

Introduction ●○○	Imputation ○○○○○○○○○○○○○○○○○○	Hierarchical ASE ○○○○○○ ○○○○○○○○○○○○○○○○○○	Model Comparison ○○○○○○○○○○○○	Conclusions ○	References
---------------------	----------------------------------	--	----------------------------------	------------------	------------

2014 SISG MODULE 4: Bayesian Statistics for Genetics

Lecture 10: Imputation, Hierarchical Mixture Models, Model Comparison

Jon Wakefield

Departments of Statistics and Biostatistics
University of Washington

Introduction ○○●	Imputation ○○○○○○○○○○○○○○○○○○	Hierarchical ASE ○○○○○○ ○○○○○○○○○○○○○○○○○○	Model Comparison ○○○○○○○○○○○○	Conclusions ○	References
---------------------	----------------------------------	--	----------------------------------	------------------	------------

Outline

Introduction and Description of Data

Methods for Imputation

Hierarchical Modeling of Allele-Specific Expression Data

Motivation

Modeling

Model Comparison

Conclusions

Introduction

- In this lecture we consider two topics.
- **First**, we consider methods for **imputation** of missing genotypes.
- We describe a number of the more common Bayesian approaches to this problem.
- **Second**, we return to the Allele-Specific Expression (ASE) data introduced in **Lecture 3**.
- An in-depth analysis of the ASE data will be given – we will analyze with a **Bayes hierarchical mixture model**.
- The material in the second part of the lecture is necessarily more advanced than in previous lectures, but will give a sense of the power of Bayesian methods, and will highlight some of the practical issues in the use of these methods.
- Finally, we will briefly review a number of procedures to carry out **model comparison**.

Motivation for Imputation

- **Imputation** is the prediction of missing genotypes.
- **Imputation** is used in both GWAS and in fine-mapping studies.
- The technique is becoming increasingly popular since it can:
 - Increase **power** in GWAS.
 - Facilitate **meta-analysis** in which data it is required to combine information from different panels which have different sets of SNPs. In this way power can be increased.
 - Fine-map **causal variants**, see Figure 1. Imputed SNPs that show large associations can be better candidates for replication studies.
- The key idea in the approaches we describe is the use of data on haplotypes from a relevant population to build a **prior model** for the missing data, basically the models leverage **linkage disequilibrium**.

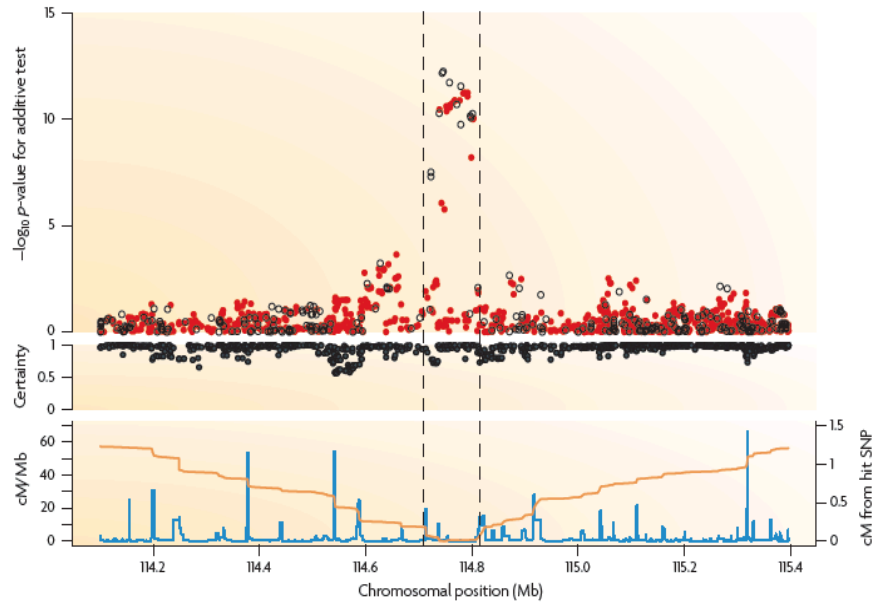


Figure 1 : Imputation for the TCF7L2 gene, from Marchini *et al.* (2007). Imputed SNP signals are in red and observed SNPs in black.

Box 1 | How genotype imputation works

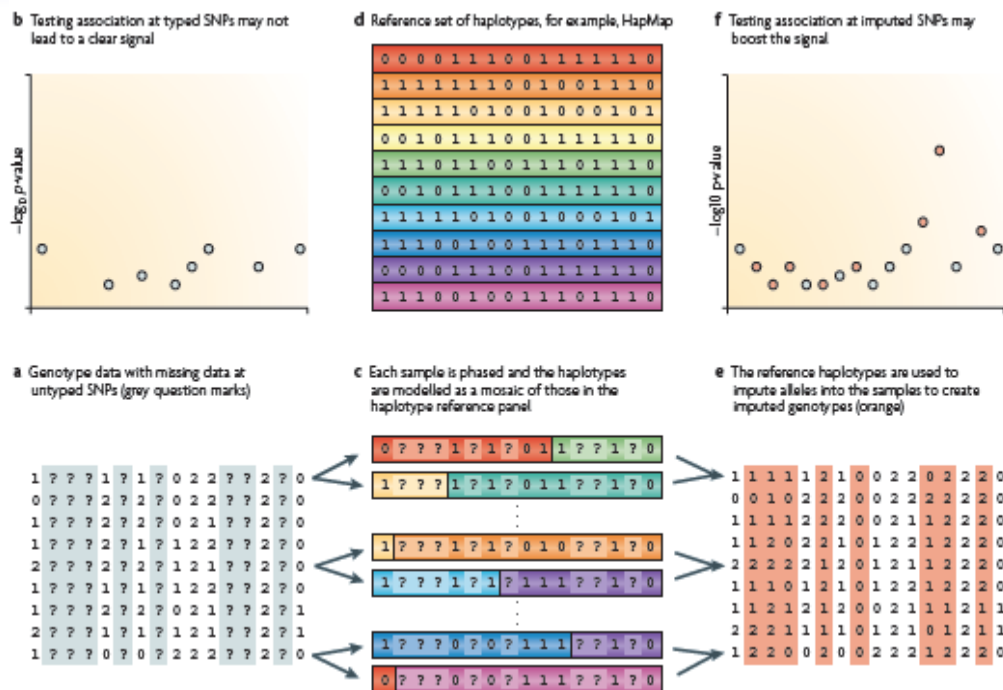


Figure 2 : Imputation overview from Marchini and Howie (2010).

The Statistical Framework

- Suppose we wish to estimate the association between a phenotype and m genetic markers in n individuals.
- Let G_{ij} represent the **genotype** of individual i at SNP j with G_{ij} unobserved for some SNPs.
- We consider diallelic SNPs so that G_{ij} can take the value 0, 1 or 2 depending on whether the pair of constituent SNPs are $\{0, 0\}$, $\{0, 1\}$, $\{1, 0\}$ or $\{1, 1\}$.
- If G_{ij} is **observed** then for SNP j we simply model

$$p(y_i | G_{ij})$$

- For example, if the phenotype y_i is **continuous**, we might assume a normal model:


$$E[Y_i] = \beta_0 + \beta_1 G_{ij},$$

and if y_i is **binary**, a logistic model is an obvious candidate:

$$\frac{p_i}{1 - p_i} = \exp(\beta_0 + \beta_1 G_{ij})$$

where p_i is the probability of disease for individual i .

The Statistical Framework

- Let $\mathbf{H} = (\mathbf{H}_1, \dots, \mathbf{H}_N)$ represent **haplotype** information at m SNPs in a relevant reference-panel, with N distinct haplotypes. 
- If G_{ij} is **unobserved** then for SNP j we have the model

$$p(y_i | \mathbf{H}, \mathbf{G}_i) = \sum_{k=0}^2 p(y_i | G_{ij} = k) \times \Pr(G_{ij} = k | \mathbf{H}, \mathbf{G}_i)$$

- The big question is how to obtain the predictive distribution

$$\Pr(G_{ij} = k | \mathbf{H}, \mathbf{G}_i).$$

- A common approach is to take as prior a **Hidden Markov Model (HMM)**.
- We digress to discuss HMMs.

Hidden Markov Models

- **Example: Poisson Time Series** A common problem is how to model count data over time. A Poisson model is the obvious choice but how to introduce:

1. overdispersion and
2. dependence over time.

- Consider the model:

Stage 1: $Y_t | \lambda_t \sim \text{Poisson}(\lambda_t)$, $t = 1, 2, \dots$

Stage 2: $\lambda_t | Z_t \sim_{iid} \begin{cases} \lambda_0 & \text{if } Z_t = 0 \\ \lambda_1 & \text{if } Z_t = 1 \end{cases}$

Stage 3: $Z_t | p \sim_{iid} \text{Bernoulli}(p)$.

