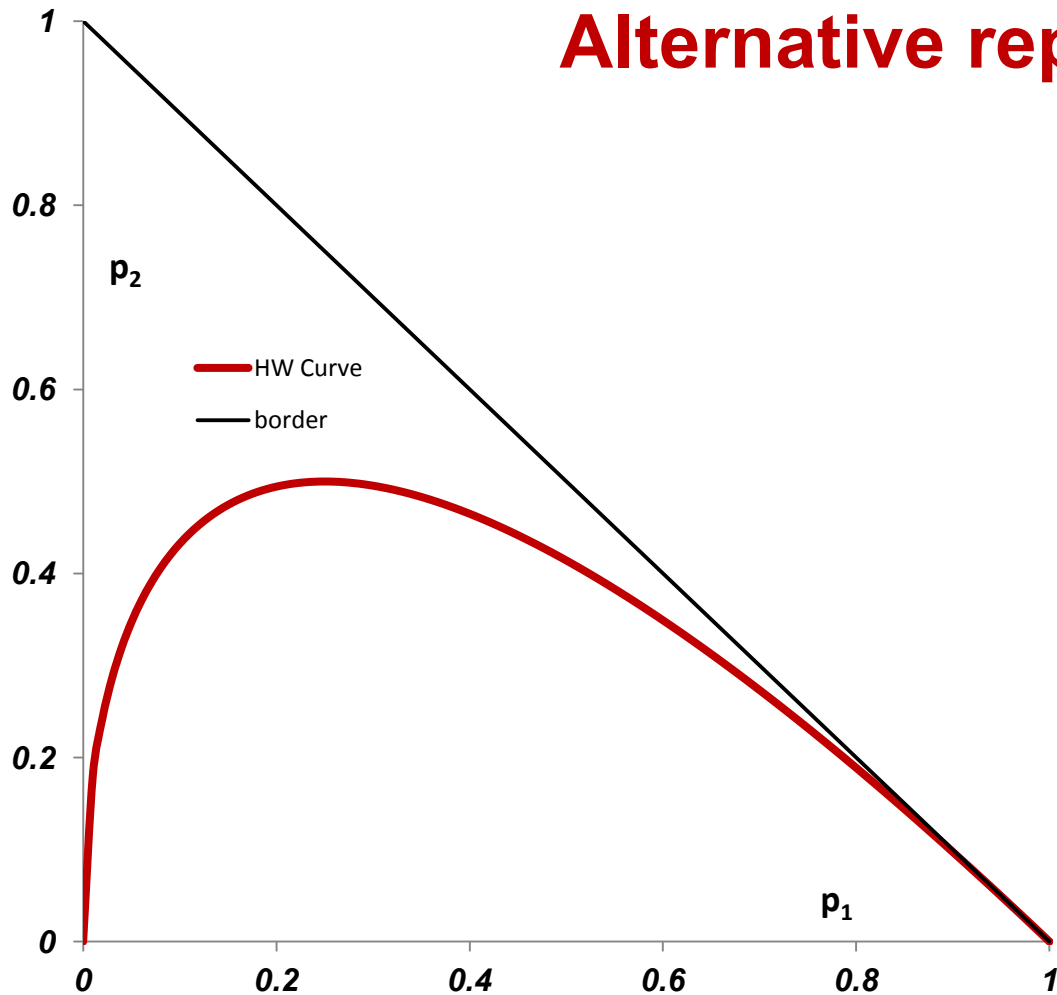
HW Equilibrium hypothesis subspace of Dim 1 in \mathcal{R}^3



HW Equilibrium hypothesis subspace of Dim 1 in \mathfrak{R}^2



HW Equilibrium hypothesis subspace of Dim 1 in \mathfrak{R}^2 : Alternative representation



Conjugate Bayesian Operation: Prior/Posterior

$$f(\theta | \mathbf{d}) = \Gamma(n + \mathbf{a}) \prod_{i=1}^3 \frac{\pi_i^{(n_i + a_i - 1)}}{\Gamma(n_i + a_i)} \quad \text{for } \theta \in \Theta$$

Taking $\mathbf{A}_i = n_i + a_i$ for $(a_1; a_2; a_3)$ the prior parameter:

$$\hat{\theta} = \mathbf{E}(\theta | \mathbf{d}) = \left(\frac{\mathbf{A}_1}{\mathbf{A}}, \frac{\mathbf{A}_2}{\mathbf{A}}, \frac{\mathbf{A}_3}{\mathbf{A}} \right) = (\hat{\theta}_1, \hat{\theta}_2, \hat{\theta}_3) \quad \&$$

$$\Sigma = \frac{\mathbf{1}}{\mathbf{A} + \mathbf{1}} \left(\left(\begin{array}{ccc} \hat{\theta}_1 & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \hat{\theta}_2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \hat{\theta}_3 \end{array} \right) - \left(\begin{array}{c} \hat{\theta}_1 \\ \hat{\theta}_2 \\ \hat{\theta}_3 \end{array} \right) \left(\hat{\theta}_1, \hat{\theta}_2, \hat{\theta}_3 \right) \right)$$

Disequilibrium Coefficient

- The disequilibrium coefficient, denoted by λ , is now defined and will be explained and motivated in the next section. In fact this coefficient is a modification of the one presented by Pereira and Rogatko (1984). This new is based on the correlation coefficient proposed by Yule (1912) to measure association of two binary variables.

