



---

**Statistical genetics:  
statistical concepts in a nutshell**

## Content

1. Motivation
2. Approaches to statistics
3. Hypothesis testing
4. **Model comparison**
5. SNP association studies



## Likelihood ratio test

- ▶ The null hypothesis  $\mathcal{H}_0$  states that model  $\mathcal{M}_0$  fits the observed data  $\mathbf{y}$  as well as  $\mathcal{M}_1$  does, while the alternative hypothesis  $\mathcal{H}_1$  states the opposite
- ▶ The test statistic is defined as

$$t(\mathbf{y}) = -2 \log \left\{ \frac{p(\mathbf{y} | \hat{\boldsymbol{\theta}}_0, \mathcal{M}_0)}{p(\mathbf{y} | \hat{\boldsymbol{\theta}}_1, \mathcal{M}_1)} \right\},$$

where  $\mathcal{M}_0$  and  $\mathcal{M}_1$  are respectively the simpler (null) and more complex (alternative) model and  $\hat{\boldsymbol{\theta}}_0$  and  $\hat{\boldsymbol{\theta}}_1$  the corresponding maximum likelihood estimates



## Likelihood ratio test

- ▶ The distribution of the test statistic

$$t(\mathbf{y}) = -2 \log \left\{ \frac{p(\mathbf{y} | \hat{\boldsymbol{\theta}}_0, \mathcal{M}_0)}{p(\mathbf{y} | \hat{\boldsymbol{\theta}}_1, \mathcal{M}_1)} \right\},$$

under the assumption that  $\mathcal{H}_0$  is true, converges against a  $\chi^2(\nu)$  distribution as  $n$  approaches  $\infty$  with  $\nu$  being the difference in the number of parameters between  $\mathcal{M}_1$  and  $\mathcal{M}_0$

- ▶ The observed value  $t(\mathbf{y})$  of the test statistic is calculated and the hypothesis  $\mathcal{H}_0$  that model  $\mathcal{M}_0$  fits the data as well as  $\mathcal{M}_1$  does is rejected if  $t(\mathbf{y}) > c$



## $\chi^2$ approximation to the likelihood ratio for a simple $\mathcal{H}_0$

- ▶  $\mathcal{H}_0 : \theta = \theta_0$  and  $\mathcal{H}_1 : \theta \neq \theta_0$
- ▶ Expand the log-likelihood  $\ell(\theta_0 | \mathbf{y}) = p(\mathbf{y} | \theta_0)$  as a second-order Taylor series around the maximum likelihood estimate  $\hat{\theta}$

$$\ell(\theta_0 | \mathbf{y}) \approx \ell(\hat{\theta} | \mathbf{y}) + \ell'(\hat{\theta} | \mathbf{y})(\theta_0 - \hat{\theta}) + \ell''(\hat{\theta} | \mathbf{y})(\theta_0 - \hat{\theta})^2/2$$

- ▶ Plug the expansion into  $t(\mathbf{y}) = -2\ell(\theta_0 | \mathbf{y}) + 2\ell(\hat{\theta} | \mathbf{y})$

$$\begin{aligned} t(\mathbf{y}) &\approx -2\ell(\hat{\theta} | \mathbf{y}) - \ell''(\hat{\theta} | \mathbf{y})(\theta_0 - \hat{\theta})^2 + 2\ell(\hat{\theta} | \mathbf{y}) \\ &= -\ell''(\hat{\theta} | \mathbf{y})(\theta_0 - \hat{\theta})^2 \end{aligned}$$

- ▶ By the LLN and since  $\hat{\theta}$  is a consistent estimator

$$-\frac{1}{n}\ell''(\hat{\theta} | \mathbf{y}) \xrightarrow{\mathcal{P}} -\mathbb{E}[\ell''(\theta_0 | \mathbf{y})] = \mathcal{I}(\theta_0)$$

 $\chi^2$  approximation to the likelihood ratio for a simple  $\mathcal{H}_0$  cont'd

