

Joint modeling of linkage and association for binary traits

J Concepción Loredo-Osti

joint work with

Haiyan Yang and Michael Woods

Memorial University
St. John's, Newfoundland and Labrador, Canada
jcloredoosti@mun.ca

August 7, 2018



Linkage and association

- Linkage and association methods are widely used in genetic analysis.
- Linkage describes a relationship between phenotype and loci while association describes a relationship between phenotype and alleles.
- Linkage is a consequence of co-segregation, a fundamental genetic process.
- Association is plainly a statistical statement about co-occurrence of alleles and phenotype, embellished with some population genetics arguments that most of times are of no use in the analysis.

- Linkage is a phenomenon to be studied within families.
- However, when two or more supposedly unrelated individuals from a particular population share the same phenotype, it is natural to think that the population genetic processes have affected all the members of the population and not only unrelated singletons.
- On the other hand, association can be studied not only in families but in unrelated individuals as well.

- This indicates that both linkage and association can be modelled together and, in fact, there are some methods that use both linkage and association (TDT, FBAT).
- Of course, joint modelling of linkage and association does not consist of performing a linkage study followed by an association study with case-controls (or any other design) and a denser set of markers on the regions where 'a linkage peak' was found.
- A 'good' methodology takes into consideration both linkage and association in the development of estimation and testing procedures and provides some flexibility to deal with general family configurations and 'unrelated' individuals.

Logistic model with over-dispersion

- Suppose that we have a random sample y_1, y_2, \dots, y_T , with $y_t \in \{0, 1\}$, $t = 1, 2, \dots, T$ and each y_t has an associated vector of covariates $\{\mathbf{z}_t \in \mathbb{R}^q, t = 1, 2, \dots, T\}$. Additionally, assume that for each t there is an independent random variable $\xi_t \in \mathbb{R}$, $t = 1, 2, \dots, T$ such that given ξ_t and \mathbf{z}_t , the random variable y_t has a Bernoulli distribution with parameter $\pi_t(\xi_t)$ where

$$\begin{aligned}\pi_t(\xi_t) &= \Pr(y_t = 1 \mid \mathbf{z}_t, \xi_t) \\ &= \frac{e^{\boldsymbol{\theta}'\mathbf{z}_t + \xi_t}}{1 + e^{\boldsymbol{\theta}'\mathbf{z}_t + \xi_t}} \quad \text{for } t = 1, 2, \dots, T.\end{aligned}$$

- Define $\zeta_t = e^{\theta' \mathbf{z}_t}$, $\zeta_t(\xi) = e^{\xi} \zeta_t$ and $\lambda_t(\xi) = (1 + \zeta_t(\xi))^{-1}$ so that $\pi_t(\xi) = \zeta_t(\xi) \lambda_t(\xi)$. Thus, if every ξ_t has a density parameterized by ϑ , say $f_{\vartheta}(\cdot)$, for $t = 1, 2, \dots, T$

$$\Pr(y_t | \mathbf{z}_t, \xi) = \zeta_t^{y_t}(\xi) \lambda_t(\xi) \quad (1)$$

$$\Pr(y_t | \mathbf{z}_t) = \int_{-\infty}^{\infty} \Pr(y_t | \mathbf{z}_t, \xi) f_{\vartheta}(\xi) d\xi \quad (2)$$

$$= y_t + (-1)^{y_t} \int_{-\infty}^{\infty} \lambda_t(\xi) f_{\vartheta}(\xi) d\xi$$

$$= \pi_t,$$

i.e., π_t is a function of y_t and \mathbf{z}_t only.

Consequently, the conditional on $\{\mathbf{z}_t\}$,

$$\Pr(\mathbf{y} | \{\mathbf{z}_t\}) = \prod_{t=1}^T \pi_t \quad (3)$$

and the log-likelihood can be written as

$$\ell(\boldsymbol{\theta}, \boldsymbol{\vartheta}) = \sum_{t=1}^T \log \pi_t$$

with $\boldsymbol{\theta}$ and $\boldsymbol{\vartheta}$ being the parameters of the logistic model and over-dispersion distribution, respectively.