Definition

- Since we have the complete specification of the parameter θ and the fact that λ is simply a function of θ , we may obtain also the specification of the distribution of λ .
- **Definition: The Hardy-Weinberg disequilibrium coefficient is defined as follows:**

$$\lambda = \frac{\sqrt{\pi_1 \pi_3} - \frac{1}{2} \pi_2}{\sqrt{\pi_1 \pi_3} + \frac{1}{2} \pi_2} \in [-1;1]$$

Motivation

Dirichlet prior parameter =

$(a_1; a_{12}; a_{21}; a_3)$, data = $d =$

$(n_1; n_{12}; n_{21}; n_3)$, & Dirichlet

posterior parameter =

$(A_1; A_{12}; A_{21}; A_3)$; $A_i = a_i + n_i$

Table 1: Frequency data, n , (parameters, π) in a multinomial classification

| $t \setminus g$ | g | G |
|-----------------|---------------------|---------------------|
| t | $n_1 (\pi_1)$ | $n_{12} (\pi_{12})$ |
| T | $n_{21} (\pi_{21})$ | $n_3 (\pi_3)$ |

Cross product ratio is $\psi = \frac{\pi_1 \pi_3}{\pi_{12} \pi_{21}}$ with

$$E(\psi | d) = \frac{A_1 A_2}{(A_{12} - 1)(A_{21} - 1)} \quad \& \quad V(\psi | d) = \frac{A_1 A_3 (A_1 + A_{12} - 1)(A_3 + A_{21} - 1)}{(A_{12} - 1)(A_{21} - 1)(A_{12} - 2)(A_{21} - 2)}$$

Yule's Association Coefficient

- Yule (1912) understood that although ψ could be considered an association coefficient, it is unbounded & unbalanced: negative association occurs if $\psi \in (0;1)$, independence if $\psi=1$, and positive association if $\psi \in (1; \infty)$. Hence, defined

$$\dot{\lambda} = \frac{\sqrt{\pi_1\pi_3} - \sqrt{\pi_{12}\pi_{21}}}{\sqrt{\pi_1\pi_3} + \sqrt{\pi_{12}\pi_{21}}} = \frac{\sqrt{\psi} - 1}{\sqrt{\psi} + 1} \in (-1; +1)$$

First Example

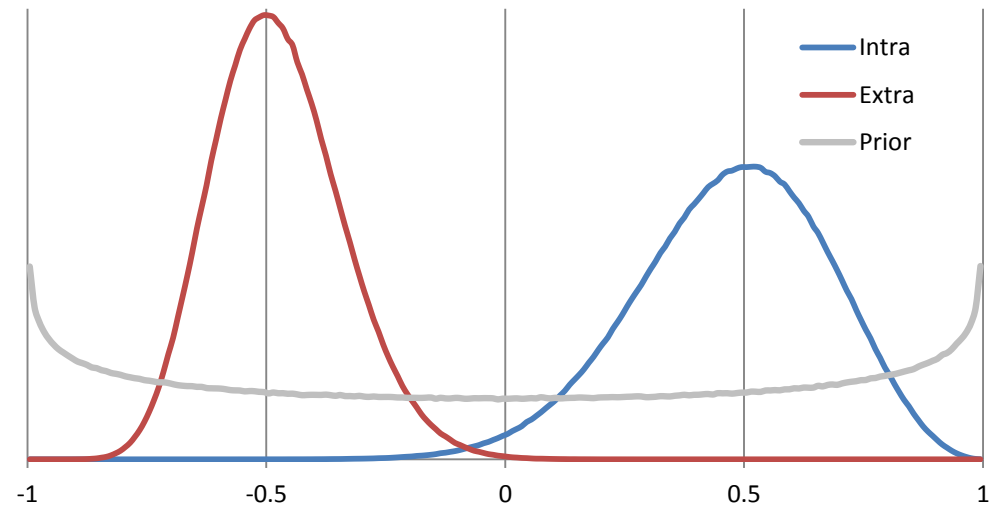
TABLE 2: Outcomes of the diagnostic tests for 93 children

| Results | Extra-hepatic | | | Intra-hepatic | | | Total |
|----------------|---------------|-----------|-----------|---------------|-----------|-----------|-----------|
| | E^+ | E^- | Sum | E^+ | E^- | Sum | |
| $\varepsilon+$ | 5 | 9 | 14 | 28 | 12 | 40 | 54 |
| $\varepsilon-$ | 28 | 6 | 34 | 1 | 4 | 5 | 39 |
| Sum | 33 | 15 | 48 | 29 | 16 | 45 | 93 |