- ▶ The maximum likelihood estimator  $\hat{\theta}$  is asymptotically normal:

$$\sqrt{n\mathcal{I}(\theta_0)}(\hat{\theta} - \theta_0) \xrightarrow{\mathcal{D}} \text{Normal}(0, 1)$$

- ▶ The distribution of the test statistic  $t(\mathbf{y})$  converges against a  $\chi^2(1)$  distribution as  $n$  approaches  $\infty$

$$n\mathcal{I}(\theta_0)(\hat{\theta} - \theta_0)^2 \xrightarrow{\mathcal{D}} \chi^2(1)$$

because the square of a standard normal random variable is  $\chi^2$  distributed with 1 degree of freedom



## Bayesian information criterion (BIC)

- ▶ Approximate the prior predictive distribution under  $\mathcal{M}_i$  with Laplace's method

$$p(\mathbf{y} | \mathcal{M}_i) \approx (2\pi)^{d/2} \det(Q)^{-1/2} p(\mathbf{y} | \hat{\theta}, \mathcal{M}_i) p(\hat{\theta} | \mathcal{M}_i),$$

where  $\hat{\theta}$  is the maximum likelihood estimate and  $Q$  the negative Hessian of the log-likelihood evaluated at  $\hat{\theta}$

- ▶ The BIC is derived by writing

$$\begin{aligned} -2 \log p(\mathbf{y} | \mathcal{M}_i) &\approx -d \log(2\pi) + \log \det(Q) - 2 \log p(\hat{\theta} | \mathcal{M}_i) \\ &\quad - 2 \log p(\mathbf{y} | \hat{\theta}, \mathcal{M}_i) \end{aligned}$$

- ▶ By the LLN and since  $\hat{\theta}$  is a consistent estimator

$$Q = -\frac{n}{n} \ell''(\hat{\theta} | \mathbf{y}) \xrightarrow{P} -n \mathbb{E}[\ell''(\theta_0 | \mathbf{y})] = n \mathcal{I}(\theta_0)$$



## Bayesian information criterion (BIC) cont'd

- ▶  $-2$  times the log prior predictive distribution is thus

$$\begin{aligned} -2 \log p(\mathbf{y} | \mathcal{M}_i) &\approx -d \log(2\pi) + d \log n + \log \det\{\mathcal{I}(\theta_0)\} \\ &\quad - 2 \log p(\hat{\theta} | \mathcal{M}_i) - 2 \log p(\mathbf{y} | \hat{\theta}, \mathcal{M}_i) \end{aligned}$$

- ▶ Dropping all terms that remain fixed as the sample size approaches  $\infty$  results in

$$-2 \log p(\mathbf{y} | \mathcal{M}_i) \approx \text{BIC}_i = -2 \log p(\mathbf{y} | \hat{\theta}, \mathcal{M}_i) + d \log n$$

- ▶ Small BIC values correspond to better models
- ▶ The Akaike Information Criterion is equal to

$$\text{AIC}_i = -2 \log p(\mathbf{y} | \hat{\theta}, \mathcal{M}_i) + 2d$$

- ▶ Small AIC values also correspond to better models





## Bayesian model averaging

- ▶ Considering that some models perform equally well, it seems reasonable to base inference on several models by using Bayesian model averaging
- ▶ Model uncertainty is then accounted for by including information from all models weighted by their posterior model probability
- ▶ The model-averaged posterior of some quantity of interest  $\Delta$  with the same interpretation across models is given by

$$p(\Delta | \mathbf{y}) = \sum_{k=1}^K p(\Delta | \mathcal{M}_k, \mathbf{y})p(\mathcal{M}_k | \mathbf{y})$$