- An alternative model replaces Stage 3 with a (first-order) **Markov chain model**, i.e., $\Pr(Z_t | Z_1, \dots, Z_{t-1}) = \Pr(Z_t | Z_{t-1})$:

$$\Pr(Z_t = 0 | Z_{t-1} = 0) = p_0$$

$$\Pr(Z_t = 1 | Z_{t-1} = 1) = p_1$$

- Z_t is an unobserved (hidden) state.
- As an example we consider the number of major earthquakes (magnitude 7 and above) for the years 1990–2006.
- We illustrate the fit of this model with **two** or **three** underlying states.

Example: Earthquake Data

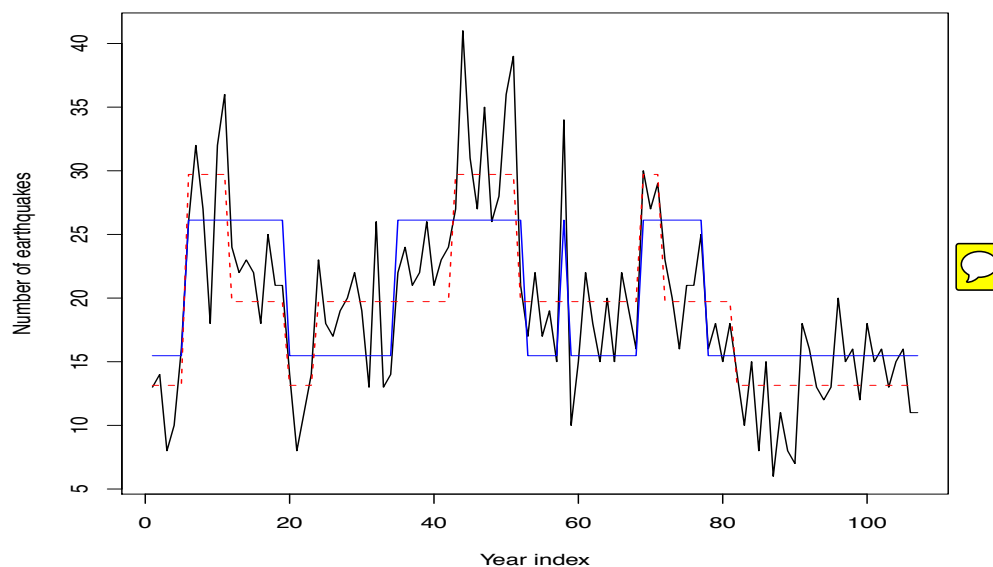


Figure 3 : The earthquake data along with the underlying states for the two and three state HMMs, in **blue** and **red**, respectively.

IMPUTE v1

- Marchini *et al.* (2007) consider a HMM for the vector of genotypes for individual i :

$$\Pr(\mathbf{G}_i | \mathbf{H}, \theta, \rho) = \sum_{\mathbf{Z}_i = (\mathbf{Z}_i^{(1)}, \mathbf{Z}_i^{(2)})} \Pr(\mathbf{G}_i | \mathbf{Z}_i, \theta) \times \Pr(\mathbf{Z}_i | \mathbf{H}, \rho)$$

where $\mathbf{Z}_i^{(1)} = \{Z_{i1}^{(1)}, \dots, Z_{iJ}^{(1)}\}$ and $\mathbf{Z}_i^{(2)} = \{Z_{i1}^{(2)}, \dots, Z_{iJ}^{(2)}\}$.

- The $(\mathbf{Z}_i^{(1)}, \mathbf{Z}_i^{(2)})$ are the pair of haplotypes for SNP j from the reference panel that are copied to form the genotype vector. These are the **hidden states**.
- The term $\Pr(\mathbf{Z}_i | \mathbf{H}, \rho)$ models how the pair of copied haplotypes for individual i changes along the sequence. This probability changes according to a **Markov chain** with the switching of states depending on the **fine-scale recombination rate** ρ .
- The term $\Pr(\mathbf{G}_i | \mathbf{Z}, \theta)$ allows the observed genotypes to differ from the pair of copied haplotypes through **mutation**; the mutation parameter is θ .
- **IMPUTE v2** (Howie *et al.*, 2009) is a more flexible version that alternates between phasing and haploid imputation.

fastPHASE and BIMBAM

- We describe the model of Scheet and Stephens (2006).
- A **Hidden Markov Model (HMM)** is used to determine $\Pr(G_{ij} = k | \alpha, \theta, r)$.
- The basic idea is that haplotypes tend to cluster into groups of similar haplotypes; suppose there are K clusters.
- The unobserved **hidden state** is the haplotype cluster from which this SNP arose from. Each cluster has an associated set of allele frequencies θ_{kj} .
- With K underlying states we have, for SNP j , α_{kj} being the probability of arising from haplotype k , with

$$\sum_{k=1}^K \alpha_{kj} = 1.$$

- The model is

$$\Pr(\mathbf{G}_i | \alpha, \theta, r) = \sum_{\mathbf{Z}} \Pr(\mathbf{G}_i | \mathbf{Z}_i, \theta) \times \Pr(\mathbf{Z}_i | \alpha, r)$$

with Z_{ij} the haplotype of origin for individual i and SNP j .

- A **Markov chain** is constructed for Z_{ij} with the strength of dependence being based on the recombination rate r at a given location.
- Given $Z_{ij} = k$, the genotype assigned depends on the allele frequencies of the k -th haplotype at the j -th SNP.

Use in Association Studies

- The simplest approach to using imputed SNPs is to substitute \hat{G}_{ij} (a number between 0 and 2) into the phenotype association model.
- A set of probabilities $\Pr(G_{ij} = k | \mathbf{G}, \mathbf{H})$ for $k = 0, 1, 2$ are produced and these may be used to average over the uncertainty in the phenotype model.
- Within **BIMBAM** the unknown genotype is sampled from its posterior distribution, within an MCMC framework.
- Other approaches:
 - **MACH**: similar methodology to IMPUTE (Li *et al.*, 2010).
 - **Beagle**: uses a graphical model for haplotypes Browning and Browning (2009).

Practical Issues

- One may attempt to match the haplotype panel (e.g. from HapMAP 2) with the study individuals.
- An alternative approach is to use all available haplotypes, and assigning equal prior probabilities to each.
- Many studies, for example Huang *et al.* (2009), have examined SNP imputation accuracy in different populations.

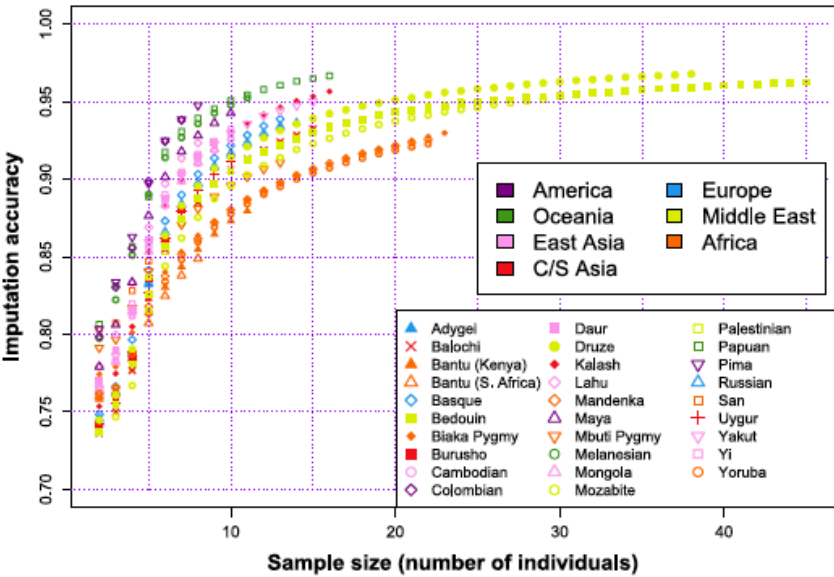


Figure 4 : Imputation accuracy as a function of sample size, from Huang *et al.* (2009).

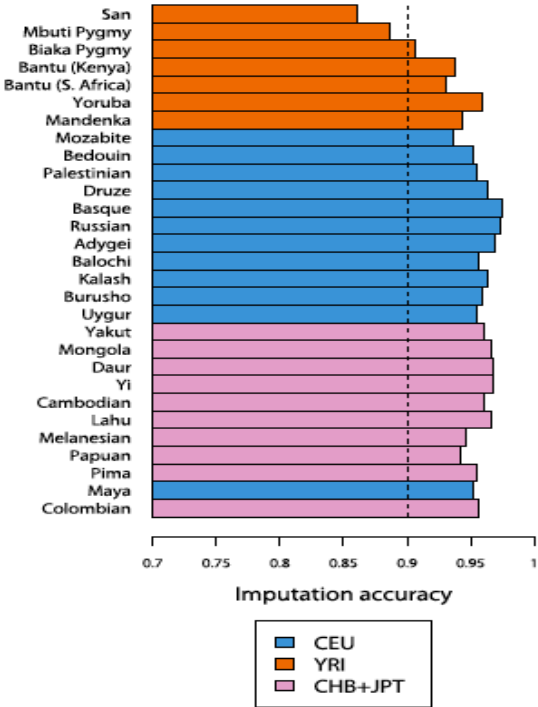


Figure 5 : Imputation accuracy for different populations with a reference-panel of 120 haplotypes. From Huang *et al.* (2009).