- Under mild assumptions regarding $f_{\vartheta}()$,

$$\begin{aligned}\frac{\partial \ell(\boldsymbol{\theta}, \boldsymbol{\vartheta})}{\partial \boldsymbol{\theta}} &= \sum_{t=1}^T \mathbf{z}_t (y_t - E_{\xi_t}(\pi_t(\xi_t) \mid \mathbf{z}_t, y_t)) \\ &= \sum_{t=1}^T \mathbf{z}_t \left(\frac{y_t - \mu_t}{\mu_t(1 - \mu_t)} \right) E_{\xi_t}(\text{Var}(y_t \mid \mathbf{z}_t, \xi_t))\end{aligned}\quad (4)$$

where $\mu_t = E(y_t \mid \mathbf{z}_t)$ and $\text{Var}(y_t \mid \mathbf{z}_t, \xi_t) = \pi_t(\xi_t)(1 - \pi_t(\xi_t))$,
and

$$\frac{\partial \ell(\boldsymbol{\theta}, \boldsymbol{\vartheta})}{\partial \boldsymbol{\vartheta}} = \sum_{t=1}^T E_{\xi_t} \left(\frac{\partial f_{\vartheta}(\xi_t)}{\partial \boldsymbol{\vartheta}} \mid \mathbf{z}_t, y_t \right) \quad (5)$$

- To evaluate the expectations involved in the previous expressions we can use the following relationship. For $r = -1, 0, 1, 2, \dots$, define

$$P_t^{(r)} = \int_{-\infty}^{\infty} \pi_t^{r+1}(\xi) f_{\theta}(\xi) dz,$$

so that for $r = 0, 1, 2, \dots$,

$$\pi_t E_{\xi_t} (\pi_t^r(\xi_t) | \mathbf{z}_t, y_t) = (1 - y_t) P_t^{(r-1)} + (-1)^{1-y_t} P_t^{(r)}$$

and

$$\mu_t = P_t^{(0)}$$

as well as

$$E_{\xi_t} (\text{Var}(y_t | \mathbf{z}_t, \xi_t)) = P_t^{(0)} - P_t^{(1)}$$

Normal random effects

- If we assume that $f_{\theta}(\cdot)$ is the normal distribution with null mean and variance $\sigma^2 > 0$ variance. Then

$$\begin{aligned}\frac{\partial \ell(\boldsymbol{\theta}, \sigma^2)}{\partial \sigma^2} &= -\frac{T}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{t=1}^T \mathbb{E}_{\xi_t} (\xi_t^2 | \mathbf{z}_t, y_t) \\ &= -\frac{1}{2\sigma^2} \sum_{t=1}^T \frac{y_t - \mu_t}{1 - \mu_t} + \frac{1}{2\sigma^4} \sum_{t=1}^T \frac{y_t - \mu_t}{\mu_t(1 - \mu_t)} \mathbb{E}_{\xi_t} (\xi_t^2 \pi_t(\xi_t))\end{aligned}$$

and the remaining components of the Fisher information matrix can be obtained in a similar way.

A model for association

- To accommodate for more than one family, lets replace the sub-index ' t ' in y_t , \mathbf{z}_t , ξ_t and π_t by ' ij ', to represent the j th individual from the i th family, $i = 1, 2, \dots, N$ and $j = 1, 2, \dots, n_i$.
- Furthermore, assume that $\mathbf{z}_{ij} = (1, g_{ij}, \mathbf{x}'_{ij})' \in \mathbb{R}^{S+2}$ with g_{ij} being the marker phenotype and \mathbf{x}_{ij} a set of covariates and, as before $y_{ij} \in \{0, 1\}$.

- Thus

$$\ell(\boldsymbol{\theta}, \boldsymbol{\vartheta}) = \sum_{i=1}^N \sum_j^{n_i} \log \pi_{ij}$$

etc...

- When the ξ_{ij} represent non-genetic over-dispersion, this is a simple model for association.