## Bayesian model averaging cont'd

- ▶ The posterior probability of model  $\mathcal{M}_k$  is given by

$$p(\mathcal{M}_k | \mathbf{y}) = \frac{p(\mathbf{y} | \mathcal{M}_k)p(\mathcal{M}_k)}{\sum_{l=1}^K p(\mathbf{y} | \mathcal{M}_l)p(\mathcal{M}_l)},$$

where  $p(\mathbf{y} | \mathcal{M}_k)$  denotes the marginal likelihood of model  $\mathcal{M}_k$  and  $p(\mathcal{M}_k)$  the prior probability that model  $\mathcal{M}_k$  is true

- ▶ Assuming uniform prior model probabilities, the approximate posterior model probability of model  $\mathcal{M}_k$  is using the BIC given by

$$\hat{p}(\mathcal{M}_k | \mathcal{D}) = \frac{\exp \left\{ -\frac{1}{2} \text{BIC}_k \right\}}{\sum_{l=1}^K \exp \left\{ -\frac{1}{2} \text{BIC}_l \right\}}$$



## Bayesian model averaging for logistic regression

- ▶ The model-averaged posterior inclusion probability of a predictor  $x_i$  is given by

$$p(\beta_i \neq 0 \mid \mathcal{D}) = \sum_{k=1}^K \mathbb{1}_{\mathcal{M}_k}(\beta_i) p(\mathcal{M}_k \mid \mathcal{D})$$

The model-averaged posterior mean of  $\beta_i$  is given by

$$\mathbb{E}[\beta_i \mid \mathcal{D}] = \sum_{k=1}^K \mathbb{E}[\beta_i \mid \mathcal{M}_k, \mathcal{D}] p(\mathcal{M}_k \mid \mathcal{D})$$



## Bayesian model averaging for logistic regression cont'd

- ▶ The model-averaged posterior variance of  $\beta_i$  is given by

$$\mathbb{V}[\beta_i | \mathcal{D}] = \sum_{k=1}^K \left\{ \left( \mathbb{V}[\beta_i | \mathcal{M}_k, \mathcal{D}] + \mathbb{E}[\beta_i | \mathcal{M}_k, \mathcal{D}]^2 \right) \times p(\mathcal{M}_k | \mathcal{D}) \right\} - \mathbb{E}[\beta_i | \mathcal{D}]^2$$

- ▶ If the sample size of the observed data  $\mathbf{y}$  is large, then the posterior  $p(\beta_i | \mathbf{y}, \mathcal{M}_k)$  of  $\beta_i$  under model  $\mathcal{M}_k$  is asymptotically normal
- ▶ The mean is equal to the maximum likelihood estimator and variance equal to respective diagonal element of the inverse of the observed information matrix evaluated at the maximum likelihood estimator



## Bayesian model averaging for logistic regression cont'd

- ▶ Dobutamine stress echocardiography study at the UCLA School of Medicine from 1991 until it closed in 1996
- ▶ The aim of the study was to assess if measurements taken during the stress echocardiography may be used to predict cardiac death, heart attack or coronary heart disease

## Top 5 posterior model probabilities

Rank	Model	Posterior model probability <sup>†</sup>	Cumulative posterior model probability	Posterior model odds
01	$\mathcal{M}_1$ : posSE, dobEF, hxofHT, restwma	0.0948	0.0948	1.00
02	$\mathcal{M}_2$ : posSE, dobEF	0.0864	0.1812	1.10
03	$\mathcal{M}_3$ : posSE, dobEF, restwma	0.0818	0.2631	1.16
04	$\mathcal{M}_4$ : posSE, dobEF, hxofHT	0.0797	0.3427	1.19
05	$\mathcal{M}_5$ : posSE, dobEF, hxofHT, ecg	0.0719	0.4147	1.32