Jeffrey's prior

$$(\pi_1; \pi_{12}; \pi_{21}; \pi_3) \approx D\left(\frac{1}{2}; \frac{1}{2}; \frac{1}{2}; \frac{1}{2}\right)$$

Figure 1: Jeffreys' prior and posterior densities for extra- & intra-hepatic groups



Statistics

$$\text{Bayes Estimates: } E(\lambda | d) = \begin{cases} -0.4750 & \text{for Extra - Hepatic} \\ 0.4723 & \text{for Intra - Hepatic} \end{cases}$$

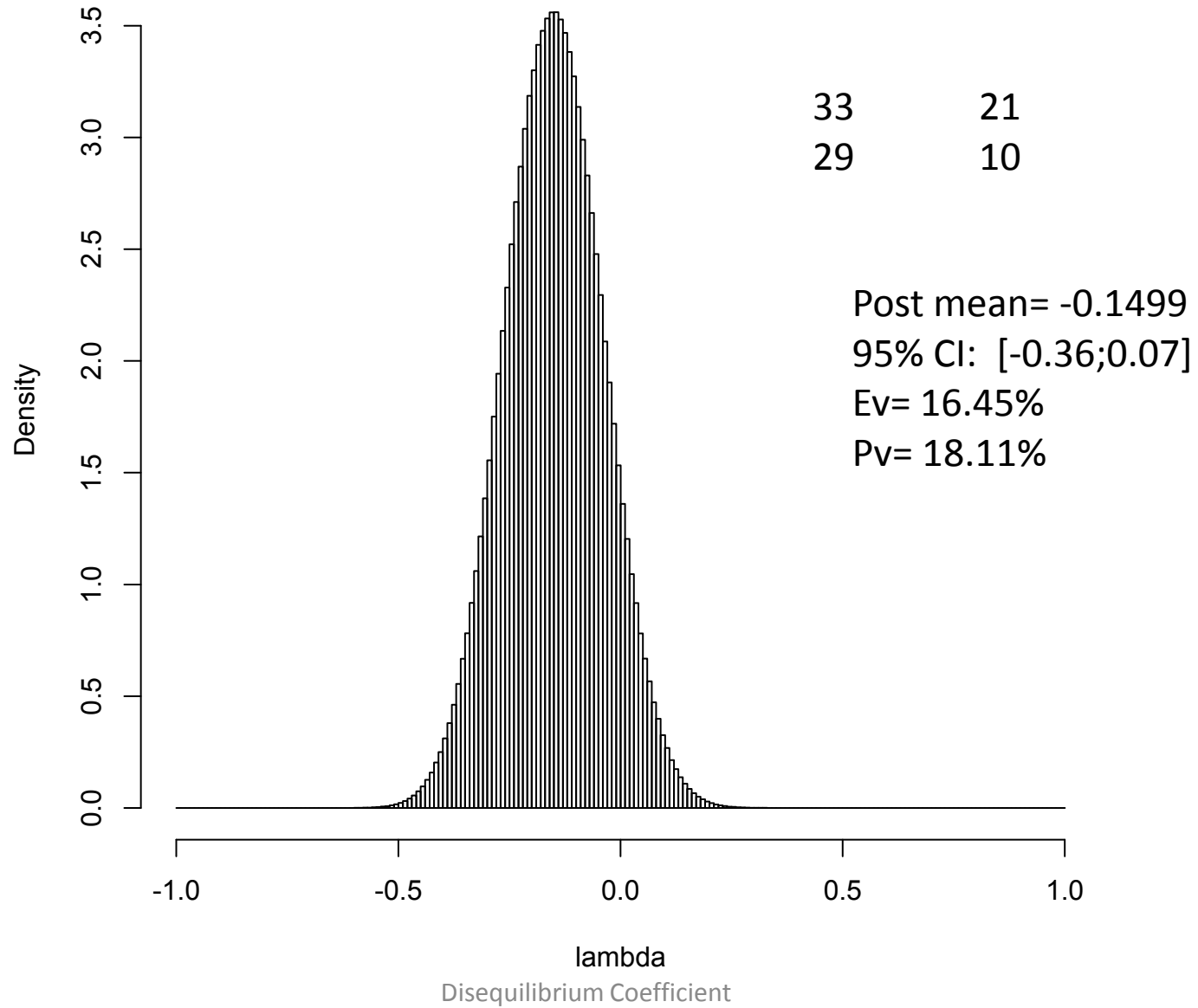
MLEstimates : for EH is (-.49) and for IH is (.51).

$$95\% \text{ credible sets : } \begin{cases} (-.72; -.21) & \text{for Extra - Hepatic} \\ (+.08; +.85) & \text{for Intra - Hepatic} \end{cases}$$

$$\text{Significance Index: Bayes (ev)} \begin{cases} 0.13\% & \text{for Extra - Hepatic} \\ 1.96\% & \text{for Intra - Hepatic} \end{cases}$$

$$\text{Significance Index: } \chi^2 \text{ (pv)} \begin{cases} 0.12\% & \text{for Extra - Hepatic} \\ 1.81\% & \text{for Intra - Hepatic} \end{cases}$$