Mixed model with genetic random effects

- In this case, we assume that the vector of random effects for the i th family, ξ_i has null mean and variance $\sigma^2 \Sigma_i$, then

$$\Pr(\mathbf{y}_i | \{\mathbf{z}_{ij}\}, \xi_i) = \prod_{j=1}^{n_i} \eta_{ij}^{y_{ij}}(\xi_{ij}) \lambda(\xi_{ij})$$

and

$$\Pr(\mathbf{y}_i | \{\mathbf{x}_{ij}\}) = \int_{\mathbb{R}^{n_i}} \Pr(\mathbf{y}_i | \{\mathbf{z}_{ij}\}, \xi_i) f(\xi_i) d\xi_i$$

- Thus, under some regularity conditions,

$$\frac{\partial \Pr(\mathbf{y}_i | \{\mathbf{z}_{ij}\}, \xi_i)}{\partial \boldsymbol{\theta}} = \Pr(\mathbf{y} | \{\mathbf{z}_{ij}\}, \xi_i) \sum_{j=1}^{n_i} \mathbf{z}_{ij} (y_{ij} - \pi_{ij}(\xi_{ij}))$$

and, consequently,

$$\frac{\partial \ell_i(\boldsymbol{\theta}, \sigma^2)}{\partial \boldsymbol{\theta}} = \mathbb{E}_{\xi_i} \left(\sum_{j=1}^{n_i} \mathbf{z}_{ij} (y_{ij} - \pi_{ij}(\xi_{ij})) \mid \{\mathbf{z}_{ij}\}, \mathbf{y}_i \right)$$

- When $\xi_i \sim \mathcal{N}(\mathbf{0}, \sigma^2 \Sigma_i)$ with $\sigma^2 > 0$, then

$$\frac{\partial \ell_i(\boldsymbol{\theta}, \sigma^2)}{\partial \sigma^2} = -\frac{n_i}{2\sigma^2} + \frac{1}{2\sigma^4} \mathbb{E}_{\xi_i} (\xi_i' \Sigma^{-1} \xi_i \mid \{\mathbf{z}_{ij}\}, \mathbf{y}_i)$$

- Furthermore, if Σ_i can be written as $\Sigma_i = \mathbf{I} + \beta \Phi_i$, for some $\beta \geq 0$ and a n.n.d. Φ_i , then

$$\frac{\partial \ell_i(\boldsymbol{\theta}, \sigma^2, \beta)}{\partial \beta} = \frac{1}{2\sigma^2} \mathbb{E}_{\xi_i} (\xi_i' \Sigma^{-1} \Phi_i \Sigma^{-1} \xi_i - \sigma^2 \text{tr}(\Sigma^{-1} \Phi_i) \mid \{\mathbf{z}_{ij}\}, \mathbf{y}_i)$$

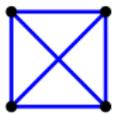
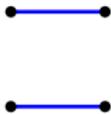
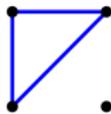
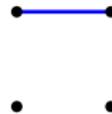
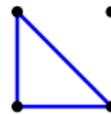
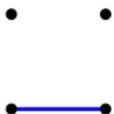
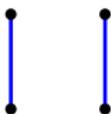
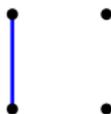
- We can think of Φ_i as the Malécot kinship matrix at any given locus for the i th family, say $\Phi_i^{(m)}$, and β as the signal-to-noise ratio at such a locus.
- Computation of $\varphi_{ij}^{(m)}$ involves the estimation of IBD at the m th locus.