## Bayesian model averaging estimate

Predictor	Posterior inclusion probability	Posterior mean	Posterior standard deviation	95% equal tail interval	
				lower	upper
intercept	1.000	-0.34999	1.1031	-2.5120	1.8120
posSE	0.978	1.1126	0.2975	0.5295	1.6957
dobEF	0.882	-0.03617	0.0177	-0.0655	-0.0166
hxofHT	0.546	0.42712	0.4550	0.1586	1.4061
restwma	0.492	0.43209	0.5078	0.1647	1.5915
ecg	0.403	0.31211	0.4315	0.1419	1.4064
hxofMI	0.208	0.11074	0.2502	-0.0121	1.0750
hxofDM	0.147	0.06297	0.1811	-0.0753	0.9344
baseEF	0.132	-0.00277	0.0132	-0.0807	0.0388



## Summary

- ▶ Likelihood ratio tests require nested models and rely on asymptotic approximations
- ▶ The BIC and AIC are tools that help balance model complexity and fit, which is evaluated through the maximized likelihood. The BIC prefers simpler models with small amounts of data, but becomes willing to accept more complex ones with increasing amount of data.
- ▶ Bayesian model averaging is a powerful tool to account for model uncertainty. The BIC may be used to approximate posterior model probabilities

## Content

1. Motivation
2. Approaches to statistics
3. Hypothesis testing
4. Model comparison
5. **SNP association studies**





## Relationship between SNP genotype and phenotype

- ▶ Genetic information is stored in the DNA in form of 4 nucleotide bases
- ▶ The human reference genome is approximately 3 giga bases long and any 2 humans differ in their genetic code by a small fraction
- ▶ Single Nucleotide Polymorphism (SNP) is a form of genetic variation at a genetic site at which the nucleotide base between 2 humans differs





## Relationship between SNP genotype and phenotype cont'd

- ▶ The genotype is, among other factors, a strong influence on the phenotype
- ▶ An association between genotype and phenotype may be presumed for disease susceptibility, drug treatment or crop yields
- ▶ In case of drug treatments, some people react normally to the treatment, whereas others show none or life-threatening effects
- ▶ A particular set of SNPs may be characteristic for these phenotypes
- ▶ The goal of genome-wide association studies is then to reveal SNP patterns that permit disease susceptibility screens or personalized drug treatments



## Relationship between SNP genotype and phenotype cont'd

- ▶ DNA carried by the chromosomes is present in 2 copies
- ▶ Without considering DNA copy number variation, the genotype at a biallelic SNP is either
  - ▶  $AA$  - 2 copies of the common allele
  - ▶  $AB$  - 1 copy of each allele
  - ▶  $BB$  - 2 copies of the rare allele

where allele refers to the particular nucleotide base.

- ▶ The frequency distribution of a SNP genotype and phenotype can be visualized in a contingency table

	<b>AA</b>	<b>AB</b>	<b>BB</b>	
<b>Case</b>	$n_{AA}^{\text{Case}}$	$n_{AB}^{\text{Case}}$	$n_{BB}^{\text{Case}}$	$n^{\text{Case}}$
<b>Control</b>	$n_{AA}^{\text{Control}}$	$n_{AB}^{\text{Control}}$	$n_{BB}^{\text{Control}}$	$n^{\text{Control}}$
	$n_{AA}$	$n_{AB}$	$n_{BB}$	$n$



## General genotype count model: prospective model

- ▶ The counts may be modeled directly

	<b>AA</b>	<b>AB</b>	<b>BB</b>	
<b>Case</b>	$n_{AA}^{\text{Case}}$	$n_{AB}^{\text{Case}}$	$n_{BB}^{\text{Case}}$	$n^{\text{Case}}$
<b>Control</b>	$n_{AA}^{\text{Control}}$	$n_{AB}^{\text{Control}}$	$n_{BB}^{\text{Control}}$	$n^{\text{Control}}$
	$n_{AA}$	$n_{AB}$	$n_{BB}$	$n$

- ▶ In a prospective model, the phenotype is the random variable, whereas the genotype variable is supposed to be known
- ▶ Under the null hypothesis  $\mathcal{H}_0$ , there exists no association between both variables and thus

$$p(\mathbf{y} | \theta, \mathcal{H}_0) = \binom{n}{n^{\text{Case}}} \theta^{n^{\text{Case}}} (1 - \theta)^{n^{\text{Control}}}$$