Posterior of lambda



Disequilibrium coefficient based in the Yule's association coefficient

| | | Father's Allel | |
|----------------|---|--------------------|--------------------|
| | | A | D |
| Mother's Allel | A | $n_1(\pi_1)$ | $n_{12}(\pi_{12})$ |
| | D | $n_{21}(\pi_{21})$ | $n_3(\pi_3)$ |

Having the secondary diagonal cells be confounded in one, we replace the table by

| | | Father's Allel | |
|----------------|---|------------------|------------------|
| | | A | D |
| Mother's Allel | A | $n_1(\pi_1)$ | $.5n_2(.5\pi_2)$ |
| | D | $.5n_2(.5\pi_2)$ | $n_3(\pi_3)$ |

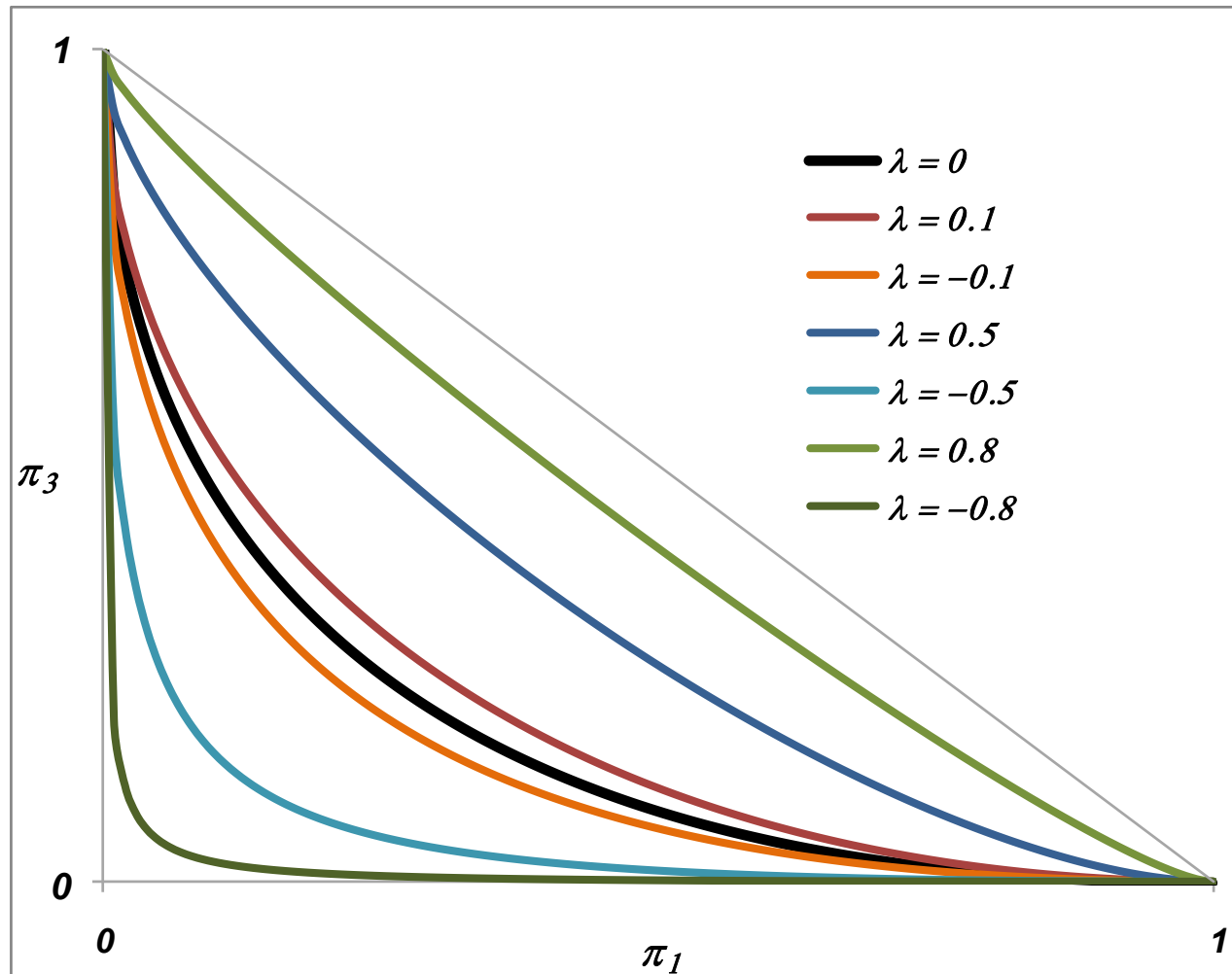
$$n_2 = n_{12} + n_{21} \quad \& \quad \pi_2 = \pi_{12} + \pi_{21}$$

$$\lambda = \frac{\sqrt{\pi_1\pi_3} - \frac{1}{2}\pi_2}{\sqrt{\pi_1\pi_3} + \frac{1}{2}\pi_2} =$$