IBD state		State description		Condensed states
$g_i = (a, b)$	$g_j = (c, d)$	Partition	Ewens'	\mathcal{S}
● ■	● ■	(a, b, c, d)	$(0,0,0,1)$	$\mathcal{J}_1 = (1, 1, 1)$
● ■	● ■	$(a, b)(c, d)$	$(0,2,0,0)$	
● ■	● ■	$(a, b, c)(d)$	$(1,0,1,0)$	} $\mathcal{J}_3 = (1, 0, 1)$
● ■	● ■	$(a, b, d)(c)$	$(1,0,1,0)$	
● ■	● ■	$(a, b)(c)(d)$	$(2,1,0,0)$	
● ■	● ■	$(a, c, d)(b)$	$(1,0,1,0)$	} $\mathcal{J}_5 = (0, 1, 1)$
● ■	● ■	$(a)(b, c, d)$	$(1,0,1,0)$	
● ■	● ■	$(a)(b)(c, d)$	$(2,1,0,0)$	$\mathcal{J}_6 = (0, 1, 0)$
● ■	● ■	$(a, c)(b, d)$	$(0,2,0,0)$	} $\mathcal{J}_7 = (0, 0, 2)$
● ■	● ■	$(a, d)(b, c)$	$(0,2,0,0)$	
● ■	● ■	$(a, c)(b)(d)$	$(2,1,0,0)$	
● ■	● ■	$(a, d)(b)(c)$	$(2,1,0,0)$	} $\mathcal{J}_8 = (0, 0, 1)$
● ■	● ■	$(a)(b, c)(d)$	$(2,1,0,0)$	
● ■	● ■	$(a)(b, d)(c)$	$(2,1,0,0)$	
● ■	● ■	$(a)(b)(c)(d)$	$(4,0,0,0)$	$\mathcal{J}_9 = (0, 0, 0)$

Table. IBD states for a pair of individuals at a given locus.

The kinship coefficient φ_{ij} between the i th and j th individuals can be written as

$$\varphi_{ij} = \Delta_{1,ij} + \frac{1}{2} (\Delta_{3,ij} + \Delta_{5,ij} + \Delta_{7,ij}) + \frac{1}{4} \Delta_{8,ij}$$

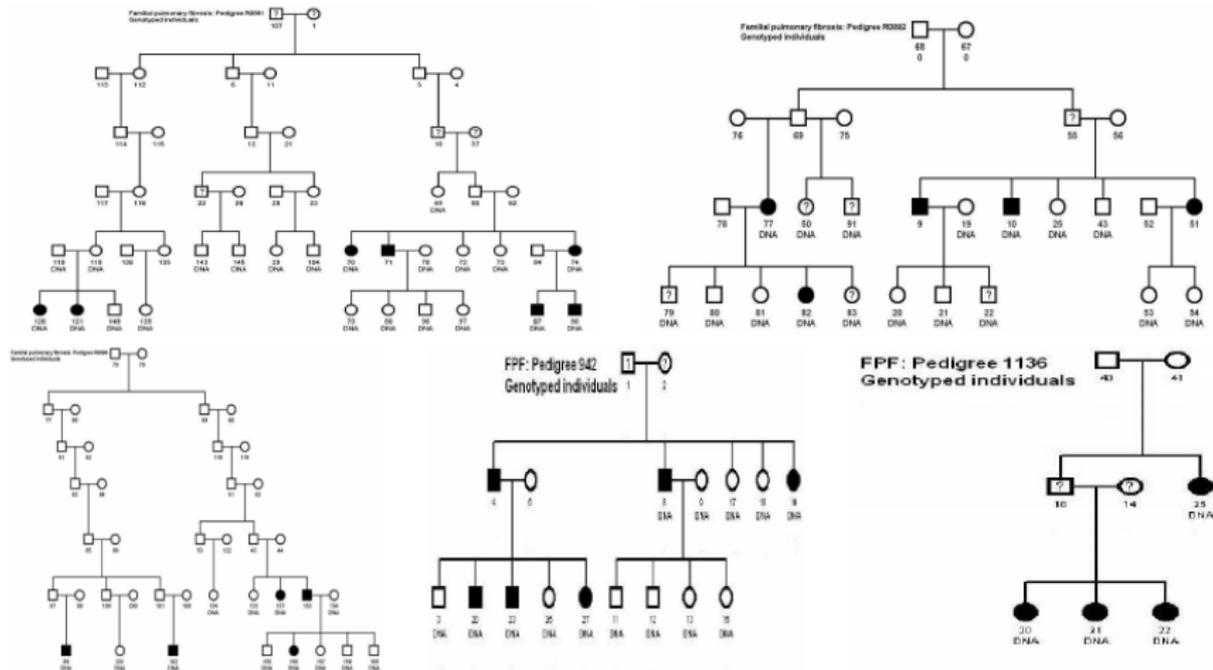
 \mathcal{S}_1  \mathcal{S}_2  \mathcal{S}_3  \mathcal{S}_4  \mathcal{S}_5  \mathcal{S}_6  \mathcal{S}_7  \mathcal{S}_8  \mathcal{S}_9