Introduction ○○○	Imputation ○○○○○○○○○○○○○○○●○	Hierarchical ASE ○○○○○○○ ○○○○○○○○○○○○○○○○○○○	Model Comparison ○○○○○○○○○○○	Conclusions ○	References
---------------------	---------------------------------	--	---------------------------------	------------------	------------

Table 1. Association Analysis results.

Locus	SNPname	Type	Effect Allele/ Other	Freq Effect Allele	Effect (SE) ^a	P-value	Genomic Annotation	Variance explained by the locus	Top GWAS SNP	Effect Allele/ Other	Freq Effect Allele	Effect (SE) ^a	P-value	r ²	Adjusted P-value	Variance explained by the locus
PCSK9	rs11591147	MetaboChip	T/G	0.037	−0.380 (0.048)	2.90×10 ^{−15}	missense (R46L)	1.19%	rs11206510	C/T	0.243	−0.106 (0.023)	5.71×10 ^{−07}	0.101	0.013	0.23%
	rs2479415	1000G	C/T	0.413	0.076 (0.019)	7.50×10 ^{−05}	8 Kb from PCSK9									
SORT1	rs583104	MetaboChip	T/G	0.177	0.149 (0.024)	1.28×10 ^{−09}	31 Kb from SORT1 ^b	0.63%	rs599839	G/A	0.276	−0.148 (0.025)	1.43×10 ^{−09}	0.991	0.90	0.61%
B3GALT4	rs28361085	1000G	C/T	0.073	0.114 (0.036)	0.00169	146 Kb from B3GALT3	0.22%	rs2254287	G/C	0.492	0.005 (0.018)	0.771	0.413	0.84	0.02%
B4GALT4	rs34507110	1000G	G/A	0.154	0.122 (0.030)	4.99×10 ^{−05}	83 Kb from B4GALT4	0.48%	rs12695382	A/G	0.075	−0.074 (0.035)	0.035	0.795	0.48	0.03%
APOB	rs547235	1000G	A/G	0.187	−0.144 (0.024)	1.69×10 ^{−09}	140 Kb from APOB	0.51%	rs562338	A/G	0.173	−0.139 (0.025)	1.43×10 ^{−8}	0.878	0.98	0.43%
LDLR	rs73015013	MetaboChip	T/C	0.138	−0.155 (0.027)	1.12×10 ^{−08}	9 kb from LDLR	1.17%	rs6511720	T/G	0.132	−0.160 (0.027)	1.71×10 ^{−08}	0.934	0.97	0.59%
	rs72658864	MetaboChip	C/T	0.005	0.626 (0.136)	3.90×10 ^{−06}	missense (V578A)									
APOC1/C2/E	rs7412	MetaboChip	T/C	0.037	−0.563 (0.048)	1.80×10 ^{−31}	missense (R176C) APOE	3.33%	rs4420638 ^c	G/A	0.097	0.218 (0.031)	4.67×10 ^{−12}	0.0003	6.41×10 ^{−10}	1.07%
	rs429358	Affy+Sanger	C/T	0.071	0.260 (0.036)	5.82×10 ^{−11}	missense (C130R) APOE									

The left panel shows the association results at 7 loci. For each gene, the strongest variant is listed first, and any second detected independent signal is listed with results from the conditional analysis (Materials and Methods). The column Type indicates whether the SNP was directly genotyped (MetaboChip) or imputed using 1000G reference haplotype (1000G) or the Sardinian reference panel (Affy+Sanger). The right panel shows the association results for the GWAS SNPs previously described [5], the correlation with the top SNP listed in the left panel, and its p-value in the conditional analysis (Adjusted P-value).
^aEffect sizes are standardized (see Materials and Methods), and represent the change in trait LDL-C values associated with each copy of the reference allele, measured in standard deviation units.
^bSNP rs583104 is also 1 Kb from PSRC1 transcript.
^cr² = 0.967 with MetaboChip second-independent SNP, rs429358. After adjusting for the two independent SNPs, rs7412 and rs429358, the p-value for rs4420638 was 0.5.
doi:10.1371/journal.pgen.1002198.t001

Figure 6 : Example from Sanna *et al.* (2011). Imputation carried out using the MACH software.

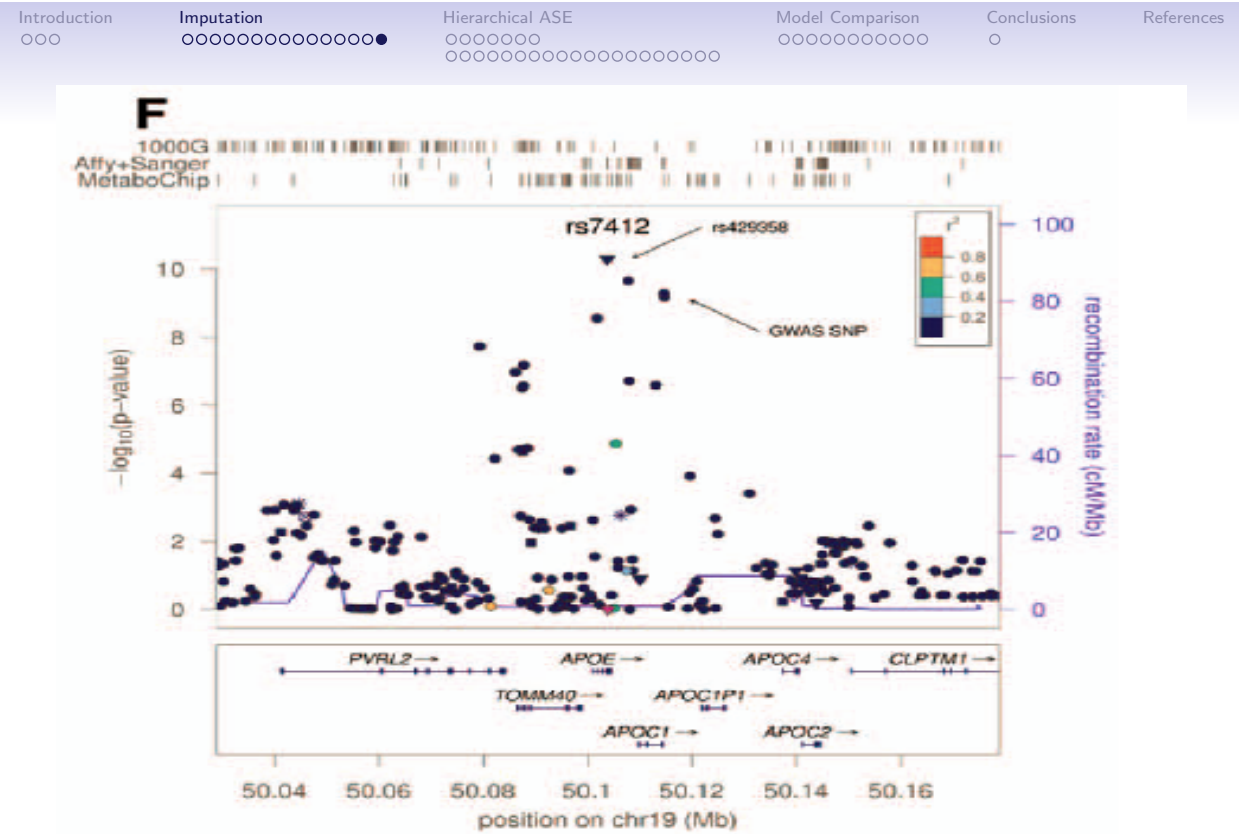


Figure 7 : Example from Sanna *et al.* (2011).

Specifics of ASE Experiment

Details of the data:

- Two “individuals” from genetically divergent yeast strains, BY and RM, are mated to produce a diploid hybrid.
- Three replicate experiments: same individuals, but separate samples of cells.
- Two technologies: Illumina and ABI SOLiD. Each of a few trillion cells are processed.
- Pre- and post-processing steps are followed by fragmentation to give millions of 200–400 base pair long molecules, with short reads obtained by sequencing.
- Strict criteria to call each read as a match are used, to reduce read-mapping bias.
- Data from 25,652 SNPs within 4,844 genes.

Allele Specific Expression via RNA-Seq

Additional data:

- **Genomic DNA** is sequenced in the diploid hybrid, which has one copy of each gene from BY and from RM.
- The only **difference** between the genomic DNA and the main experiment is that we expect the genomic DNA to always be present 50:50 (one copy each of BY and RM), whereas for the main experiment it is only 50:50 if there is no ASE.
- For both genomic DNA and RNA we obtain counts at SNPs, at each of BY and RM.