Pereira & Rogatko (1984) introduced $\rho = \psi^{-1}$.

$$= \frac{\sqrt{\rho^{-1}} - 1}{\sqrt{\rho^{-1}} + 1} \in [-1; 1] \text{ for } \rho = \frac{\pi_2^2}{4\pi_1\pi_3}$$

Disequilibrium Curves



Molecular Biology Examples

- **Using data from published studies we calculated the disequilibrium coefficient and plotted the posterior density curves. To date, the APOE- $\epsilon 4$ is the only recognized risk factor to Alzheimer Disease (AD). It is well established that whenever one carries a $\epsilon 4$ allele, his AD risk increases in an allele dose dependent manner.**
- **Previous reports have suggested that additional factors within the APOE locus, in other genes and from the environment might also modulate risk.**
- **Table 3 presents all statistical results and looking at Figure 3a one should see that controls are in HWE. As expected, the density curve includes the zero with high density (tails having large areas). On the other hand, the density curve for the disease state seems to be out of HWE since the zero has low frequency (tails having small areas).**

FIGURE 3 – Posterior Density distributions of cases and controls genotyped for SNPs in APOE gene (a, b, c) and IL6 gene (d).

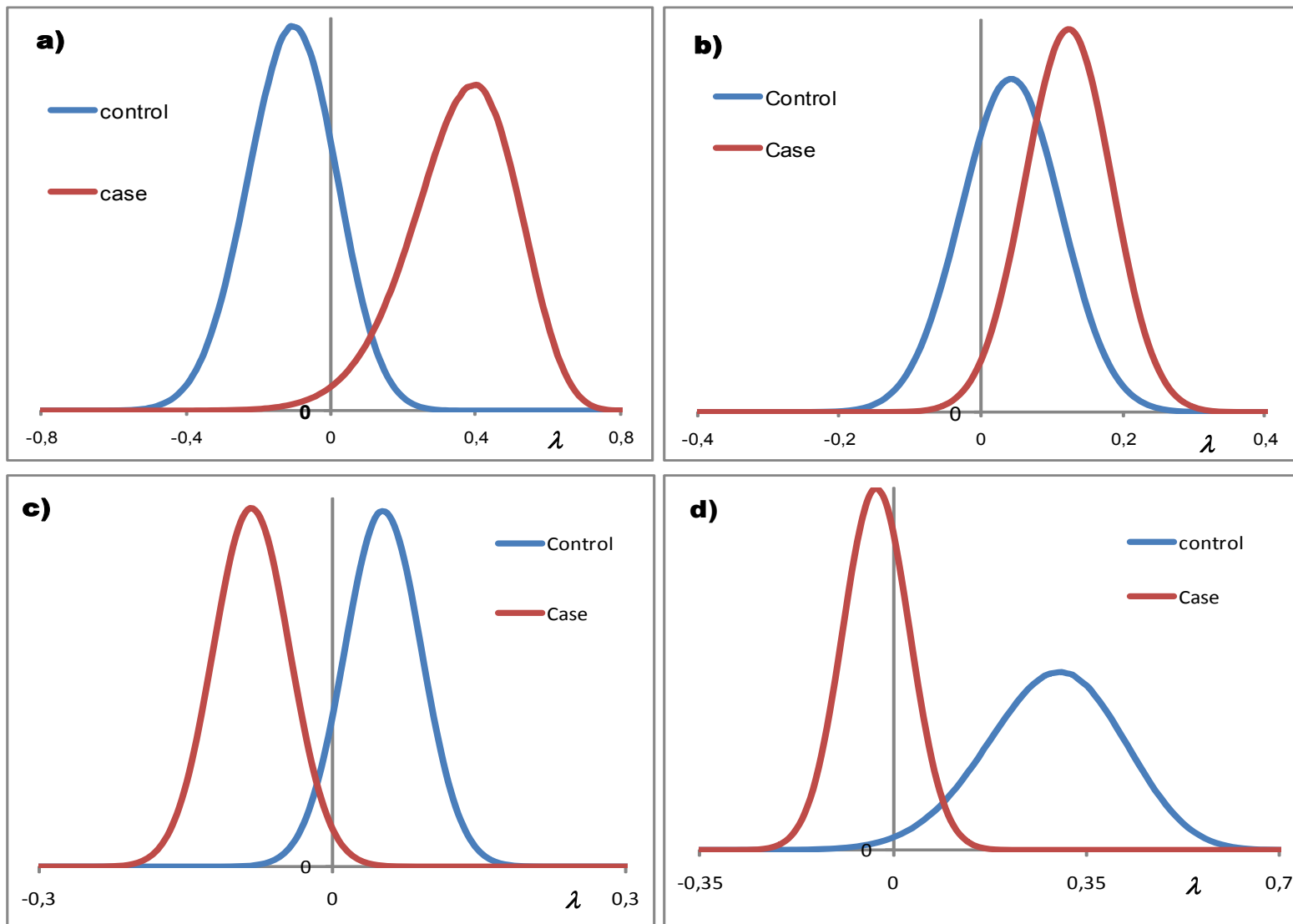


TABLE 3 – Statistical end-points for evaluating disequilibrium.

