

2019 SISG MODULE 8: Bayesian Statistics for Genetics

Lecture 9: Bayesian and Frequentist Testing

Jon Wakefield

Departments of Statistics and Biostatistics
University of Washington

Review of p -Values and Bayes Factors

Multiple Testing

Conclusions

Appendix: Substantive Prior Information

Review of p -Values and Bayes Factors

The Statistical Set-Up

We review frequentist and Bayesian test procedures.

- ▶ We begin with a very simple situation in which we have a single parameter of interest θ .
- ▶ Assume the null of interest is

$$H_0 : \theta = 0$$

with θ , for example, a treatment difference, or a log odds ratio, or a log hazard ratio.

- ▶ We assume an analysis yields a statistic T for which large values indicate departures from the null.
- ▶ For example, the squared **Wald statistic**, $T = \hat{\theta}^2 / V$, with V the asymptotic variance of the MLE¹.
- ▶ An alternative is the **likelihood ratio statistic**.

¹ $T=Z^2$ where Z is the **Z-score**

Types of Testing

- ▶ The observed p -value is given by:

$$p = \Pr(T > t_{\text{obs}} | H_0)$$

where t_{obs} is a number that is evaluated for the data at hand.

- ▶ To report p only, gives a **pure significance test**.
 - ▶ A small p -value can arise because:
 - ▶ H_0 is true but we were “unlucky”.
 - ▶ H_0 is not true.
- to decide which explanation is responsible depends crucially on the **prior** belief on whether H_0 is true or not.

Key question: How small is small?

Types of Testing

- ▶ A **test of significance** sets a cut-off value (e.g. $\alpha = 0.05$) and rejects H_0 if $p < \alpha$.

Again: How to pick α ?

- ▶ A type I error is to reject H_0 when it is true, and a test of significance controls the type I error (whereas a pure significance test does not).
- ▶ A type II error occurs when H_1 is true but H_0 is not rejected.
- ▶ A **hypothesis test** goes one step further and specifies an alternative hypothesis.
- ▶ A decision is then taken as to which of H_0 and H_1 is chosen.
- ▶ The celebrated **Neyman-Pearson lemma** shows that for fixed α -level the likelihood ratio statistic maximizes the power.
- ▶ Wouldn't it be more reasonable to **balance** type I and type II errors?

The Dangers of Fixed Significance Levels

- ▶ **Example:** Sample, Y_1, \dots, Y_n of size n from $N(\theta, 1)$,

$$H_0 : \theta = 0, \quad H_1 : \theta = 1.$$

Obvious that we should reject H_0 for $\bar{Y}_n > k(n)$, a constant².

- ▶ The table below illustrates the problems of choosing a fixed α , regardless of sample size — **imbalance** in α and β as a function of n :

n	α	β	$k(n)$
1	0.01	0.91	2.33
25	0.01	0.0038	0.46
100	0.01	8×10^{-15}	0.23

- ▶ **Also:** Statistical versus practical significance.
- ▶ For both p -values and α levels we need thresholds that **decrease** as a function of the sample size n . Pearson (1953, p. 68), “...the quite legitimate device of reducing α as n increases”.

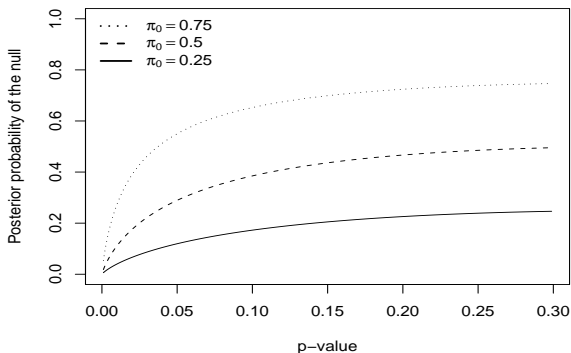
²Note that the threshold for $T = [\bar{Y}_n / (1/\sqrt{n})]^2$ is constant

A quite remarkable result!

- ▶ With $\pi_0 = \Pr(H_0)$, Sellke *et al.* (2001) show that:

$$\Pr(H_0 | \text{data}) \geq \left\{ 1 - \frac{1}{2.72 p \log p} \times \frac{1 - \pi_0}{\pi_0} \right\}^{-1} \quad (1)$$

- ▶ A small p -value doesn't translate to a small probability that the null is not true.



Why does anyone use p -values?

- ▶ Historically, it was usual to carry out well-powered (single) experiments, and the prior on the alternative was not small.
- ▶ With respect to (1) and with $\pi_0 = 0.5$:
 - ▶ $p\text{-value} = 0.05$ gives $\Pr(H_0 | \text{data}) > 0.29$.
 - ▶ $p\text{-value} = 0.01$ gives $\Pr(H_0 | \text{data}) > 0.11$.
- ▶ Scientists well-calibrated in their own discipline?
- ▶ Perhaps, but if you're going to be subjective, why not be formal about it?
- ▶ **Aside:** Reason for lack of replication in observational epidemiology? Along with confounding, data dredging, measurement error,...

Calibrating α -Levels

- ▶ We want $\Pr(H_0 | \text{data})$, where “data” corresponds to the event $T > t_{\text{fix}}$, but to obtain this we must specify alternatives – consider a simple alternative, say $H_1 : \theta = \theta_1$.
- ▶ Then,

$$\begin{aligned} \text{Posterior Odds of } H_0 &= \frac{\Pr(H_0 | \text{data})}{\Pr(H_1 | \text{data})} \\ &= \frac{\Pr(T > t_{\text{fix}} | H_0)}{\Pr(T > t_{\text{fix}} | H_1)} \times \frac{\Pr(H_0)}{\Pr(H_1)} \\ &= \frac{\alpha}{1 - \beta} \times \text{Prior Odds of } H_0 \end{aligned}$$

- ▶ For **ranking** associations (which does not involve the prior odds if constant across tests): must consider the **power**, $\Pr(\text{data} | H_1)$.
- ▶ For **calibration**: must consider the **prior odds** of H_0 .

A Sanity Check via a Simple Example

- ▶ The model:

$$Y_i | \theta \sim_{iid} \mathbf{N}(\theta, \sigma^2), \quad \sigma^2 \text{ known,}$$

$$i = 1, \dots, n.$$

- ▶ The distribution of the **MLE** is:

$$\hat{\theta} = \bar{Y} \sim \mathbf{N}(\theta, V)$$

with $V = \sigma^2/n$,

$$T = \frac{n\bar{Y}^2}{\sigma^2}.$$

- ▶ Null and alternative hypotheses are

$$H_0 : \theta = 0, \quad H_1 : \theta \neq 0.$$

A Sanity Check via a Simple Example

- ▶ Under H_1 assume the prior $\theta \sim N(0, W)$.
- ▶ Recall from previous lectures that the evidence in the data for a pair of hypotheses is summarized in the **Bayes factor**:

$$\text{BF} = \frac{p(\mathbf{y}|H_0)}{p(\mathbf{y}|H_1)} = \frac{\prod_{i=1}^n \mathbf{N}(y_i|0, \sigma^2)}{\int_{\theta} \prod_{i=1}^n \mathbf{N}(y_i|\theta, \sigma^2) \times \mathbf{N}(\theta|0, W) d\theta}$$

where \mathbf{N} is shorthand for the density of a normal random variable.