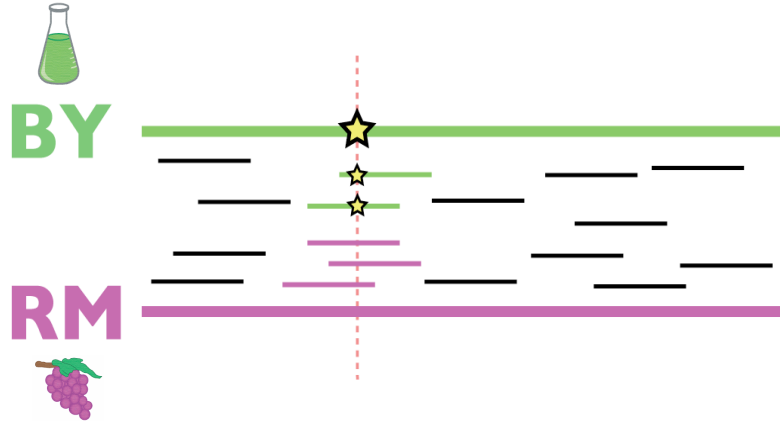


Figure 8 : Mapping of RNA short reads to BY and RM.

- **Aim of the Experiment:** Estimate the proportion of genes that display ASE.
- Let p be the probability of a map to BY at a particular SNP.
- Additionally, we would like to classify genes into:
 - Genes that do not show ASE.
 - Genes that show:
 - Constant ASE across SNPs.
 - Variable ASE across SNPs, i.e. p varies within gene.

Subsequently, we will examine genes displaying ASE to investigate the mechanism.

- A **hierarchical model** is feasible since we have **within gene** and **between gene** variability.
- Further, a **mixture model** is suggested, with a mixture of genes that do not display ASE (so there p 's are 0.5) and that do display ASE.

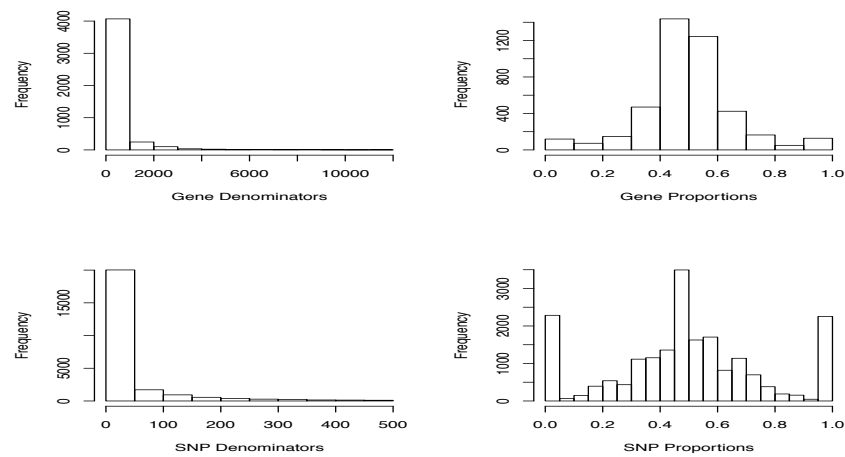


Figure 9 : Summaries for RNA BY/RM yeast data * 739 SNP denominators > 500.

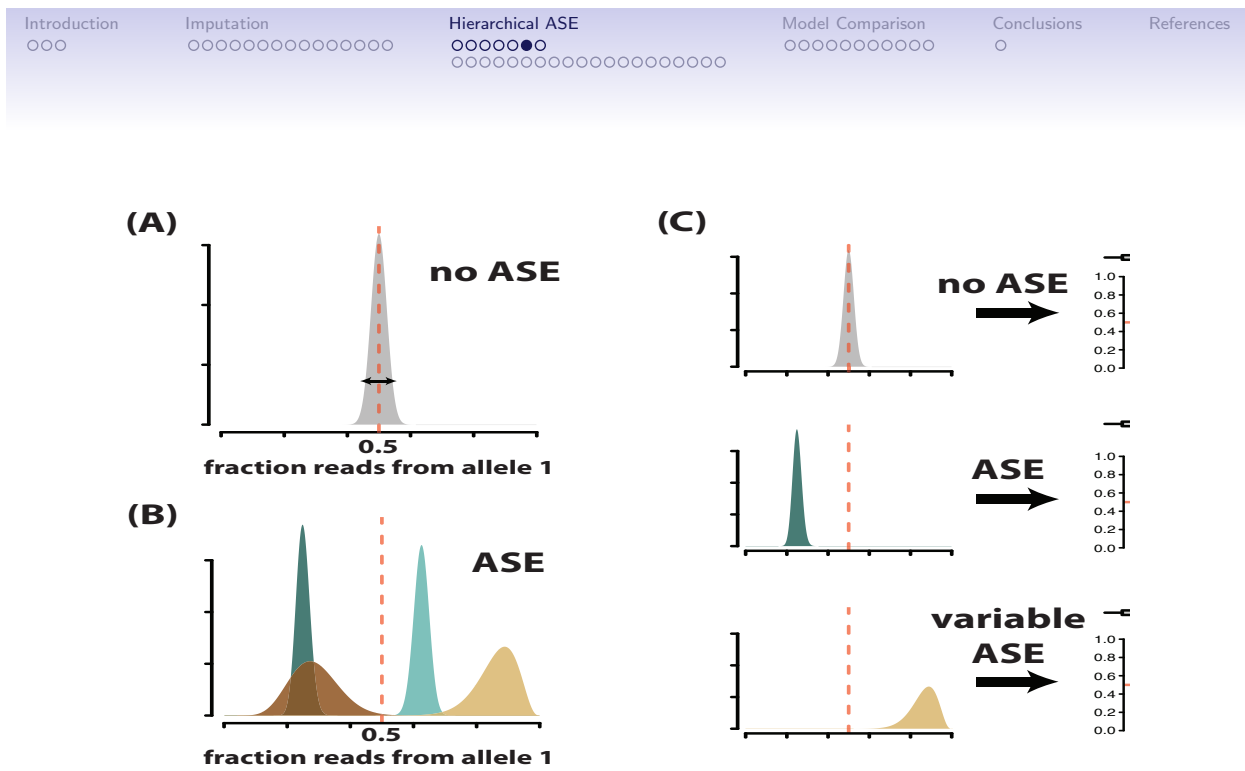


Figure 10 : Schematic of the hierarchical model.

Approach to Modelling RNASeq Data

Overview, three models fitted:

1. **Model 1:** Two component mixture model to filter out aberrant SNPs using **genomic DNA** data.
2. **Model 2:** Using the filtered genomic DNA data, fit a hierarchical SNP within gene model, to determine the “null” distribution of counts. Specifically: “wobble” in p about 0.5, and SNP “wobble” in p within genes.
Absence of ASE is not experimentally equivalent to $Y_i \sim \text{binomial}(N_i, p = 0.5)$ because of the steps involved in the experiment.
3. **Model 3:** For the RNA Seq data develop a two-component mixture model where each gene either displays no ASE, or ASE, with null component determined from the analysis of the genomic DNA data (**Model 2**).

Model 1: Filtering Model for Genomic DNA

Two-component mixture model for SNPs:

1. **Majority** of SNP counts arise from a beta-binomial distribution with p “close to” 0.5.
2. **Minority** of SNP counts arise from a beta-binomial distribution with p “not close to” 0.5 due to sequencing bias at these SNPs.
 - Data: y_j and N_j are counts at SNP j for $j = 1, \dots, m$ SNPs.
 - Note: Ignores gene information – don’t want to impose too much structure at this point.
 - SNPs that are more likely to arise from component 2 are then removed from further analyses.

Filtering Model for Genomic DNA

- *Stage 1: SNP Count Likelihood:*

$$y_j | p_j \sim \text{binomial}(N_j, p_j), \quad j = 1, \dots, N.$$

- *Stage 2: Between-SNP Prior:*

$$p_j | a, b, c, \pi_0 = \begin{cases} \text{beta}(a, a) & \text{with probability } \pi_0 \\ \text{beta}(b, c) & \text{with probability } 1 - \pi_0 \end{cases}$$

- *Stage 3: Hyperpriors:* Constrain $b < 1$, $c < 1$ to give U-shaped beta distribution.

$$a \sim \text{lognormal}(4.3, 1.8)^*$$

$$b \sim \text{uniform}(0, 1)$$

$$c \sim \text{uniform}(0, 1)$$

$$\pi_0 \sim \text{uniform}(0, 1)$$

* 80% interval for p : [0.43, 0.57]. Separate a, b, c, π_0 for each technology.

Implementation for Genomic DNA

- Integrate p_j from model to give:

$$y_j | a, b, c, \pi_0 \sim \pi_0 \times \text{beta-binomial}(N_j, a, a) + (1 - \pi_0) \times \text{beta-binomial}(N_j, b, c).$$

- This is a mixture of two distributions:
 1. The first distribution is for the majority of signals close to 0.5. The size of a denotes how close is close.
 2. The second distribution is for the minority of aberrant SNPs.