## General genotype count model: prospective model cont'd

- ▶ Assume that the prior of  $\theta$ , which represents the probability of being a case, is a Beta distribution

$$p(\theta) = \frac{1}{\mathcal{B}(\alpha, \beta)} \theta^{\alpha-1} (1 - \theta)^{\beta-1},$$

where the Beta function is equal to

$$1/\mathcal{B}(\alpha, \beta) = \Gamma(\alpha)\Gamma(\beta)/\Gamma(\alpha + \beta)$$

- ▶ The prior predictive distribution is then

$$\begin{aligned} p(\mathbf{y} | \mathcal{H}_0) &= \binom{n}{n^{\text{Case}}} \frac{1}{\mathcal{B}(\alpha, \beta)} \int_0^1 \theta^{n^{\text{Case}} + \alpha - 1} (1 - \theta)^{n^{\text{Control}} + \beta - 1} d\theta \\ &= \binom{n}{n^{\text{Case}}} \frac{\mathcal{B}(n^{\text{Case}} + \alpha, n^{\text{Control}} + \beta)}{\mathcal{B}(\alpha, \beta)} \end{aligned}$$



## General genotype count model: prospective model cont'd

- ▶ Under the alternative hypothesis  $\mathcal{H}_1$ , the 3 genotypes are assumed to be independent and thus

$$p(\mathbf{y} \mid \tau_{AA}, \tau_{AB}, \tau_{BB}, \mathcal{H}_1) = \binom{n}{n^{\text{Case}}} \prod_{i \in \{AA, AB, BB\}} \tau_i^{n_i^{\text{Case}}} \times (1 - \tau_i)^{n_i^{\text{Control}}}$$

- ▶ Assume that the prior of  $\tau_i$ , which represents the probability of being a case given that the genotype is  $i$ , has also a Beta distribution



## General genotype count model: prospective model cont'd

- ▶ The prior predictive distribution is then

$$\begin{aligned}
 p(\mathbf{y} | \mathcal{H}_1) &= \binom{n}{n^{\text{Case}}} \prod_{i \in \{AA, AB, BB\}} \frac{1}{\mathcal{B}(\alpha, \beta)} \times \\
 &\quad \int_0^1 \tau_i^{n_i^{\text{Case}} + \alpha - 1} (1 - \tau_i)^{n_i^{\text{Control}} + \beta - 1} d\tau_i \\
 &= \binom{n}{n^{\text{Case}}} \prod_{i \in \{AA, AB, BB\}} \frac{\mathcal{B}(n_i^{\text{Case}} + \alpha, n_i^{\text{Control}} + \beta)}{\mathcal{B}(\alpha, \beta)}
 \end{aligned}$$



## General genotype count model: prospective model cont'd

- ▶ The prospective Bayes factor is

$$\begin{aligned}
 B_{10} &= \frac{p(\mathbf{y} \mid \mathcal{H}_1)}{p(\mathbf{y} \mid \mathcal{H}_0)} \\
 &= \frac{\mathcal{B}(\alpha, \beta)}{\mathcal{B}(n^{\text{Case}} + \alpha, n^{\text{Control}} + \beta)} \times \\
 &\quad \prod_{i \in \{AA, AB, BB\}} \frac{\mathcal{B}(n_i^{\text{Case}} + \alpha, n_i^{\text{Control}} + \beta)}{\mathcal{B}(\alpha, \beta)}
 \end{aligned}$$

- ▶ Data-dependent hyperparameters may be used:

$$(\alpha, \beta) = \lambda (n^{\text{Case}}/n, n^{\text{Control}}/n) ,$$

which are uninformative in distinguishing  $\mathcal{H}_0$  from  $\mathcal{H}_1$  and where  $\lambda$  is used to scale the effect size.



Thank you for your attention