Testing: decision theory

A reminder of the ingredients for decision theory;

- ▶ **Loss function** $L(\theta, d)$: how bad it would be if the truth were θ but you took decision d . (Optimists: note we could equivalently define **utility** as $-L(\theta, d)$ — how good it would be — economists do this)
- ▶ **Expected posterior loss** $E[L(\theta, d)]$ — loss for some decision d averaged over posterior uncertainty

The Bayes rule is the decision d that minimizes $E[L(\theta, d)]$ — but for testing, d is 0 or 1, so this means checking whether

$$E[L(\theta, d = 0)] \leq E[L(\theta, d = 1)],$$

i.e., do we expect less loss deciding $d = 0$ or $d = 1$?

		Truth	
		$\theta = 0$	$\theta \neq 0$
Decision	$d = 0$	0	L_1
	$d = 1$	L_2	0

With respect to this table, the posterior expected cost associated with the decision d is

$$E[L(\theta, d)] = L(\theta = 0, d) \Pr(\theta = 0|\mathbf{y}) + L(\theta \neq 0, d) \Pr(\theta \neq 0|\mathbf{y}).$$

The two possible decisions (report $\theta = 0$ or $\theta \neq 0$) the **expected losses** are:

$$E[L(\theta, d = 0)] = 0 \times \Pr(\theta = 0|\mathbf{y}) + L_2 \Pr(\theta \neq 0|\mathbf{y})$$

$$E[L(\theta, d = 1)] = L_1 \Pr(\theta = 0|\mathbf{y}) + 0 \times \Pr(\theta \neq 0|\mathbf{y})$$

Testing

We now have to find the decision that minimizes the posterior expected loss, as a function of $\Pr(\theta \neq 0|\mathbf{y}) = \Pr(\theta|\mathbf{y})$.

A little rearrangement leads to reporting $\theta \neq 0$ if

$$\Pr(\theta \neq 0|\mathbf{y}) \geq \frac{L_1}{L_1 + L_2} = \frac{1}{1 + L_2/L_1} = \frac{1}{1 + R},$$

or equivalently

$$\Pr(\theta = 0|\mathbf{y}) < \frac{1}{1 + R}.$$

Examples:

If $L_1 = L_2$ ($R = 1$), report $\theta \neq 0$ if $\Pr(\theta \neq 0|\mathbf{y}) \geq \frac{1}{2}$.

If $L_1 = 3 \times L_2$ ($R = 1/3$), report $\theta \neq 0$ if $\Pr(\theta \neq 0|\mathbf{y}) \geq \frac{3}{4}$.

If $L_2 = 3 \times L_1$ ($R = 3$), report $\theta \neq 0$ if $\Pr(\theta \neq 0|\mathbf{y}) \geq \frac{1}{4}$.

A Sanity Check via a Simple Example

- ▶ We take $W = \sigma^2$, which corresponds to the “unit information prior” of Kass and Wasserman (1995) (this choice not so important).
- ▶ With a prior odds, PO, and ratio of costs of type II to type I errors, R , this gives the decision rule to reject H_0 :

$$\begin{aligned}\text{Posterior Odds} &= \text{BF} \times \text{PO} \\ &= \sqrt{1+n} \times \exp\left(-\frac{T}{2} \frac{n}{1+n}\right) \times \text{PO} < R\end{aligned}$$

- ▶ Notice how this depends on T and n .

A Bayesian Test Statistics Threshold

- ▶ Rearrangement gives a threshold for rejection of:

$$T > \frac{2(1+n)}{n} \log \left(\frac{\text{PO}}{R} \sqrt{1+n} \right)$$

- ▶ For relatively large prior odds on the null **PO**: require **T** to be larger (more evidence).
- ▶ For relatively large cost of Type II errors **R** (so that we are averse to type II error, i.e. missing signals): require **T** to be smaller (less evidence).
- ▶ Not such a simply summarization for **n** but, beyond a certain point, as **n** gets larger, we require larger **T** (more evidence).
- ▶ The above should be contrasted with the usual frequentist approach of

$$T > \text{const}$$

with the constant usually chosen to control the type I error.

A Bayesian Test Statistic Threshold

- ▶ The table below evaluates the probability of rejection given H_0 . We assume $R = 1$.
- ▶ For $\pi_0 = 0.5$ and $n = 20, 50, 100$ the thresholds give ≈ 0.05 — the situation in which this infamous threshold was first derived?

	$\pi_0 = 0.25$	$\pi_0 = 0.50$	$\pi_0 = 0.95$
$n = 10$	0.64	0.10	0.0025
$n = 20$	0.35	0.074	0.0022
$n = 50$	0.18	0.045	0.0016
$n = 100$	0.12	0.031	0.0011
$n = 1000$	0.030	0.0085	0.00034

Calibration with p -values

- ▶ The ABF can be inverted to give a rule for Z^2 that depends on PO, R and n (as with the simple example presented previously).
- ▶ For more details, see Wakefield (2009).
- ▶ Figure 1 shows the behavior of this rule as a function of the sample size n , and for different choices of the prior on the alternative π_1 and the ratio of costs of type II to type I errors.

- ▶ The curves have the expected ordering and, as n gets large, a greater and greater level of evidence is required.

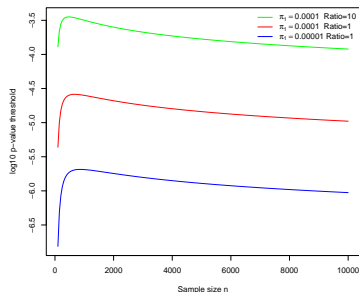


Figure 1: Threshold for rejection, on the $\log_{10}(p)$ -value scale, vs sample size.

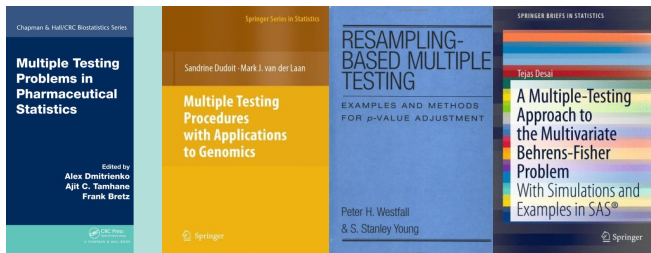
This is as we would expect because as the sample size increases we want both Type I and Type II errors to go to zero.

Multiple Testing

Motivation for Multiple Testing

We have covered testing procedures, both frequentist and Bayesian, in the context of single tests.

How to proceed, when multiple tests are envisaged, is a big topic:



A lot of interest lately, given the advent of technologies that allow huge numbers of experiments to be performed.