Implementation for Genomic DNA

- Likelihood:**

$$\Pr(\mathbf{y}|a, b, c, \pi_0) = \prod_{j=1}^N \binom{N_j}{Y_j} \left\{ \pi_0 \frac{\Gamma(2a)}{\Gamma(a)^2} \frac{\Gamma(y_j + a)\Gamma(N_j - y_j + a)}{\Gamma(N_j + 2a)} + (1 - \pi_0) \frac{\Gamma(b + c)}{\Gamma(b)\Gamma(c)} \frac{\Gamma(y_j + b)\Gamma(N_j - y_j + c)}{\Gamma(N_j + b + c)} \right\}$$

- Posterior:**

$$p(a, b, c, \pi_0|\mathbf{y}) \propto \Pr(\mathbf{y}|a, b, c, \pi_0) \times p(a)p(b)p(c)p(\pi_0).$$

- Implementation:** Markov chain Monte Carlo.

- Recall: Sequencing bias lead to aberrant SNPs, and these errors are likely to be repeated in the main experiment.
- SNPs falling in the second mixture component were removed from further analyses.

Posterior Distributions

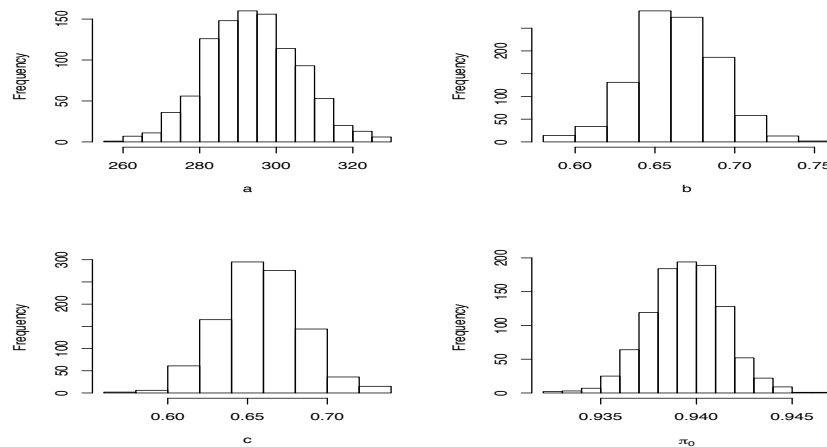


Figure 11 : Posteriors for genomic filtering model for Illumina platform.

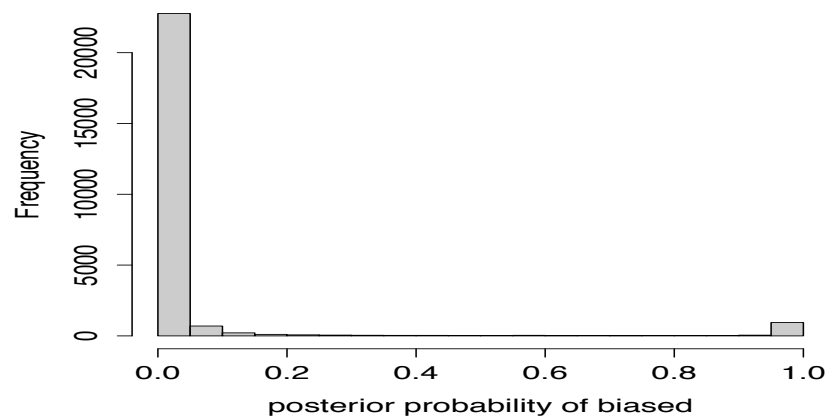


Figure 12 : Posterior probabilities of biased genomic DNA SNPs: 1,295 removed from 25,262.

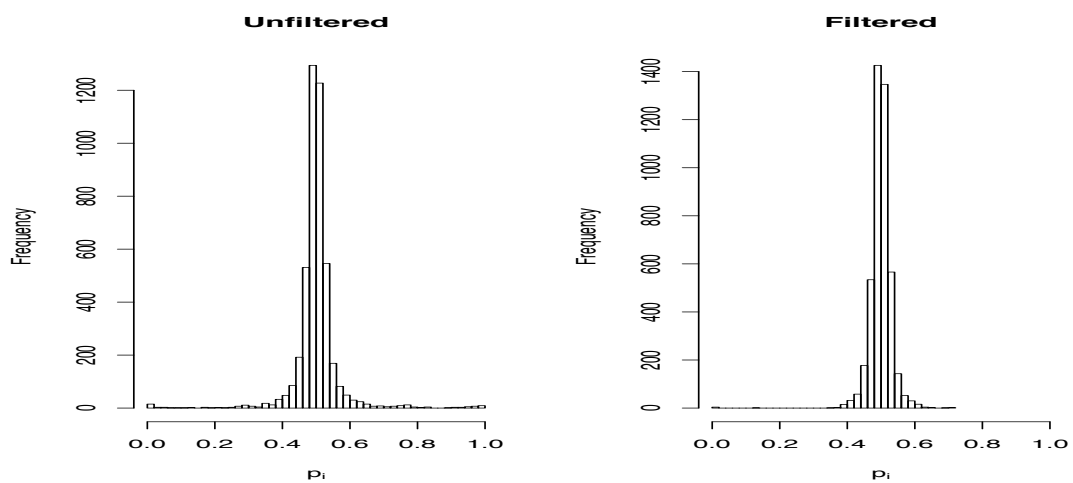


Figure 13 : Original and filtered data, for Illumina platform.

Model 2: Calibration Model for Genomic Data

- With aberrant SNPs removed, the next step is to calibrate the null component.

- *Stage 1: Within-Genes Likelihood:*

$$Y_{ij}|p_{ij} \sim \text{binomial}(N_{ij}, p_{ij}).$$

where p_{ij} is the probability of an outcome from the first genetic background.

- *Stage 2: Within-Genes Prior:*

$$p_{ij}|\alpha_i, \beta_i \sim \text{beta}(\alpha_i, \beta_i)$$

so that α_i, β_i determine the distribution of variants within gene i .

Calibration Model for Genomic Data

- α_i and β_i are not straightforward to interpret (as we saw in Lecture 3).
- We reparameterize $(\alpha_i, \beta_i) \rightarrow (p_i, e_i)$ with mean and dispersion parameters (recall $\alpha_i + \beta_i$ is a prior sample size):

$$p_i = \frac{\alpha_i}{\alpha_i + \beta_i}$$

$$e_i = \frac{1}{1 + \alpha_i + \beta_i}$$

- Moments of ASE parameters:

$$E[p_{ij}|p_i, e_i] = p_i$$

$$\text{var}(p_{ij}|p_i, e_i) = p_i(1 - p_i)e_i$$

- Moments of data:

$$E[Y_{ij}|p_i, e_i] = N_{ij}p_i$$

$$\text{var}(Y_{ij}|p_i, e_i) = N_{ij}p_i(1 - p_i)[1 + (N_{ij} - 1)e_i]$$

- As $e_i \rightarrow 0$ we approach the binomial model.
- As $e_i \rightarrow 1$ we have increased overdispersion (variability within gene).

Calibration Model for Genomic data

- **Stage 3: Within-Gene Likelihood:**

$$p_i|a \sim \text{beta}(a, a)$$

$$e_i|d \sim \text{beta}(1, d)$$

Note: prior on within-gene dispersion is monotonic decreasing from 0 (corresponding to no variability).

- **Stage 4: Hyperpriors:** Require priors on $a > 0, d > 0$.
- We take

$$a \sim \text{lognormal}(4.3, 1.8)$$

$$d \sim \text{exponential}(0.0001)$$

- The latter prior determines the within-gene variability within-gene variability in genomic DNA – chosen by examination of resultant p_{ij} 's.
- Separate a, d for each technology.

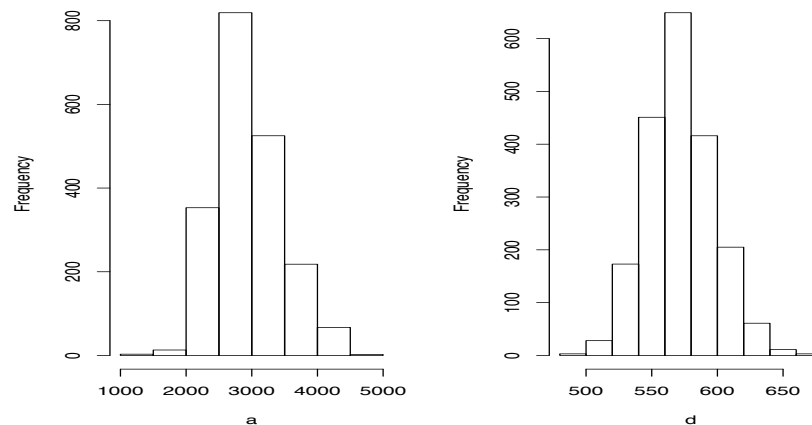


Figure 14 : Posteriors for the RNA-Seq data, Illumina platform.

Model 3: Model for RNA-Seq Data

- Data are modeled as a two-component mixture: the first “null” component having a known distribution, from the genomic DNA analysis on the filtered data.

- Stage 1: Within-Genes Likelihood:*

$$Y_{ij}|p_{ij} \sim \text{binomial}(N_{ij}, p_{ij}).$$

where p_{ij} is the probability of an outcome from the first genetic background.