| study | sample | Fenotype | | | λ estimate | | 95% Credible | | e-value | p-value |
|------------|---------|----------|-----|-----|--------------------|--------|--------------|--------|---------|----------|
| | | AA | AD | DD | Bayes | MLE | LimInf | LimSup | FBST | χ^2 |
| AD 121 (1) | case | 4 | 18 | 94 | 0.357 | 0.366 | 0.066 | 0.649 | 0.016 | 0.018 |
| | control | 6 | 53 | 74 | -0.112 | -0.114 | -0.352 | 0.128 | 0.361 | 0.362 |
| AD 121 (2) | case | 57 | 118 | 100 | 0.122 | 0.123 | 0.003 | 0.242 | 0.045 | 0.046 |
| | control | 58 | 97 | 48 | 0.042 | 0.042 | -0.095 | 0.179 | 0.550 | 0.550 |
| AD 121 (3) | case | 120 | 361 | 194 | -0.084 | -0.084 | -0.16 | -0.007 | 0.032 | 0.032 |
| | control | 206 | 309 | 142 | 0.051 | 0.051 | -0.026 | 0.128 | 0.198 | 0.197 |
| AD 108 (4) | case | 110 | 148 | 44 | -0.031 | -0.031 | -0.149 | 0.088 | 0.611 | 0.611 |
| | control | 34 | 22 | 12 | 0.290 | 0.295 | 0.051 | 0.529 | 0.018 | 0.022 |