As with testing, this topic is **controversial**.

Motivating Example

- ▶ We follow a running example with data from a microarray study of 102 men, 52 with prostate cancer and 50 normal controls (Efron and Hastie, 2016).
- ▶ Gene expression levels were measured for $m = 6033$ genes.
- ▶ A two-standard t-test was carried out.

Motivating Example

- ▶ A transformation was made so that the resultant statistic z_i , has distribution under the null:

$$H_{0i} : z_i \sim N(0, 1),$$

for $i = 1, \dots, m$ genes.

- ▶ Under the alternative:

$$H_{1i} : z_i \sim N(\mu_i, 1),$$

for $i = 1, \dots, m$ genes.

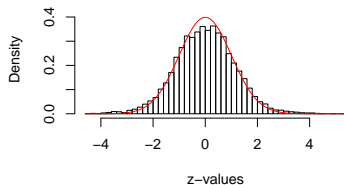


Figure 2: Histogram of z-values for prostate microarray study, with $N(0, 1)$ distribution in red.

- ▶ The aim is to find genes with non-zero μ_i .

Framework for Multiple Testing

Possibilities with m tests and when K are flagged as requiring further attention:

	Non-Flagged	Flagged	
H_0	A	B	m_0
H_1	C	D	m_1
	$m - K$	K	m

- ▶ m_0 is the number of true nulls.
- ▶ B is the number of type I errors.
- ▶ C is the number of type II errors.

Problem: To select a rule that will determine K .

We discriminate between:

- ▶ A sensible **criterion**.
- ▶ How the criterion should **depend on sample size**.

The Family-Wise Error Rate

	Non-Flagged	Flagged	
H_0	A	B	m_0
H_1	C	D	m_1
	$m - K$	K	m

- ▶ The **family-wise error rate** (FWER) is the probability of making at least one Type I error, i.e.

$$\Pr(B \geq 1 \mid \text{all } H_0 \text{ true}).$$

- ▶ Let B_i be the event that the i -th null is incorrectly rejected, so that $B = \cup_{i=1}^m B_i$ is the total number of incorrectly rejected nulls.

The Family-Wise Error Rate

- ▶ The **FWER** is given by:

$$\begin{aligned}\text{FWER} = \Pr(B \geq 1 \mid \text{all } H_0 \text{ true}) &= \Pr(\cup_{i=1}^m B_i \mid \text{all } H_0 \text{ true}) \\ &\leq \sum_{i=1}^m \Pr(B_i \mid \text{all } H_0 \text{ true}) \\ &= m\alpha^*\end{aligned}$$

where α^* is the level for each test.

- ▶ This is true regardless of whether the tests are independent or not.
- ▶ Bonferroni takes $\alpha^* = \alpha/m$ to give $\text{FWER} \leq \alpha$.
- ▶ **Example:** For control at $\alpha = 0.05$ with $m = 500K$ tests take $\alpha^* = 0.05/500,000 = 10^{-7}$.
- ▶ Such stringent rules lead to a loss of power, but not ridiculous if you think there is a reasonable chance that all nulls could be true (but α should depend on n , in particular should decrease as n gets larger and larger).

Sidak Correction

- ▶ If tests are **independent**:

$$\begin{aligned}\Pr(B \geq 1) &= 1 - \Pr(B = 0) \\ &= 1 - \Pr(\cap_{i=1}^m B'_i) \\ &= 1 - \prod_{i=1}^m \Pr(B'_i) \\ &= 1 - (1 - \alpha^*)^m\end{aligned}$$

- ▶ So to achieve $\text{FWER} = \alpha$ take p -value threshold as $\alpha^* = 1 - (1 - \alpha)^{1/m}$ — the **Sidak correction** (Sidak, 1967).
- ▶ **Example**: with $m = 500K$ tests take

$$\alpha^* = 1 - (1 - 0.05)^{1/500,000} = 1.03 \times 10^{-7}.$$

Bayes Bonferroni

There is a prior that results in a Bayesian Bonferroni-type correction³.

If the prior probabilities of each of the nulls are independent with $\pi_{0i} = \pi_0$ for $i = 1, \dots, m$.

Then the prior probability that all nulls are true is

$$\Pi_0 = \Pr(H_1 = 0, \dots, H_m = 0) = \pi_0^m$$

which we refer to as **prior P_1** , and let $\alpha_{i,B}$ be the posterior probability of the null under this prior for gene i .

Example if $\pi_0 = 0.5$ and $m = 10$, $\Pi_0 = 0.00098$, which may not reflect the required prior belief.

³The following describes a very idealized setting where the data model and prior are both normal

Suppose instead that we wish to fix the prior probability that all of the nulls are true at Π_0 .

A simple way of achieving this is to take $\pi_{0i} = \Pi_0^{1/m}$, a specification we call prior P_2 .

Westfall *et al.* (1995) show that for independent tests

$$\begin{aligned}\alpha_{i,B}^* &= \Pr(H_i = 0 \mid \mathbf{y}_i, P_2) \\ &\approx m \times \Pr(H_i = 0 \mid \mathbf{y}_i, P_1) \\ &= m \times \alpha_{i,B}.\end{aligned}$$

So a Bayesian version of a Bonferroni-like result is recovered.

As we have seen before, the posterior probability on the null, is strongly dependent on the prior on the null.

Expected Number of False Discoveries

We describe an alternative criterion.

For $i = 1, \dots, m$ tests let B_i again be the 1/0 random variable representing whether the null was incorrectly rejected or not, so that $B = \cup_{i=1}^m B_i$.

The **expected number of false discoveries** (EFD), with significance level α for each test, is given by

$$\text{EFD} = E[B] = \sum_{i=1}^m E[B_i] = m\alpha$$

if all nulls are true.

Expected False Discoveries

For m_0 true nulls: $E[B] = m_0\alpha$, but m_0 is unknown, so all we can say is

$$\text{EFD} = E[B] \leq m\alpha.$$

- ▶ In a GWAS context suppose $m = 500K$ and $\alpha = 0.05$; this gives $\text{EFD} \leq 25,000$, so conventional levels will clearly not work!
- ▶ We can easily put an upper bound on the EFD.
- ▶ For example, if we set $\alpha = 1/m$ the expected number of false discoveries is **bounded** by 1.
- ▶ With $\alpha = 5/m$ the expected number of false discoveries is **bounded** by 5.
- ▶ Compare to Bonferroni which controls the FWER via α/m .

False Discovery Rate

A very popular criterion is the **false discovery rate (FDR)**.

	Non-Flagged	Flagged	
H_0	A	B	m_0
H_1	C	D	m_1
	$m - K$	K	m

Define the false discovery proportion (FDP) as the **proportion of incorrect rejections**:

$$\text{FDP} = \begin{cases} \frac{B}{K} & \text{if } K > 0 \\ 0 & \text{if } K = 0 \end{cases}$$

Then the **false discovery rate (FDR)**, the **expected proportion of rejected nulls that are actually true nulls**, is given by

$$\text{FDR} = E[\text{FDP}].$$

False Discovery Rate

We describe an algorithm for **controlling the FDR**.