- Stage 2: Within-Genes Prior:*

$$p_{ij}|\alpha_i, \beta_i \sim \text{beta}(\alpha_i, \beta_i)$$

so that α_i, β_i determine the distribution of variants within gene i .

- Stage 3: Between-Genes Prior:* We again reparameterize $(\alpha_i, \beta_i) \rightarrow (p_i, e_i)$:

$$p_i, e_i|f, g, h, \pi_0 \sim \begin{cases} \text{beta}(\hat{a}, \hat{a}) \times \text{beta}(1, \hat{d}) & \text{with probability } \pi_0 \\ \text{beta}(f, g) \times \text{beta}(1, h) & \text{with probability } 1 - \pi_0 \end{cases}$$

with \hat{a}, \hat{d} from genomic DNA analysis.

Stage 4: Hyperpriors: Require priors on $\pi_0, f > 0, g > 0, h > 0$.

- Uniform prior on π_0 .
- f and g describe beta distribution of p_i for genes displaying ASE – want this distribution to be centered around symmetry.
- Reparameterize as

$$q = \frac{f}{f+g} \quad r = \frac{1}{1+f+g}$$

so that $E[p_i] = q, \text{var}(p_i) = q(1-q)r$.

- Through experimentation:

$$q \sim \text{beta}(100, 100) \quad r \sim \text{beta}(1, 20)$$

- For h , the distribution of within-gene variability in ASE:

$$h \sim \text{exponential}(0.03).$$

- Separate f, g, h for each technology.

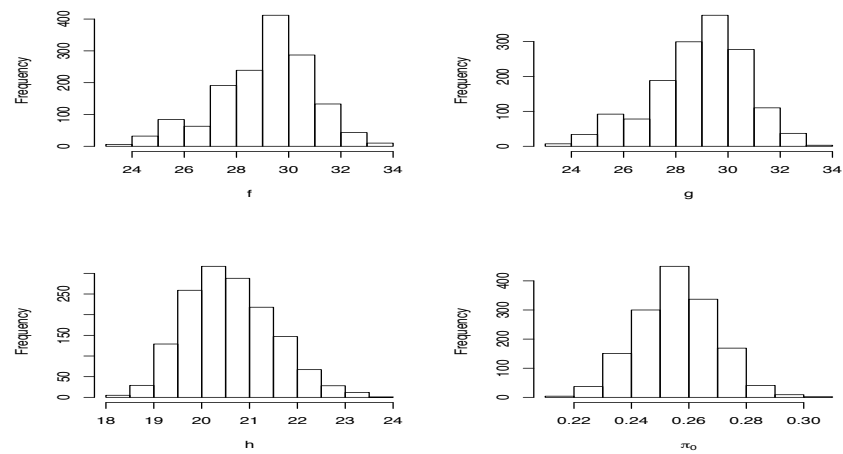


Figure 15 : Posteriors for the RNA-Seq data, Illumina platform.

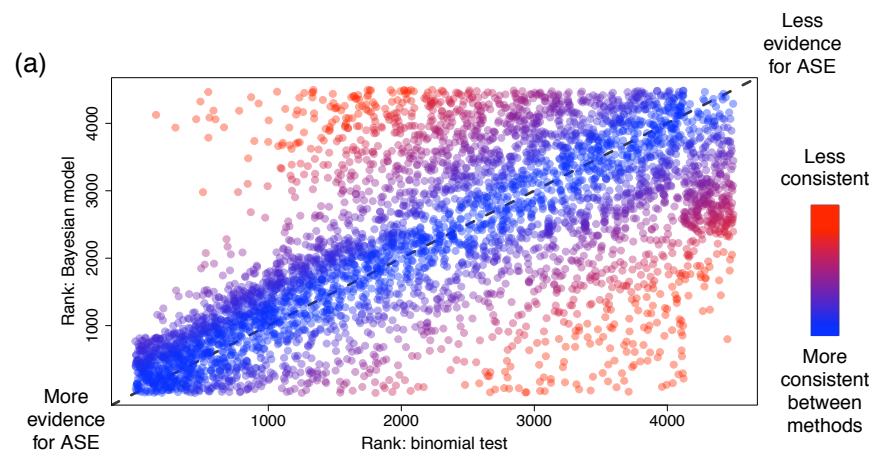


Figure 16 : Comparison of rankings from binomial test and hierarchical model.

Introduction ○○○	Imputation ○○○○○○○○○○○○○○○○○○○○	Hierarchical ASE ○○○○○○○○ ○○○○○○○○○○○○○○○○○○●○○○○○	Model Comparison ○○○○○○○○○○○○○	Conclusions ○	References
---------------------	------------------------------------	--	-----------------------------------	------------------	------------

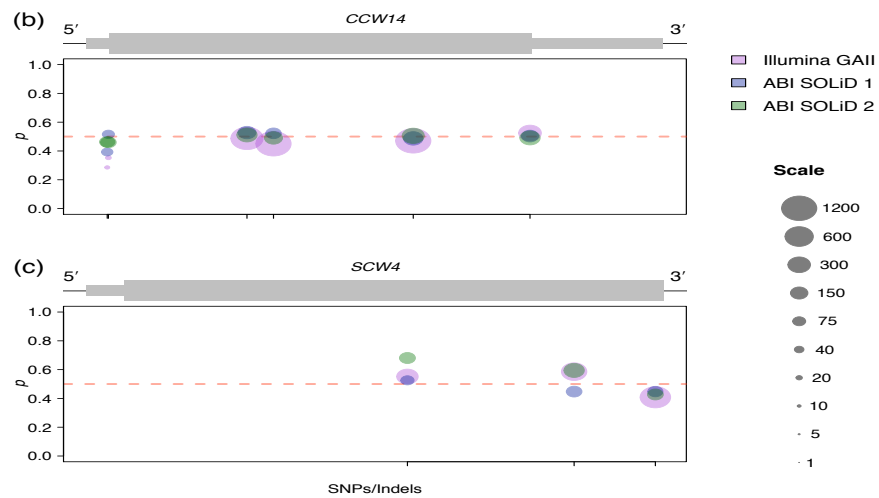


Figure 17 : Examples of opposite conclusions: In (b) the p -value said ASE and Bayes not. (large sample size, Bayes allows wobble). In (c) the p -value said no ASE, Bayes analysis yes.

Introduction ○○○	Imputation ○○○○○○○○○○○○○○○○○○○○	Hierarchical ASE ○○○○○○○○ ○○○○○○○○○○○○○○○○○○●○○○○○	Model Comparison ○○○○○○○○○○○○○	Conclusions ○	References
---------------------	------------------------------------	--	-----------------------------------	------------------	------------

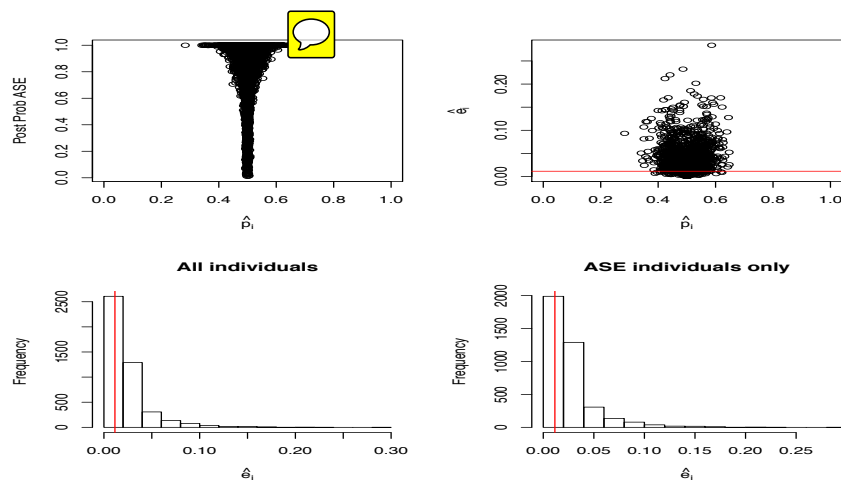


Figure 18 : Between-gene variability p_i and within-gene variability e_i .

Varying ASE within genes

- One mechanism: Imagine a gene with an exon and an intron, and that we have SNPs in both.
- At each exonic SNP we see approximately the same number of BY and RM reads.
- Now suppose the intron is not spliced out for the BY allele, but it is spliced out efficiently for the RM allele. At each intronic SNP we will still see the same number of BY reads as in the exon (everything else being equal), but approximately 0 RM reads, leading to variable ASE across the gene
- In the figure: The “thin” part of the gene (YML024W) is an intron, while the “thick” part is an exon.
- For the RM allele (magenta) the intron is not spliced out, while it is mostly spliced out in the BY allele (green).