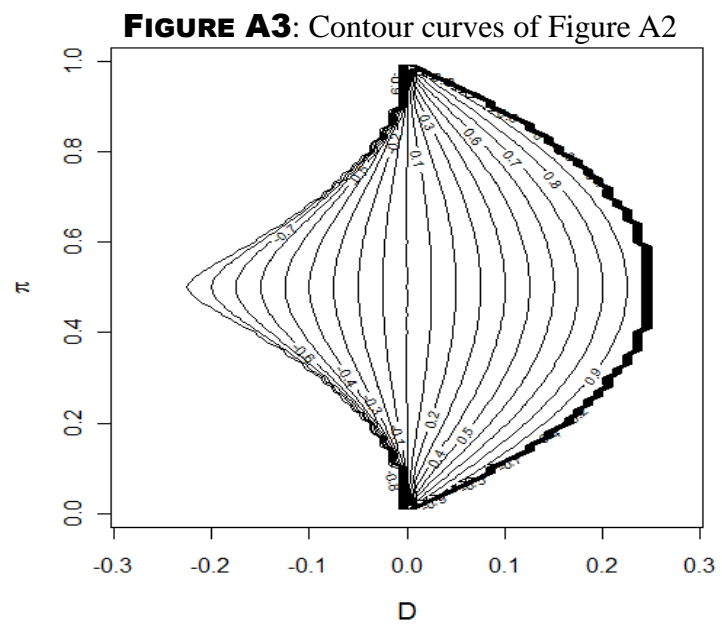
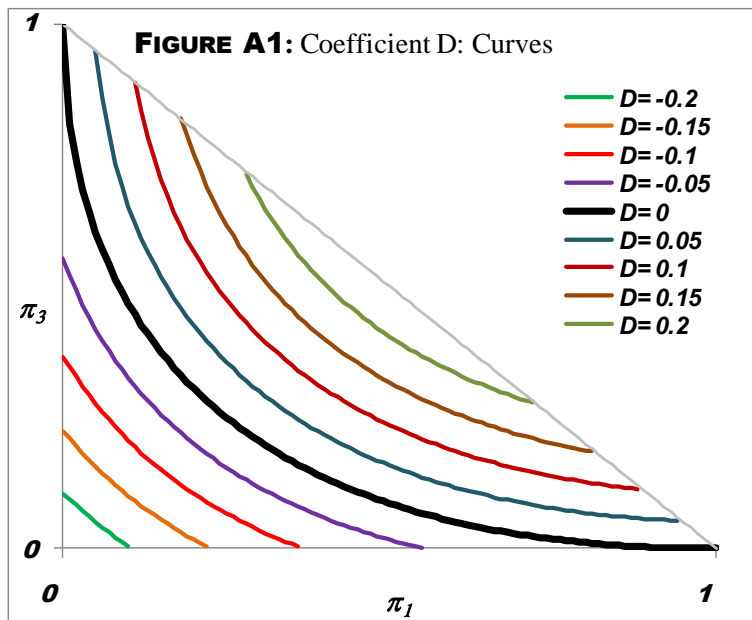
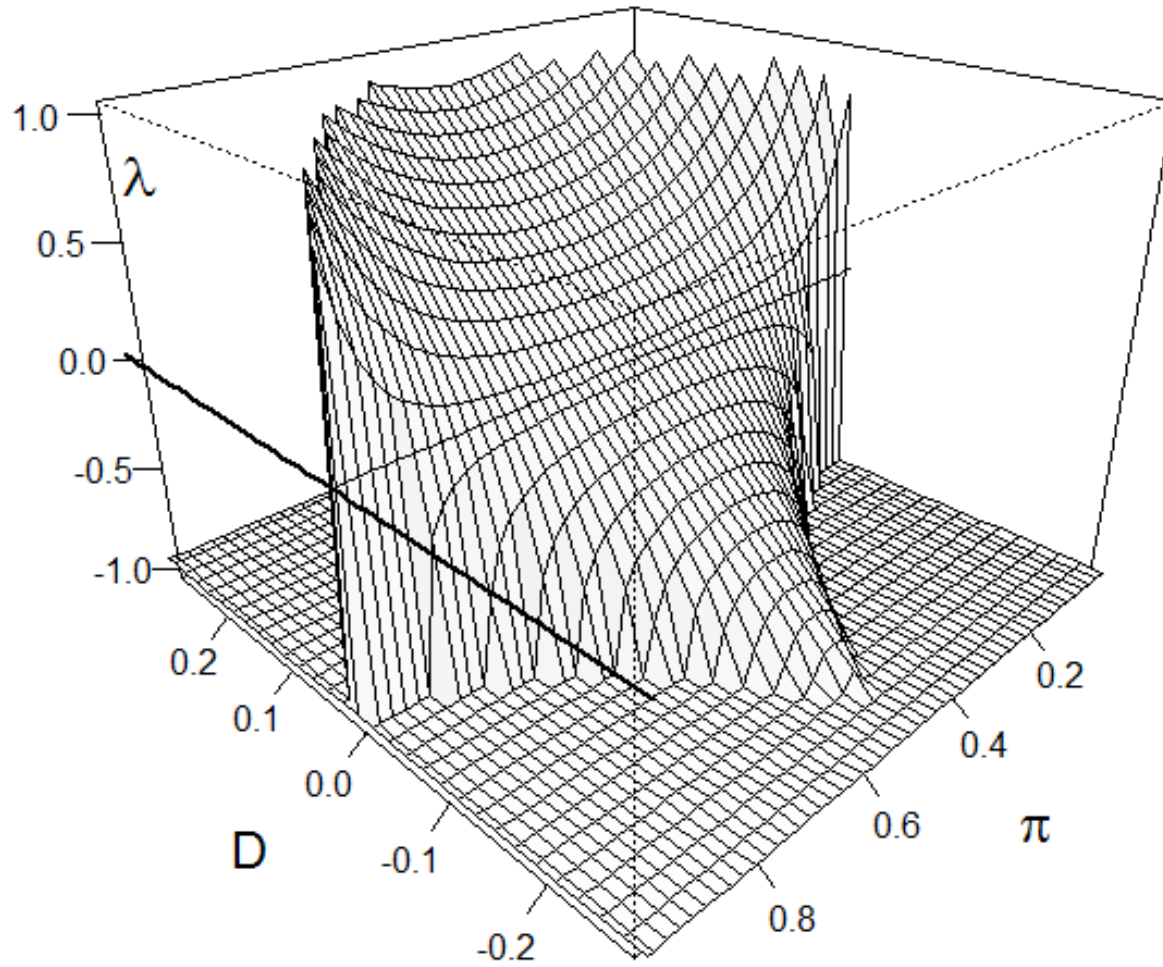


FIGURE A2: λ as a function of (π, D)



HW Bibliography

- 1908 Hardy GH, Mendelian proportions in a mixed population, *Science* 1908; **28**:49-50.
- 1908 Weinberg W, Über den Nachweis der Vererbung beim Menschen, *Jahreshefte des Vereins für vaterländische Naturkunde in Württemberg* 1908; **64**:369-382.
- 1980 Emigh TH, A comparison of tests for Hardy-Weinberg equilibrium, *Biometrics* 1980; **36**: 627-42.
- 1984 Pereira CAB, Rogatko A. The Hardy–Weinberg equilibrium law under a Bayesian perspective. *Brazilian Journal of Genetics* 1984; **4**:689-707.
- 1987 Louis EJ, Dempster ER. An exact test for Hardy–Weinberg and multiple alleles. *Biometrics* 1987; **43**:805-11.
- 1988 Lindley DV. Statistical inference concerning Hardy–Weinberg equilibrium. In *Bayesian Statistics 3*, Bernardo JM, DeGroot MH, Lindley DV, Smith AFM (eds). Clarendon Press [Oxford University Press]: Oxford, U.K., 1988; 307-26.
- 1989 Hernández JL and Weir BS, A disequilibrium coefficient approach to Hardy-Weinberg testing, *Biometrics* 1989; **45**:53-70.
- 1991 Singer JM, Peres CA and Harle CE, A note on the Hardy-Weinberg equilibrium in generalized ABO systems, *Statistics and Probability Letters* 1991; **11**:173-175.
- 1992 Chow M and Fong DKH, Simultaneous estimation of the Hardy-Weinberg proportions, *Canadian J Statistics* 1992; **20**:291-6.
- 1992 Guo SW, Thompson EA. Performing the exact test of Hardy–Weinberg proportion for multiple alleles. *Biometrics* 1992; **48**:361-72.
- 1998 Ayres KL and Balding DJ, Measuring departures from Hardy-Weinberg: A Markov chain Monte Carlo method for estimating the inbreeding coefficient, *Heredity* 1998; **80**:769-77.
- 1998 Shoemaker J, Painter I, Weir BS. A Bayesian characterization of Hardy–Weinberg disequilibrium. *Genetics* 1998; **149**:2079-88.