Consider the following procedure for independent p -values:

1. Let $P_{(1)} < \dots < P_{(m)}$ denote the ordered p -values.
2. Define $l_i = i\alpha/m$ and $R = \max\{i : P_{(i)} < l_i\}$ where α is the value at which we would like FDR control.
3. Then define the p -value threshold as $P_T = P_{(R)}$.
4. Reject all H_{0i} for which $P_i \leq P_T$.

Benjamini and Hochberg (1995) show that if this procedure is applied, then regardless of how many nulls are true (m_0) and regardless of the distribution of the p -values when the null is false

$$\text{FDR} \leq \frac{m_0}{m} \alpha < \alpha.$$

This algorithm was originally proposed by Simes (1986) to control the FWER.

Holm's Procedure

Holm's procedure Holm (1979) offers a modest improvement over Bonferroni.

Let

$$p_{(1)} \leq p_{(2)} \leq \cdots \leq p_{(i)} \leq \cdots \leq p_{(m)},$$

with corresponding null hypotheses $H_{0(i)}$.

Then, proceed as follows:

1. Let i_0 be the smallest index i such that

$$p_{(i)} > \frac{\alpha}{m - i + 1}.$$

2. **Reject** all null hypotheses $N_{0(i)}$ for $i < i_0$ and **accept** all with $i \geq i_0$.

It can be shown that Holm's procedure controls FWER at level α and is slightly less conservative.

Prostate Cancer Example

- ▶ We begin by plotting, in Figure 3 the observed p -values versus those expected under the null, i.e. $i/(m+1)$ for $i = 1, \dots, m = 6033$.
- ▶ Hard to tell what is going on here...

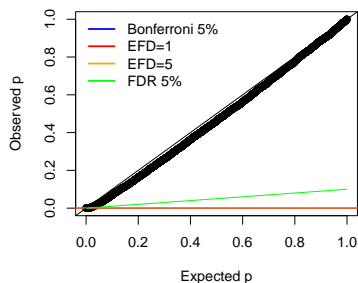


Figure 3: Observed versus expected p -values.

Prostate Cancer Example

- ▶ We stretch the scale in Figure 5 to show $-\log_{10}$ the observed p -values versus expected p -values.
- ▶ On this scale, a value of 2 corresponds to a p -value of 0.01, and a value of 3 corresponds to a p -value of 0.001.
- ▶ We see that the **FWER** is very conservative ($p = 0.05/m = 8.3 \times 10^{-6}$, or $-\log_{10}(p) = 5.1$) and only flags 3 genes as being significant (Holm's procedure gives the same 3).

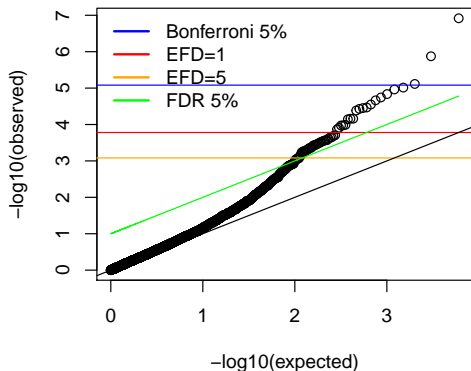


Figure 4: Observed versus expected p -values, on $-\log_{10}$ scale.

Prostate Cancer Example

- ▶ The **EFD=1** gives a p -value threshold of $1/6033 = 0.00017$, or $-\log_{10}(p) = 3.78$ and gives 21 flagged genes.
- ▶ The **EFD=5** gives a p -value threshold of $5/6033 = 0.00083$, or $-\log_{10}(p) = 3.08$ and gives 54 flagged genes.
- ▶ The **FDR** control at 5% gives the green diagonal line and flags 21 genes.

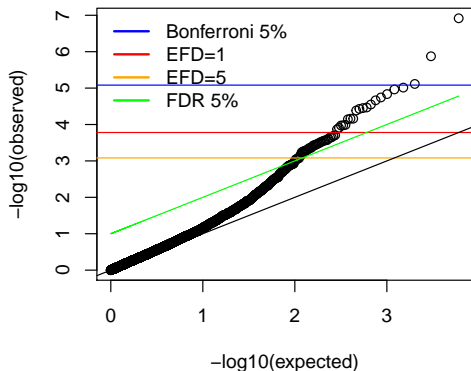


Figure 5: Observed versus expected p -values, on $-\log_{10}$ scale.

If all the nulls are true then $B = K$ (all rejections are false) and

$$\text{FDR} = E \left[\frac{B}{K} \right] = 1 \times \Pr(B \geq 1) = \text{FWER}.$$

FDR in this form and with extensions, e.g. Storey and Tibshirani (2003) has been successfully used in the microarrays field, where the number of non-null associations is not small.

Unfortunately less successful in a GWAS, because the proportion of nulls is very close to 1.

The algorithm of Benjamini and Hochberg (1995) begins with a desired FDR and then provides the p -value threshold.

Storey (2002) proposed an alternative method by which, for any fixed rejection region, a criteria closely related to FDR, the **positive false discovery rate**

$$\text{pFDR} = E[B/K \mid K > 0],$$

may be estimated⁴.

We assume rejection regions of the form $T > t_{\text{fix}}$ and consider the pFDR associated with regions of this form, which we write as $\text{pFDR}(t_{\text{fix}})$.

⁴this handles the event $K = 0$ differently to the previously-defined FDR

We define, for $i = 1, \dots, m$ tests, the random variables $H_i = 0/1$ corresponding to null/alternative hypotheses and test statistics T_i .

Then, with $\pi_0 = \Pr(H = 0)$ and $\pi_1 = 1 - \pi_0$ independently for all tests:

$$\text{pFDR}(t_{\text{fix}}) = \frac{\Pr(T > t_{\text{fix}} \mid H = 0) \times \pi_0}{\Pr(T > t_{\text{fix}} \mid H = 0) \times \pi_0 + \Pr(T > t_{\text{fix}} \mid H = 1) \times \pi_1}.$$

Consideration of the **false discovery odds**:

$$\frac{\text{pFDR}(t_{\text{fix}})}{1 - \text{pFDR}(t_{\text{fix}})} = \frac{\Pr(T > t_{\text{fix}} \mid H = 0)}{\Pr(T > t_{\text{fix}} \mid H = 1)} \times \frac{\pi_0}{\pi_1}$$

explicitly shows the weighted trade-off of type I and type II errors, with weights determined by the prior on the null/alternative.

Storey (2003) rigorously shows that

$$\text{pFDR}(t_{\text{fix}}) = \Pr(H = 0 \mid T > t_{\text{fix}}).$$

giving a Bayesian interpretation.

In terms of p -values, the rejection region corresponding to $T > t_{\text{fix}}$ is of the form $[0, \gamma]$.