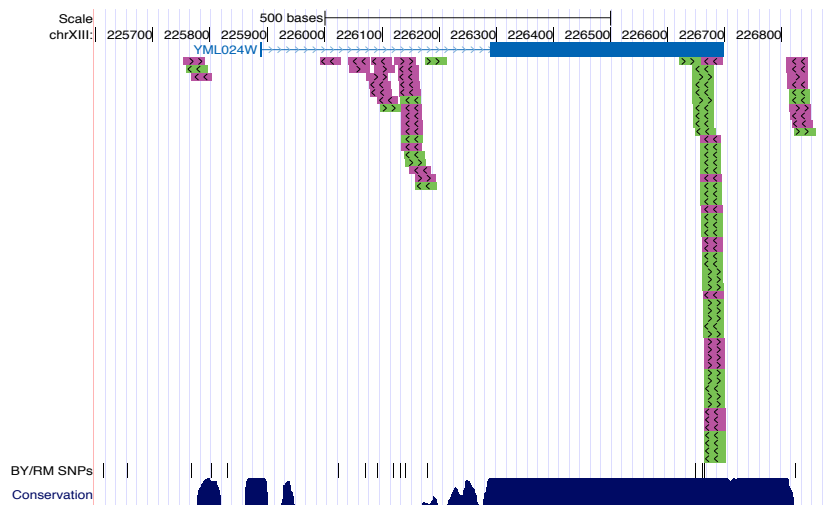


Figure 19 : Example of a gene displaying variable ASE within a gene. Green = RM, magenta = BY.


Introduction ○○○	Imputation ○○○○○○○○○○○○○○○○○○	Hierarchical ASE ○○○○○○○ ○○○○○○○○○○○○○○○○○○●	Model Comparison ○○○○○○○○○○○○○	Conclusions ○	References
---------------------	----------------------------------	--	-----------------------------------	------------------	------------

Conclusions for Mixture Model

- For the ASE data we used the DNA experiment to calibrate the prior.
- More details of this experiment and the model can be found in Skelly *et al.* (2011).
- Implementation was via Markov chain Monte Carlo, but we had to write our own code.

Introduction ○○○	Imputation ○○○○○○○○○○○○○○○○○○	Hierarchical ASE ○○○○○○○ ○○○○○○○○○○○○○○○○○○	Model Comparison ●○○○○○○○○○○○○○	Conclusions ○	References
---------------------	----------------------------------	---	------------------------------------	------------------	------------


Model Comparison

- Markov chain Monte Carlo in particular has allowed the fitting of more and more complex models, often hierarchical in nature with layers of random effects. 
- The search for a method to find the “best” of a set of candidate models has also grown.
- Let $p(\mathbf{y}|\boldsymbol{\theta})$ represent a generic likelihood for $\mathbf{y} = [y_1, \dots, y_n]$ and let

$$D(\boldsymbol{\theta}) = -2 \log[p(\mathbf{y}|\boldsymbol{\theta})]$$

represent the **deviance**.

- For example, in an iid normal($\mu_i(\boldsymbol{\theta}), \sigma^2$) normal the deviance is

$$\frac{1}{\sigma^2} \sum_{i=1}^n [y_i - \mu_i(\boldsymbol{\theta})]^2. \quad \text{$$

- Frequentist model comparison for nested models is often carried out using likelihood ratio statistics, which corresponds to the comparison of deviances in generalized linear models (GLMs), see for example McCullagh and Nelder (1989).

Model Comparison: AIC

- One approach to model comparison is based on a model's ability to make good **predictions**.
- Such an objective, and **predicting** the actual observed data, leads to Akaike's an information criterion (AIC), derived in Akaike (1973).
- In AIC one tries to estimate the (Kullback-Leibler) distance between the true distribution of the data, and the modeled distribution of the data.
- AIC is given by

$$\text{AIC} = -2 \log[p(y|\hat{\theta})] + 2k$$

where $\hat{\theta}$ is the MLE and k is the number of parameters in the model, i.e. the size of θ .

- Small values of the AIC are favored, since they suggest low prediction error.
- The **penalty term** $2k$ penalizes the double use of the data.
- In general for prediction: overly complex models are penalized since redundant parameters "use up" information in the data.

Model Comparison: BIC

- Another approach is based on trying to identify the "true" model.
- Schwarz (1978) developed the **Bayesian Information Criterion (BIC)** which is given by

$$\text{BIC} = -2 \log[p(y|\hat{\theta})] + k \log n.$$

- BIC approximates $-2 \log p(\mathbf{y}|\theta)$ under a certain unit information prior (Kass and Wasserman, 1995).
- BIC is **consistent**¹ for finding the true model, if that model lies in the set being compared.
- AIC is not consistent for finding the true model, but recall is intended for prediction.

¹meaning the BIC hones in on the true model as the sample size increases

Model Comparison: DIC

- Spiegelhalter *et al.* (2002) introduced what has proved to be a very popular model comparison statistic, the **deviance information criterion (DIC)**.
- To define the DIC, define an “effective number of parameters as

$$\begin{aligned}
 p_D &= E_{\theta|y}\{-2\log[p(y|\theta)]\} + 2\log[p(y|\bar{\theta})] \\
 &= \bar{D} + D(\bar{\theta})
 \end{aligned}$$

where $\bar{\theta} = E[\theta|y]$ is the posterior mean, $D(\bar{\theta})$ is the deviance evaluated at the posterior mean and $\bar{D} = E[D|y]$.

- Hence, p_D is the

posterior mean deviance – deviance of posterior means.

- The DIC is given by

$$\begin{aligned}
 \text{DIC} &= D(\bar{\theta}) + 2p_D \\
 &= \bar{D} + p_D,
 \end{aligned}$$

so that we have a measure of goodness of fit + complexity.

- DIC is straightforward to evaluate using MCMC or INLA.

Model Comparison: DIC

DIC has been heavily criticized (Spiegelhalter *et al.*, 2014):

- p_D is not invariant to parameterization.
- DIC is not consistent for choosing the correct model.
- DIC has a weak theoretical justification and is not universally applicable.
- DIC has been shown to under penalize complex models (Plummer, 2008; Ando, 2007).
- See Spiegelhalter *et al.* (2014) for an interesting discussion of the history of DIC, including a summary of attempts to improve DIC.
- According to Google Scholar, as of June 20th, 2014, Spiegelhalter *et al.* (2002) has 5251 citations...

Model Comparison: CPO

- Another approach based on prediction uses the conditional predictive ordinate (CPO).
- Let

$$\mathbf{y}_{-i} = [y_1, \dots, y_{i-1}, y_{i+1}, \dots, y_n]$$

represent the vector of data with the i -th observation removed.

- The idea is to predict the density ordinate of the left-out observation, based on those that remain.
- Specifically, the CPO for observation i is defined as:

$$\begin{aligned} \text{CPO}_i &= p(y_i | \mathbf{y}_{-i}) \\ &= \int p(y_i | \boldsymbol{\theta}) p(\boldsymbol{\theta} | \mathbf{y}_{-i}) d\boldsymbol{\theta} \\ &= E_{\boldsymbol{\theta} | \mathbf{y}_{-i}} [p(y_i | \boldsymbol{\theta})] \end{aligned}$$

Model Comparison: CPO

- The CPOs can be used to look at local fit, or one can define an overall score for each model:

$$\log(\text{CPO}) = \sum_{i=1}^n \log \text{CPO}_i.$$

- Good models will have relatively high values of $\log(\text{CPO})$.
- See Held *et al.* (2010) for a discussion of shortcuts for estimation (i.e. avoidance of fitting the model n times) using MCMC and INLA.

Model Comparison: Illustration, Childhood Mortality in Tanzania

- We illustrate the use of CPO and DIC in a study of estimating childhood (under 5) mortality in regions of Tanzania.
- The data are collected via a series of 8 surveys in 21 regions covering the period 1980–2009.
- Let q_{its} be the childhood mortality in area i , time point t from survey s .
- Based on the surveys we can obtain weighted (Horvitz-Thompson) estimators \hat{q}_{its} with associated asymptotic variances V_{its} .
- We summarize the data via logit estimates

$$y_{its} = \log \left(\frac{\hat{q}_{its}}{1 - \hat{q}_{its}} \right).$$

- Let

$$\phi_{its} = \log \left(\frac{q_{its}}{1 - q_{its}} \right)$$

represent the logit of the childhood mortality.

Model Comparison: Illustration, Childhood Mortality in Tanzania

We have a three-stage hierarchical model:

- **Stage 1:** Likelihood:

$$y_{its} | \phi_{its} \sim \text{normal}(\phi_{its}, V_{its}).$$

and we compare the following six models:

$$\text{Model 1: } \phi_{its} = \mu + \alpha_t + \gamma_t + \theta_i + \eta_i + \delta_{it}$$

$$\text{Model 2: } \phi_{its} = \mu + \alpha_t + \gamma_t + \theta_i + \eta_i + \delta_{it} + \nu_s$$

$$\text{Model 3: } \phi_{its} = \mu + \alpha_t + \gamma_t + \theta_i + \eta_i + \delta_{it} + \nu_s + \nu_{is}$$

$$\text{Model 4: } \phi_{its} = \mu + \alpha_t + \gamma_t + \theta_i + \eta_i + \delta_{it} + \nu_s + \nu_{ts}$$

$$\text{Model 5: } \phi_{its} = \mu + \alpha_t + \gamma_t + \theta_i + \eta_i + \delta_{it} + \nu_s + \nu_{ts} + \nu_{is}$$

$$\text{Model 6: } \phi_{its} = \mu + \alpha_t + \gamma_t + \theta_i + \eta_i + \delta_{it} + \nu_s + \nu_{ts} + \nu_{is} + \nu_{its}$$

where α_t , θ_i , δ_{it} are independent random effects for time, area and the interaction, γ_t and η_i are random effects that carry out local smoothing in time and space and ν_s , ν_{ts} , ν_{is} , ν_{its} are independent random effects to reflect survey effects.