$$\Delta_r = \Pr(\mathcal{S} = \mathcal{S}_r)$$

IBD at a particular locus

- The random process underlying changes in the patterns of *IBD* across the genome is recombination, so the natural context for estimation of *IBD* in founders and singletons is the ancestral recombination graph, which specifies the complete ancestry of a collection of chromosomes.
- For families, the natural frameworks are *inheritance vectors*, or *descent graphs* and/or some variations of these methodologies some times written in FFT terms.

- Abecasis & Wigginton proposed a generalization of the Green & Lander algorithm to account for LD. However, this method does not scale well even for pedigrees of modest size.
- Han & Abney have developed a method for estimating the IBD at any position, given a dense genotype data and a pedigree of any size and complexity. Their method is based in a linear approximation of a HMM simplified version of the continuous-chromosome IBD process.
- Browning provides a population-based framework for the inference of *IBD*.
- Thompson and Glazner & Thompson address the issue for families and individuals within a population also through a HMM approximation of the *IBD* process.



A crude approximation

- In practice, much simpler approximations may provide satisfactory results. For example, we have estimated the set $\{\Delta_{r,ij}^{(m)}\}$ for the ij th pair of individuals at the m th marker with a simple application of the Bayes Theorem

$$\begin{aligned}\Delta_{r,ij}^{(m)} &= \Pr\left(\mathcal{S}_{ij}^{(m)} = \mathcal{J}_r \mid \mathbf{g}_i^{(m)}, \mathbf{g}_j^{(m)}\right) \\ &= \frac{\Delta_{r,ij}^o \Pr\left(\mathbf{g}_i^{(m)}, \mathbf{g}_j^{(m)} \mid \mathcal{S}_{ij}^{(m)} = \mathcal{J}_r\right)}{\sum_{s=1}^9 \Delta_{s,ij}^o \Pr\left(\mathbf{g}_i^{(m)}, \mathbf{g}_j^{(m)} \mid \mathcal{S}_{ij}^{(m)} = \mathcal{J}_s\right)}\end{aligned}$$

where the conditional *AIS* probabilities come from the next *Tables* and $\Delta_{r,ij}^o$ can be either the pedigree based estimate of the IBD coefficients or an estimate obtained for the ij th pair from properly chosen unlinked markers.

With a bi-allelic marker there are six distinguishable marker genotype pairs and their conditional probabilities given the condensed identity states, $\Pr(g_i, g_j | \mathcal{S} = \mathcal{I}_r)$, are

g_i, g_j	Condensed identity states								
	\mathcal{I}_1	\mathcal{I}_2	\mathcal{I}_3	\mathcal{I}_4	\mathcal{I}_5	\mathcal{I}_6	\mathcal{I}_7	\mathcal{I}_8	\mathcal{I}_9
aa, aa	p_a	p_a^2	p_a^2	p_a^3	p_a^2	p_a^3	p_a^2	p_a^3	p_a^4
aa, ab	0	0	$p_a p_b$	$2p_a^2 p_b$	0	0	0	$p_a^2 p_b$	$2p_a^3 p_b$
ab, aa	0	0	0	0	$p_a p_b$	$2p_a^2 p_b$	0	$p_a^2 p_b$	$2p_a^3 p_b$
aa, bb	0	$p_a p_b$	0	$p_a p_b^2$	0	$p_a^2 p_b$	0	0	$p_a^2 p_b^2$
ab, ab	0	0	0	0	0	0	$2p_a p_b$	$p_a p_b (p_a + p_b)$	$4p_a^2 p_b^2$
bb, aa	0	$p_a p_b$	0	$p_a^2 p_b$	0	$p_a p_b^2$	0	0	$p_a^2 p_b^2$
ab, bb	0	0	0	0	$p_a p_b$	$2p_a p_b^2$	0	$p_a p_b^2$	$2p_a p_b^3$
bb, ab	0	0	$p_a p_b$	$2p_a p_b^2$	0	0	0	$p_a p_b^2$	$2p_a p_b^3$
bb, bb	p_b	p_b^2	p_b^2	p_b^3	p_b^2	p_b^3	p_b^2	p_b^3	p_b^4