Let P be the random p -value resulting from a test.

Under the null, $P \sim U(0, 1)$, and so

$$\begin{aligned} \text{pFDR}(t_{\text{fix}}) &= \frac{\Pr(P \leq \gamma \mid H = 0) \times \pi_0}{\Pr(P \leq \gamma)} \\ &= \frac{\gamma \times \pi_0}{\Pr(P \leq \gamma)}. \end{aligned} \tag{2}$$

From this expression, the crucial role of π_0 is evident.

q-values

- ▶ Storey (2002) estimates (2), using uniformity of p -values under the null, to produce the estimates

$$\hat{\pi}_0 = \frac{\#\{p_i > \lambda\}}{m(1 - \lambda)} \quad (3)$$

$$\hat{Pr}(P \leq \gamma) = \frac{\#\{p_i \leq \gamma\}}{m} \quad (4)$$

with λ chosen via the bootstrap to minimize the mean-squared error for prediction of the pFDR.

- ▶ The expression (3) calculates the empirical proportion of p -values to the right of λ , and then inflates this to account for the proportion of null p -values in $[0, \lambda]$.

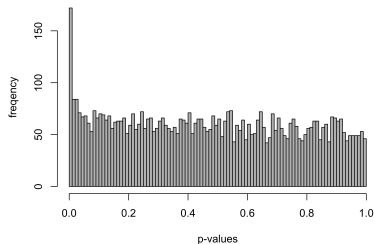


Figure 6: Histogram of p -values for prostate cancer example.

- ▶ π_0 is estimated as 0.854 for the prostate cancer data.
- ▶ 71 genes flagged at 10% FDR level.

This method highlights the benefits of allowing the **totality of p -values** to estimate fundamental quantities of interest such as π_0 .

The **q -value** is the minimum FDR that can be attained when a particular test is called significant.

We give a derivation of the q -value and, following Storey (2002), first define a set of nested rejection regions

$$\{t_\alpha\}_{\alpha=0}^1$$

where α is such that

$$\Pr(T > t_\alpha \mid H = 0) = \alpha.$$

Then,

$$p\text{-value}(t) = \inf_{t_\alpha: t \in t_\alpha} \Pr(T > t_\alpha \mid H = 0)$$

is the p -value corresponding to an observed statistic t .

The q -value is defined as

$$q\text{-value}(t) = \inf_{t_\alpha: t \in t_\alpha} \Pr(H = 0 \mid T > t_\alpha) \quad (5)$$

Therefore, for each observed statistic t_j there is an associated q -value.

The q -value is the minimum FDR that can be attained when calling that feature significant.

The q -values are estimated as

$$\hat{q}_i(p_i) = \min_{t \geq p_i} \widehat{\text{pFDR}}(t).$$

Note that in some papers (Storey and Tibshirani, 2003) the q -values is defined in terms of the FDR, since often m is large and $\Pr(D > 0) \approx 1$ and $\text{FDR} \approx \text{pFDR}$.

It can be shown that,

$$\Pr(H_0 \mid T > t_{\text{obs}}) < \Pr(H_0 \mid T = t_{\text{obs}}) \quad (6)$$

so that the evidence for H_0 given the exact ordinate is **always greater** than that corresponding to the tail area.

When one decides upon a value of FDR (or pFDR) to use in practice, the sample size should again be taken into account, since for large sample size one would not want to tolerate as large an FDR as with a small sample size.

Again, we would prefer a procedure that was consistent.

However, as in the single test situation, there is no prescription for deciding how FDR should decrease with increasing sample size.

Prostate cancer

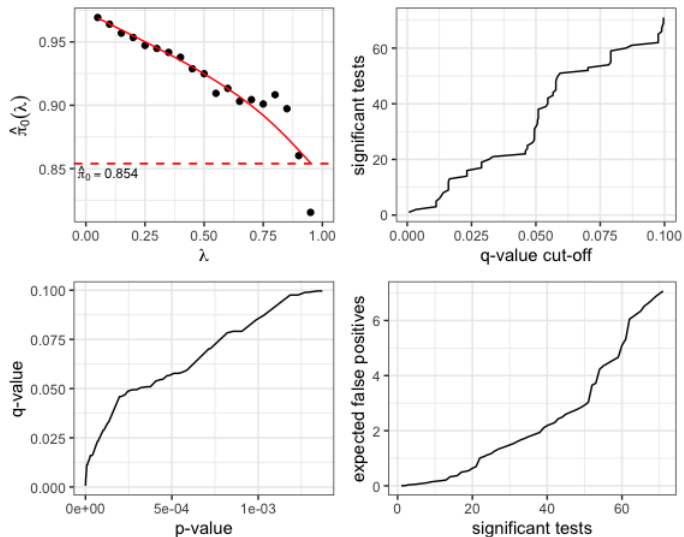


Figure 7: q-value plots for prostate cancer data.

Bayesian False Discoveries/Non-Discoveries

- ▶ In a Bayesian approach, based on Bayes factors we have a rule to flag a single association as **noteworthy** if:

$$\begin{aligned} \text{Posterior Odds} &= \text{Bayes Factor} \times \text{Prior Odds} \\ &< R \end{aligned}$$

where R is the ratio of costs of type II to type I errors.

- ▶ In a multiple testing situation in which m associations are being examined nothing, in principle, changes.
- ▶ We simply apply the same rule m times, perhaps changing the priors if we have different priors for different associations.
- ▶ The choice of threshold, R , and hence the procedure, does not depend on: **the number of tests being carried out**⁵.

⁵unless the prior on the null, or the ratio of costs of errors depends on the number of tests

Bayesian False Discoveries/Non-Discoveries

- ▶ As we have seen, the Bayes factor depends, crucially, on the **sample size**.
- ▶ In contrast, multiple testing based on p -values (e.g. Bonferroni/Sidak) does not depend on the sample size but, crucially, on the **number of tests m** .
- ▶ We have already noted that p -value calibration is very difficult, and we would like a procedure by which p -value thresholds **decrease to zero** with increasing sample size.
- ▶ The same would also be required of EFD or FDR based procedures.

To summarize in the case of normal test statistics:

The Bayesian decision is based on the Z score and on the sample size, n , but not on the number of tests, m .

In contrast:

The Bonferroni decision is based on the Z score and on the number of tests, m , but not on the sample size, n .

Bayesian Multiple Testing

In a Bayesian context, for a single test:

- ▶ If we call a hypothesis **noteworthy** then $\Pr(H_0 | \text{data})$ is the probability of a **false discovery**.
- ▶ If we call a hypothesis **not rejected** then $\Pr(H_1 | \text{data})$ is the probability of a **false non-discovery**.

A Key Point: A Bayesian analysis of a single SNP alone, or the same SNP from multiple SNPs will produce the same decision (assuming the prior is the same).