- Pereira CAB and Stern JM (1999), Evidence and credibility: full Bayesian significance test of precise hypothesis, *Entropy* 1:99-110
- 2000 Rogatko A, Slifker MJ and Babb JS, Hardy-Weinberg equilibrium diagnostics, *Theoretical Population Biology* 2000; **62**:251-7.
- 2001 Montoya-Delgado LE, Irony TZ, Pereira CAB, Whittle MR. An unconditional exact test for the Hardy-Weinberg equilibrium law: Sample-space ordering using the Bayes factor. *Genetics* 2001; **158**:875-83.
- 2004 Hosking L, Lumsden S, Lewis K, Yeo A, McCarthy L, Bansal A, Riley J, Purvis I, Xu CF (2004), Detection of genotyping errors by Hardy-Weinberg equilibrium testing, *European J Human Genetics* 2004; **12**(5):395-
- 2005 Wigginton J, Cutler D, Abecasis G. A note on exact tests of Hardy-Weinberg equilibrium. *American Journal of Human Genetics* 2005; **76**:887-93.
- 2006 Huber M, Chen Y, Dinwoodie I, Dobra A, Nicholas M. Monte Carlo algorithms for Hardy-Weinberg proportions. *Biometrics* 2006; **62**:49-53.
- 2006 Pereira CAB, Nakano F, Stern JM, Martin WR. Genuine Bayesian multiallelic significance test for the Hardy-Weinberg equilibrium law. *Genetics and Molecular Research* 2006; **5**:619-31.
- 2008 Consonni G, Gutiérrez-Peña E, Veronese P. Compatible priors for Bayesian model comparison with an application to the Hardy-Weinberg equilibrium model. *Test* 2008; **17**:585-605.
- 2008 Kelli Ryckman and Scott M Williams, Calculation and Use of the Hardy-Weinberg Model in Association Studies, *Current Protocols in Human Genetics* 2008; **57**:1.18.1-1.18.11.
- 2009 Lauretto MS, Nakano F, Faria SRJ, Pereira CAB, Stern JM. A straightforward multiallelic significance test for the Hardy-Weinberg equilibrium law. *Genetics and Molecular Biology* 2009; **32**:619-25.
- 2009 Wakefield J. Bayesian methods for examining Hardy-Weinberg equilibrium. *Biometrics* 2009; **66**:257-
- 2010 Attia J, Thakkinstian A, McElduff P, Milne E, Dawson S, Scott RJ, Klerk N, Armstrong B and Thompson J (2010), Detecting genotyping error using measures of degree of Hardy-Weinberg disequilibrium, *Statistical Applications in Genetics and Molecular Biology* 2010; **9**(1):5.
- 2011 Consonni G, Moreno E & Venturini S, Testing Hardy-Weinberg equilibrium: An objective Bayesian analysis, *Statistics in Medicine* 2010 (on line); DOI: 10.1002/sim.4084.

ADDITIONAL REFERENCES

- Bhojak TJ, DeKosky ST, Ganguli M, Kamboh MI (2000), Genetic polymorphisms in the cathepsin D and interleukin-6 genes and the risk of Alzheimer's disease, ***Neuroscience Letters*** 288(1):21-4.
- Jeffreys H (1931). ***Scientific Inference***, Cambridge University Press.
- Lambert JC et al (2002), Contribution of APOE promoter polymorphisms to Alzheimer's disease risk, ***Neurology*** 59(1):59-66.
- Pereira CAB and Stern JM (2008), Especial characterizations of standard discrete models, ***Statistical Journal*** 6(3):199-230
- Pereira CAB, Stern JM and Wechsler S (2008), Can a significance test be genuinely Bayesian? ***Bayesian Analysis*** 3(1):79-100
- R Development Core Team (2009), R: A language and environment for statistical computing, ***R Foundation for Statistical Computing***, Vienna, Austria, ISBN 3-900051-07-0, <http://www.R-project.org>.
- Weir BS (1996), ***Genetic data analysis II - methods for discrete population genetic data***, Sinauer Associates, Sunderland.