With a polymorphic marker, the relevant conditional probabilities $\Pr(g_i, g_j \mid \mathcal{S} = \mathcal{J}_r)$ are

g_i, g_j	Condensed identity states								
	\mathcal{J}_1	\mathcal{J}_2	\mathcal{J}_3	\mathcal{J}_4	\mathcal{J}_5	\mathcal{J}_6	\mathcal{J}_7	\mathcal{J}_8	\mathcal{J}_9
aa, aa	p_a	p_a^2	p_a^2	p_a^3	p_a^2	p_a^3	p_a^2	p_a^3	p_a^4
aa, ab	0	0	$p_a p_b$	$2p_a^2 p_b$	0	0	0	$p_a^2 p_b$	$2p_a^3 p_b$
ab, aa	0	0	0	0	$p_a p_b$	$2p_a^2 p_b$	0	$p_a^2 p_b$	$2p_a^3 p_b$
aa, bb	0	$p_a p_b$	0	$p_a p_b^2$	0	$p_a^2 p_b$	0	0	$p_a^2 p_b^2$
ab, ab	0	0	0	0	0	0	$2p_a p_b$	$p_a p_b (p_a + p_b)$	$4p_a^2 p_b^2$
aa, bc	0	0	0	$2p_a p_b p_c$	0	0	0	0	$2p_a^2 p_b p_c$
ab, ac	0	0	0	0	0	0	0	$p_a p_b p_c$	$4p_a^2 p_b p_c$
ab, cc	0	0	0	0	0	$2p_a p_b p_c$	0	0	$2p_a p_b p_c^2$
ab, cd	0	0	0	0	0	0	0	0	$4p_a p_b p_c p_d$

- Of course, $\Delta_r^{(m)}$ s obtained this way do not possess the beautiful properties of the estimates obtained through any of the HHM approaches mentioned above. In fact, they may not even be compatible with the legal states of Markov chain. However, for all practical purposes, they produce very much the same numerical results for this problem.
- Be aware that estimating $\Delta_{s,ij}^o$ by method of moments from SNP data is not possible because in such a case the matrix of conditional probabilities given by the bi-allelic *Table* is not of full rank. This also implies that the log-likelihood is not strictly convex so one needs to take extra precaution to find the MLE estimates. Nonetheless, kinship and inbreeding are identifiable.

Bottom line is once $\Delta_{r,ij}^{(m)}$ for a fixed locus m has been obtained, we can compute

$$\varphi_{ij}^{(m)} = \Delta_{1,ij}^{(m)} + \frac{1}{2} \left(\Delta_{3,ij}^{(m)} + \Delta_{5,ij}^{(m)} + \Delta_{7,ij}^{(m)} \right) + \frac{1}{4} \Delta_{8,ij}^{(m)}$$

where

$$\Delta_{r,ij}^{(m)} = \Pr \left(\mathcal{S}^{(m)} = \mathcal{J}_r \mid \text{'marker data'} \right)$$

so that $\Sigma = I + \beta \Phi^{(m)}$ in our model, β captures the linkage signal from the m th locus, while α does association.