Bayesian False Discoveries/Non-Discoveries

In a multiple-hypothesis testing situation (and assuming ordered so the first K are rejected), we have

$$\text{Expected number of false discoveries} = \sum_{i=1}^K \Pr(H_{0i} | \text{data}_i)$$

$$\text{Proportion of false discoveries} = \frac{1}{K} \sum_{i=1}^K \Pr(H_{0i} | \text{data}_i)$$

$$\text{Expected number of false non-discoveries} = \sum_{i=K+1}^m \Pr(H_{1i} | \text{data}_i)$$

$$\text{Proportion of false non-discoveries} = \frac{1}{m - K} \sum_{i=K+1}^m \Pr(H_{1i} | \text{data}_i).$$

In the frequentist approaches to the expected FDR is (as usual) with respect to infinite hypothetical identical situations; the above Bayesian approach we have posterior summaries (so they are dependent on the model).

Empirical Bayes method

- ▶ Efron's **local FDR** (Efron *et al.*, 2001) uses a two-groups model to estimate the proportion of null/signal as a function of Z_i .

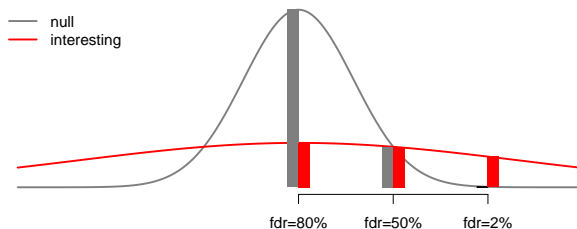


Figure 8: Local FDR.

- ▶ Estimating the 'null' component from the middle of the data, subtracting it from an overall density estimate, we can estimate local FDR, denoted $FDR(Z)$.

Local FDR Estimation

The **local FDR** corresponding to a test statistic z_0 is defined as

$$\text{FDR}(z_0) = \Pr(\text{gene } i \text{ is null} | z_i = z_0).$$

Note: not a tail area.

We have

$$\text{FDR}(z) = \frac{\pi_0 f_0(z)}{f(z)}.$$

In practice $f(z)$ is replaced by $\hat{f}(z)$, which is estimated via a Poisson model with log mean taken as a polynomial in z (so the z values are binned).

Prostate cancer example

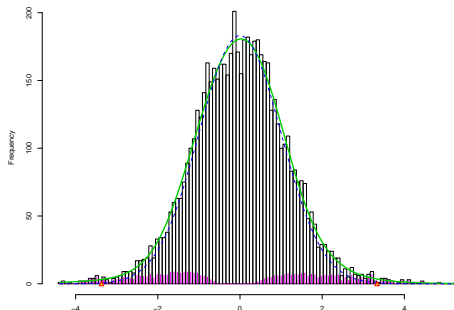


Figure 9: Local FDR for prostate cancer data. Blue curve is distribution if all null. the green solid line is the spline-based estimate of the mixture density f . Pink are non-null signals.

We find 25 genes with $\widehat{\text{FDR}}(Z_i) < 0.1$.

Multiple testing: Does Bayes help?

Efron's $\widehat{\text{FDR}}(z)$ is an 'empirical Bayes' method – it 'borrows strength' from the collection $z_i, i = 1, \dots, m$, to say what happens at specific z .

Hierarchical models also do this, using prior assumptions of **exchangeability** to motivate borrowing strength across subgroups.

As shown by Gelman *et al.* (2012)⁶, this is **not the same** as, for example, Bonferroni.

They also discuss **Type S** errors, which are sign errors, i.e., saying an association is positive when it is truly negative.

⁶In a paper entitled, 'Why We (Usually) Don't Have to Worry About Multiple Comparisons'

Multiple testing: Does Bayes help?

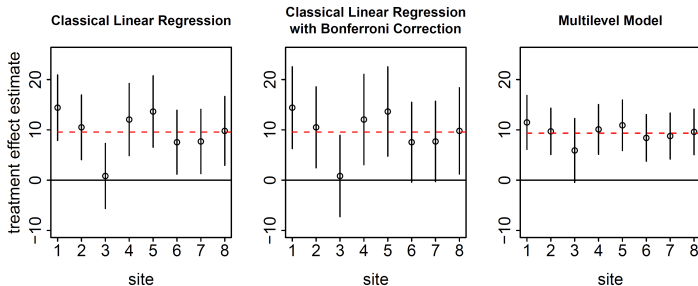


Figure 10: Point and 95% intervals, reproduction of Figure 1 from Gelman *et al.* (2012).

Compared to simpler methods, multilevel approaches do allow better inference on vectors of parameters – generally by trading some bias for reduced variance.

Bayes Mixture Model

We consider the mixture model described in Chapter 4 of Wakefield (2013).

The sampling model is $Y_i | \mu_i \sim N(\mu_i, \sigma_i^2)$, where the σ_i^2 are assumed known.

We specify a **mixture model** for the collection $[\mu_1, \dots, \mu_m]$, with

$$\mu_i = \begin{cases} 0 & \text{with probability } \pi_0 \\ \sim N(\delta, \tau^2) & \text{with probability } \pi_1 = 1 - \pi_0 \end{cases}$$

We use mixture component indicators $H_i = 0/1$ to denote the zero/normal membership model for transcript i .

Bayes Mixture Model

Collapsing over μ_i gives the three stage model:

Stage One:

$$Y_i \mid H_i, \delta, \tau, \pi_0 \sim_{ind} \begin{cases} \mathbf{N}(0, \sigma_i^2) & \text{if } H_i = 0 \\ \mathbf{N}(\delta, \sigma_i^2 + \tau^2) & \text{if } H_i = 1. \end{cases}$$

Stage Two: $H_i \mid \pi_1 \sim_{iid} \text{Bernoulli}(\pi_1)$, $i = 1, \dots, m$.

Stage Three: Independent priors on the common parameters:

$$p(\delta, \tau, \pi_0) = p(\delta)p(\tau)p(\pi_0).$$

Bayes Mixture Model

We illustrate the use of this model with

$$\begin{aligned}p(\delta) &\propto \mathbf{1}, \\p(\tau) &\propto \mathbf{1}/\tau \\p(\pi_0) &= \mathbf{1},\end{aligned}$$

so that we have improper priors for δ and τ^2 .

The latter choice still produces a proper posterior, because we have fixed variances at the first stage of the model.

Implementation is via a Markov chain Monte Carlo algorithm; Exercise 4.4 of Wakefield (2013) derives details of the algorithm.