- **Stage 2:** Normal random effects Distributions.
- **Stage 3:** Hyperpriors on μ and the random effects variances.

Model Comparison: Illustration, Childhood Mortality in Tanzania

Table 1 : Model comparison statistics for 6 models for the Tanzania data; “best” in red.

Model	No. Parameters	p_D	\bar{D}	DIC	log(CPO)
2	181	75	409	484	-295
2	189	81	382	463	-288
3	313	120	219	339	-193
4	223	91	364	454	-282
5	347	128	202	330	-182
6	920	149	185	334	-184

- Notice how much smaller the effective number of parameters is, when compared with the total number of parameteres; this is because of the shrinkage/penalization of the random effects distributions.
- Both CPO and DIC suggest that model 5 is the best:

$$\text{Model 5: } \phi_{its} = \mu + \alpha_t + \gamma_t + \theta_i + \eta_i + \delta_{it} + \nu_s + \nu_{ts} + \nu_{is}$$

- So survey effects vary across time and across areas (different teams sent out).

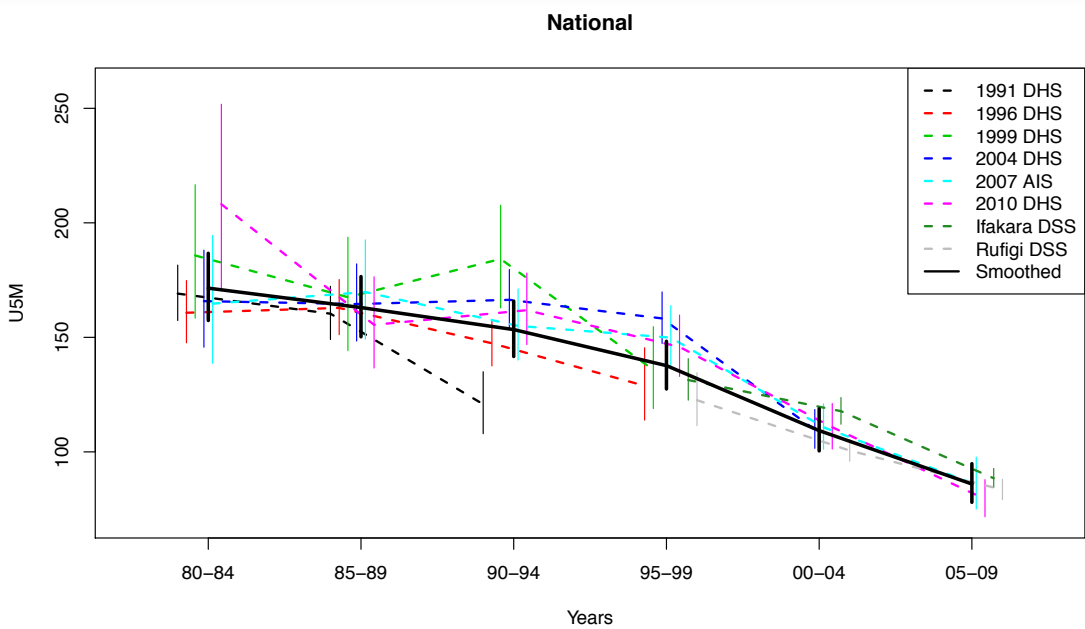


Figure 20 : Smoothed estimates of national under 5 mortality in Tanzania (solid line), different surveys denoted with dashed lines and vertical lines represent 95% interval estimates.

Introduction ○○○	Imputation ○○○○○○○○○○○○○○○○○○	Hierarchical ASE ○○○○○○○ ○○○○○○○○○○○○○○○○○○○○	Model Comparison ○○○○○○○○○○○○○	Conclusions ●	References
---------------------	----------------------------------	---	-----------------------------------	------------------	------------

Conclusions

- Hierarchical models allow complex dependencies within data to be modeled.
- Prior specification for variance components is not straightforward, and sensitivity analysis is a good idea.
- No universally agreed upon approach to carrying out model comparison.

Introduction ○○○	Imputation ○○○○○○○○○○○○○○○○○○	Hierarchical ASE ○○○○○○○ ○○○○○○○○○○○○○○○○○○○○	Model Comparison ○○○○○○○○○○○○○	Conclusions ○	References
---------------------	----------------------------------	---	-----------------------------------	------------------	------------

References

- Akaike, H. (1973). Information theory and an extension of the maximum likelihood principle. In P. B.N. and C. F., editors, *Second International Symposium on Information Theory*, pages 267–281. Akademia Kiadó, Budapest.
- Ando, T. (2007). Bayesian predictive information criterion for the evaluation of hierarchical bayesian and empirical bayes models. *Biometrika*, **94**, 443–458.
- Browning, B. and Browning, S. (2009). A unified approach to genotype imputation and haplotype-based inference for large data sets of trios and unrelated individuals. *The American Journal of Human Genetics*, **84**, 1084–1097.
- Held, L., Schrödle, B., and Rue, H. (2010). Posterior and cross-validated predictive checks: A comparison of MCMC and INLA. In T. Kneib and G. Tutz, editors, *Statistical Modeling and Regression Structures – Festschrift in Honour of Ludwig Fahrmeir*, pages 91–110. Physica-Verlag.
- Howie, B., Donnelly, P., and Marchini, J. (2009). A flexible and accurate genotype imputation method for the next generation of genome-wide association studies. *PLoS Genetics*, **5**, e1000529.
- Huang, L., Yun, L., Singleton, A., Hardy, J., Abecasis, G., Rosenberg, N., and Scheet, P. (2009). Genotype-imputation accuracy across worldwide human populations. *The American Journal of Human Genetics*, **84**, 235–250.

Introduction ○○○	Imputation ○○○○○○○○○○○○○○○○○○	Hierarchical ASE ○○○○○○○ ○○○○○○○○○○○○○○○○○○○○	Model Comparison ○○○○○○○○○○○○○	Conclusions ○	References
---------------------	----------------------------------	---	-----------------------------------	------------------	------------

Kass, R. and Wasserman, L. (1995). A reference Bayesian test for nested hypotheses and its relationship to the schwarz criterion. *Journal of the American Statistical Association*, **90**, 928–934.

Li, Y., Willer, C., Ding, J., Scheet, P., and Abecasis, G. (2010). Mach: using sequence and genotype data to estimate haplotypes and unobserved genotypes. *Genetic Epidemiology*, **34**, 816–834.

Marchini, J. and Howie, B. (2010). Genotype imputation for genome-wide association studies. *Nature Reviews Genetics*, **11**, 499–511.

Marchini, J., Howie, B., Myers, S., McVean, G., and Donnelly, P. (2007). A new multipoint method for genome-wide association studies by imputation of genotypes. *Nature Genetics*, **39**, 906–913.

McCullagh, P. and Nelder, J. (1989). *Generalized Linear Models, Second Edition*. Chapman and Hall, London.

Plummer, M. (2008). Penalized loss functions for bayesian model comparison. *Biostatistics*, **9**, 523–539.

Sanna, S., Li, B., and Mulas, A. (2011). Fine mapping of five loci associated with low density lipoprotein cholesterol fetects variants that double the explained heritability. *PLoS Genetics*, **7**, 1002198.

Scheet, P. and Stephens, M. (2006). A fast and flexible statistical model for large-scale population genotype data: applications to inferring missing genotypes and haplotypic phase. *The American Journal of Human Genetics*, **78**, 629–644.

Introduction ○○○	Imputation ○○○○○○○○○○○○○○○○○○	Hierarchical ASE ○○○○○○○ ○○○○○○○○○○○○○○○○○○○○	Model Comparison ○○○○○○○○○○○○○	Conclusions ○	References
---------------------	----------------------------------	---	-----------------------------------	------------------	------------

Schwarz, G. (1978). Estimating the dimension of a model. *Annals of Statistics*, **6**, 461–464.

Skelly, D., Johansson, M., Madeoy, J., Wakefield, J., and Akey, J. (2011). A powerful and flexible statistical framework for testing hypothesis of allele-specific gene expression from RNA-Seq data. *Genome Research*, **21**, 1728–1737.

Spiegelhalter, D., Best, N., Carlin, B., and Linde, A. V. D. (2002). Bayesian measures of model complexity and fit. *Journal of the Royal Statistical Society: Series B*, **64**, 583–639.

Spiegelhalter, D., Best, N., Carlin, B., and Linde, A. V. D. (2014). The deviance information criterion: 12 years on (with discussion). *Journal of the Royal Statistical Society: Series B*, **64**, 485–493.