- The joint linkage and association test is

$$H_0 : \alpha = 0 \text{ and } \beta = 0$$

vs $H_A : \alpha \neq 0 \text{ and/or } \beta > 0$

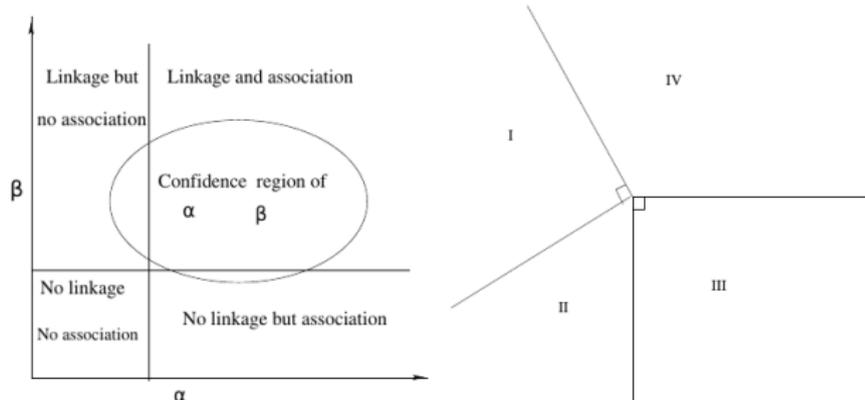
- The alternative has three cases:

$$H_{A_1} : \alpha \neq 0 \text{ and } \beta = 0$$

$$H_{A_2} : \alpha = 0 \text{ and } \beta > 0$$

$$H_{A_3} : \alpha \neq 0 \text{ and } \beta > 0$$

Testing for linkage and association



- Region I represents the LRT for H_{A_1} with angle $\pi/2$ and dist. χ_1^2 ; region II represents H_{A_2} with angle $\pi/2$ and dist. $\frac{1}{2}\chi_0^2 + \frac{1}{2}\chi_1^2$; region IV has an angle ρ and dist. $\frac{1}{2}\chi_1^2 + \frac{1}{2}\chi_2^2$; and Region III with angle $\pi - \rho$ and dist. χ_0^2 ,

$$\left(\rho = \arccos \frac{i_{\alpha\beta}}{\sqrt{i_{\alpha\alpha}i_{\beta\beta}}} \text{ under } H_0 \right).$$
- So, the LRT has dist. $\frac{3}{8}\chi_0^2 + \frac{1}{2}\chi_1^2 + \frac{1}{8}\chi_2^2$.

Testing for linkage and association

α	β	$\mu = -3$			$\mu = -5$		
		Joint Association Linkage			Joint Association Linkage		
.2	.2	.394	.142	.129	.350	.117	.147
	.4	.463	.140	.192	.416	.109	.180
	.6	.562	.150	.235	.498	.106	.278
	.8	.604	.118	.298	.615	.095	.378
.4	.2	.632	.382	.115	.461	.216	.109
	.4	.709	.362	.146	.545	.213	.181
	.6	.758	.353	.221	.628	.251	.271
	.8	.790	.372	.268	.722	.232	.371
.6	.2	.884	.746	.112	.623	.385	.124
	.4	.883	.692	.142	.701	.398	.190
	.6	.930	.695	.201	.749	.392	.275
	.8	.938	.660	.260	.836	.434	.380
.8	.2	.978	.920	.106	.773	.586	.117
	.4	.969	.901	.127	.840	.603	.196
	.6	.980	.898	.164	.882	.636	.273
	.8	.984	.887	.213	.928	.619	.368

Table: Empirical power of score tests for binary mixed models

- What about 'joint families and case/controls'?

- We started with

$$\begin{aligned}\pi_{ij}(\xi_{ij}) &= \Pr(y_{ij} = 1 \mid \mathbf{z}_{ij}, \xi_{ij}) \\ &= \frac{e^{\boldsymbol{\theta}' \mathbf{z}_{ij} + \xi_{ij}}}{1 + e^{\boldsymbol{\theta}' \mathbf{z}_{ij} + \xi_{ij}}} \quad \text{for } i = 1, 2, \dots, N, j = 1, 2, \dots, n_i.\end{aligned}$$

and $\boldsymbol{\theta} = (\mu, \alpha, \boldsymbol{\gamma}')'$, $\boldsymbol{\xi}_i \sim \mathcal{N}(0, \sigma^2 \boldsymbol{\Sigma}_i)$ with $\boldsymbol{\Sigma}_i = \mathbf{I} + \beta \boldsymbol{\Phi}_i$,
etc. . .

- This model can be easily adapted to other situations.

Thank you



CIHR **IRSC**

Canadian Institutes of
Health Research

Instituts de recherche
en santé du Canada