Bayes Mixture Model

For transcript i , we may evaluate the posterior probabilities of the alternative

$$\begin{aligned}\Pr(H_i = 1 \mid y_i) &= \mathbb{E}[H_i \mid \mathbf{y}] \\ &= \mathbb{E}_{\delta, \tau^2, \pi_0 \mid \mathbf{y}} \left[\mathbb{E}(H_i \mid \mathbf{y}, \delta, \tau^2, \pi_0) \right] \\ &= \mathbb{E}_{\delta, \tau^2, \pi_0 \mid \mathbf{y}} \left[\Pr(H_i = 1 \mid \mathbf{y}, \delta, \tau^2, \pi_0) \right] \\ &= \mathbb{E}_{\delta, \tau^2, \pi_0 \mid \mathbf{y}} \left[\frac{\rho(\mathbf{y} \mid H_i = 1, \delta, \tau^2) \times \pi_1}{\rho(\mathbf{y} \mid H_i = 1, \delta, \tau^2) \times \pi_1 + \rho(\mathbf{y} \mid H_i = 0) \times \pi_0} \right] \tag{7}\end{aligned}$$

where

$$\begin{aligned}\rho(\mathbf{y} \mid H_i = 1, \delta, \tau^2, \pi_0) &= [2\pi(\sigma_i^2 + \tau^2)]^{-1/2} \exp \left[-\frac{(y_i - \delta)^2}{2(\sigma_i^2 + \tau^2)} \right] \\ \rho(\mathbf{y} \mid H_i = 0, \delta, \tau^2, \pi_0) &= [2\pi\sigma_i^2]^{-1/2} \exp \left[-\frac{y_i^2}{2\sigma_i^2} \right].\end{aligned}$$

Bayes Mixture Model

Expression (7) averages $\Pr(H_i = 1 \mid \mathbf{y}, \delta, \tau^2, \pi_0)$ with respect to the posterior $p(\delta, \tau^2, \pi_0 \mid \mathbf{y})$, and may be simply evaluated via

$$\frac{1}{T} \sum_{t=1}^T \frac{p(\mathbf{y} \mid H_i = 1, \delta^{(t)}, \tau^{2(t)})\pi_1^{(t)}}{p(\mathbf{y} \mid H_i = 1, \delta^{(t)}, \tau^{2(t)}, \pi_0^{(t)})\pi_1^{(t)} + p(\mathbf{y} \mid H_i = 0)\pi_0^{(t)}}$$

given samples $\delta^{(t)}, \tau^{2(t)}, \pi_0^{(t)}, t = 1, \dots, T$, from the Markov chain.

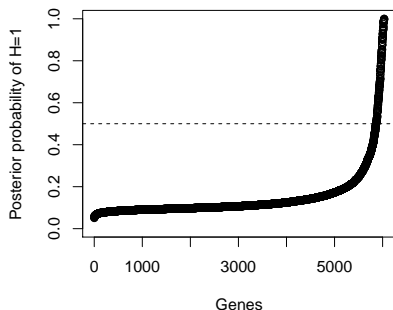


Figure 11: Posterior probability of alternative for prostate cancer.

Local false sign rates

Stephens (2017) has recently proposed an approach building on previous ideas.

The approach takes as input an estimate $\hat{\beta}_i$ and standard error s_i for the i -th signal and then (also) builds a hierarchical mixture model.

The posterior for β is

$$p(\beta_i | \hat{\beta}_i, s_i) \propto p(\hat{\beta}_i | \beta_i, s_i) \times p(\beta_i),$$

and the prior for β is assumed to be independent from a **unimodal** g , with

$$p(\beta_i) = \pi_0 \delta_0(\beta_i) + \sum_{k=1}^K \mathbf{N}(\beta_i | \mathbf{0}, \sigma_k^2),$$

where $\delta_0(\cdot)$ is a point mass at 0.

Local false sign rates

The approach centers on the **local false sign rate** LFSR_i which is the probability that we would make an error in the sign of effect i if we were forced to declare it either positive or negative (a Type S error):

$$\text{LFSR}_i = \min \left[\Pr(\beta_i \geq 0 | \hat{\pi}, \hat{\beta}, \mathbf{s}), \Pr(\beta_i \leq 0 | \hat{\pi}, \hat{\beta}, \mathbf{s}) \right].$$

Example: Suppose that

$$\Pr(\beta_i < 0 | \hat{\pi}, \hat{\beta}, \mathbf{s}) = 0.95$$

$$\Pr(\beta_i = 0 | \hat{\pi}, \hat{\beta}, \mathbf{s}) = 0.03$$

$$\Pr(\beta_i > 0 | \hat{\pi}, \hat{\beta}, \mathbf{s}) = 0.02$$

Then, $\text{LFSR}_i = \min(0.05, 0.98) = 0.05$.

For the prostate cancer data, a proportion 0.0161 of genes are associated with a non-zero effect

Conclusions

Conclusions

- ▶ Bayesian analysis is attractive in a multiple testing context, but the results are very **sensitive to the prior** on the proportion of nulls, π_0 .
- ▶ Fast methods are required for large m (e.g. in a GWAS context) of tests, which is still a drawback for many Bayesian approaches.
- ▶ **Priors** can be made a function of characteristics of the SNP (e.g. non-synonymous, previously implicated,...). See Johansson *et al.* (2012) for an example.
- ▶ Such priors can have a major impact on rankings and posterior probabilities.

What to do with multiple comparisons is a **difficult problem**:

- ▶ Apart from doing nothing, the only truly 'default' method is Bonferroni, which may not answer a relevant question, and/or may not answer it very well.
- ▶ Bonferroni is poorly-understood, as are other methods.
- ▶ If we use estimation (for example, via a hierarchical model) we can avoid multiple comparison problems (though care in the model specification needed).
- ▶ There are many summaries of techniques, see for example Efron and Hastie (2016).
- ▶ Stephens (2017) is a very good discussion of modern techniques.

References

- Benjamini, Y. and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society, Series B*, **57**, 289–300.
- Efron, B. and Hastie, T. (2016). *Computer Age Statistical Inference*. Cambridge University Press.
- Efron, B., Tibshirani, R., Storey, J., and Tusher, V. (2001). Empirical Bayes analysis of a microarray experiment. *Journal of the American Statistical Association*, **96**, 1151–1160.
- Gelman, A., Hill, J., and Yajima, M. (2012). Why we (usually) don't have to worry about multiple comparisons. *Journal of Research on Educational Effectiveness*, **5**, 189–211.
- Holm, S. (1979). A simple sequentially rejective multiple test procedure. *Scandinavian Journal of Statistics*, **6**, 65–70.
- Johansson, M., Roberts, A., Chen, D., Li, Y., Delahaye-Sourdeix, M., Aswani, N., Greenwood, M., Benhamou, S., Lagiou, P., Holcáová, I., Richiardi, L., Kjaerheim, K., Agudo, A., Castellsagué, X., Macfarlane, T., Barzan, L., Canova, C., Thakker, N. S., and A. Znaor, D. C., Healy, C., Ahrens, W., Zaridze, D., Szeszenia-Dabrowska, N., Lissowska, J., Fabiánová, E., Mates, I.,

- Bencko, V., Foretova, L., Janout, V., Curado, M., Koifman, S., Menezes, A., Wunsch-Filho, V., Eluf-Neto, J., Boffetta, P., Franceschi, S., Herrero, R., Garrote, L., Talamini, R., Boccia, S., Galan, P., Vatten, L., Thomson, P., Zelenika, D., Lathrop, M., Byrnes, G., Cunningham, H., Brennan, P., Wakefield, J., and Mckay, J. (2012). Using prior information from the medical literature in GWAS of oral cancer identifies novel susceptibility variant on chromosome 4 – the AdAPT method. *PLoS One*, **7**.
- 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.
- Pearson, E. (1953). Discussion of “Statistical inference” by D.V. Lindley. *Journal of the Royal Statistical Society, Series B*, **15**, 68–69.
- Sellke, T., Bayarri, M., and Berger, J. (2001). Calibration of p values for testing precise null hypotheses. *The American Statistician*, **55**, 62–71.
- Sidak, Z. (1967). Rectangular confidence regions for the means of multivariate normal distributions. *Journal of the American Statistical Association*, **62**, 626–633.
- Simes, R. J. (1986). An improved bonferroni procedure for multiple tests of significance. *Biometrika*, **73**, 751–754.

- Stephens, M. (2017). False discovery rates: a new deal. *Biostatistics*, **18**, 275–294.
- Storey, J. (2002). A direct approach to false discovery rates. *Journal of the Royal Statistical Society, Series B*, **64**, 479–498.
- Storey, J. (2003). The positive false discovery rate: A Bayesian interpretation and the q -value. *The Annals of Statistics*, **31**, 2013–2035.
- Storey, J. and Tibshirani, R. (2003). Statistical significance for genomewide studies. *Proceedings of the National Academy of Science*, **100**, 9440–9445.
- Wakefield, J. (2009). Bayes factors for genome-wide association studies: comparison with p -values. *Genetic Epidemiology*, **33**, 79–86.
- Wakefield, J. (2013). *Bayesian and Frequentist Regression Methods*. Springer, New York.
- Westfall, P., Johnson, W., and Utts, J. (1995). A Bayesian perspective on the bonferroni adjustment. *Biometrika*, **82**, 419–427.

Keyword Prior

- ▶ We now briefly describe a method for assigning priors to SNPs based on substantive information.
- ▶ In collaboration with scientists at IARC and at the Department of Computer Science at Sheffield University a method had been developed that searches through PubMed abstracts for pre-assigned keywords and key concepts.
- ▶ More details in Johansson *et al.* (2012).
- ▶ This information is used to assign prior probabilities of association with the phenotype for each SNP of interest.

Incorporating Prior Information

- ▶ Three prior groups were assigned, depending on the number of hits.
- ▶ The priors can subsequently be incorporated with the association results of GWAS using the previously described Bayesian framework.
- ▶ The method has acronym: **Adjusting Association Priors with Text (AdAPT)**.
- ▶ Details of the method can be found in Johansson *et al.* (2012).
- ▶ The **AdAPT** software is available here:
<http://services.gate.ac.uk/lld/gwas/service/config>

Incorporating Prior Information in a GWAS

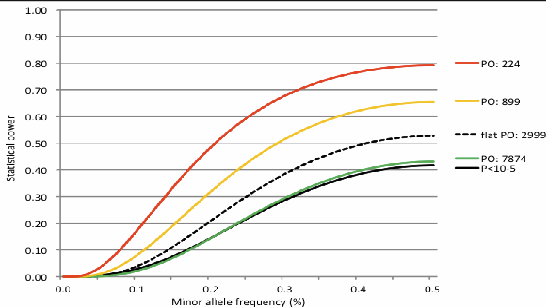
- ▶ SNPs are assigned to a group, based on the number of keywords that were found to be associated with this SNP.
- ▶ For the priors, keywords were ranked by priority: In the:
 - ▶ **1st group** G_1 : were words specific to lung cancer (eg, smoking, lung carcinoma).
 - ▶ **2nd group** G_2 : were more general words specifically relevant to lung cancer (smoking, nicotine, non-small cell carcinoma),
 - ▶ **3rd group** G_3 : were more general words (carcinogen, DNA damage).
- ▶ Each SNP was then placed in one of three prior categories:
 1. $C_1 = \{\text{not } G_1, \text{not } G_2, \text{not } G_3\}$.
 2. $C_2 = \{\text{at least one of } G_1, G_2, G_3 \text{ but not all}\}$.
 3. $C_3 = \{G_1, G_2, G_3\}$.

- ▶ We then assigned prior odds (PO) to $\Pr(H_0|C_j)/\Pr(H_1|C_j)$. Specifically for the three categories, the PO was set to 7874 (C_1), 899 (C_2) and 224 (C_3).
- ▶ These were used in the analysis to obtain the posterior odds on H_0 :

$$\frac{\Pr(H_0|y, C_j)}{\Pr(H_1|y, C_j)} = \frac{\Pr(y|H_0)}{\Pr(y|H_1)} \times \frac{\Pr(H_0|C_j)}{\Pr(H_1|C_j)}.$$

- ▶ First, the power was evaluated for the three categories, see Figure 13.
- ▶ The method was tested by comparing rankings of known **susceptibility alleles** in a previous lung cancer GWAS of 1989 cases and 2625 controls in 6 central European countries.
- ▶ The rankings of **6 SNPs** that have been independently replicated in multiple studies were calculated.

Figure 1 | Comparison of the statistical power when evaluating the noteworthiness of SNPs by BFDP and p -values.



These power calculations assume an evaluation of 300,000 SNPs, of which 300 are truly associated with the outcome and distributed evenly across three prior categories, respectively. The overall distribution of SNPs across the three prior categories is assumed to be (87.5%; 10%; 2.5%). First PO assumes one single prior category.

Figure 12: Power as a function of MAF, for three prior categories, for a single prior, and for a p -value approach.

Incorporating Prior Information: Proof of Principle Results

- ▶ The results below show that known susceptibility SNPs were ranked more highly by AdAPT BFDPs than by p-values.
- ▶ Rankings based on initial data with informative priors for the Bayes rankings:

SNP	p -value ranking	Bayes ranking
rs8034191	1	1
rs1051730	2	2
rs4324798	4	5
rs401681	73	30
rs2736100	76	32
rs3117582	121	34

Incorporating Prior Information: New Study Results

- ▶ Subsequently, the method was applied on a novel two phase GWAS of oral cancer, with 791 cases and 7,012 controls included in the discovery phase.
- ▶ A **Bayes threshold on the null** of 0.8 was assigned and 6 SNPs passed this test.
- ▶ One of these was already replicated, the replication was carried out for the remaining 5 AdAPT ranked SNPs in 1,046 cases and 2,131 controls from 4 case-control studies.
- ▶ **rs991316**, located in the ADH gene region of 4q23, displayed a statistically significant association with oral cancer risk in the replication phase (per-rare-allele log additive p -value $= 2.5 \times 10^{-3}$).
- ▶ This SNP was ranked 76th in the p -value list and so would not have been selected to carry forward, but was ranked 4th in the BFDP list.
- ▶ The combined **odds ratio** associated with having one additional rare allele was 0.84 (95% CI: 0.75–0.94).