

CHAPTER 13: HYPOTHESIS TESTING

In this chapter we discuss:

- Frequentist and Bayesian approaches to testing a single hypothesis.
- Multiple hypothesis testing, including variable selection.

Testing hypotheses is very context specific; we consider three distinct scenarios:

1. Confirmatory analyses in which an a priori hypothesis concerning a particular response/covariate relationship is of interest, and other variables have been measured and we wish to know which to adjust for.
2. Exploratory analyses where the aim is to gain clues as to structure in the data. For example, which covariates are causally related to a response, or characterizing sources of variability.
3. Prediction in which we are not concerned with causality, but merely with predicting a response given a set of variables.

385

Frequentist hypothesis testing

Suppose we are interested in the null hypothesis:

$$H_0 : \beta = 0$$

and we have a statistic T , with large values being increasingly unlikely under the null. The observed value is t_{obs} .

Possibilities for T :

- Squared Wald statistic from a regression analysis.
- Score statistic.

Under regularity conditions, $T \rightarrow_n \chi_1^2$ under the null, as $n \rightarrow \infty$. If n is not large then permutation/Monte Carlo can be used to give the distribution of the statistic under the null.

386

Interpretation of p -values

The $\Pr(T > t|H_0) = 1 - F(T)$ is uniform under the null. Let

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

denote the observed p -value — the probability of observing t_{obs} , **or a more extreme value** under the null.

Historically there has two approaches to the use of the above approach.

- Fisher and the pure test of significance. Quote the observed p -value as the measure of evidence against the null. No concept of rejecting the null in terms of the alternative, as there is no alternative!
No long-run frequency control of the type I error — $\alpha = \Pr(T > t_{\text{fix}})$, the probability of rejection of the null when it is true.
- Neyman-Pearson: specify an alternative hypothesis (with H_0 nested in H_1), and then evaluate the likelihood ratio statistic.
Use of the Neyman-Pearson lemma to find, for fixed α , the most powerful test. Decision rule: if $p < \alpha$ reject the null; under a **fixed** threshold this procedure controls the type I error.

387

Critique of Fisherian approach

- How do we decide on whether p is small or not, i.e. how to decide on a **threshold** for significance?
- With large sample sizes we will always obtain small p -values because although the p -value is still uniform under the null, we will be able to detect very subtle departures from the null with a large sample size:
 - To rectify this a confidence interval for β is often given along with the p -value, so that the scientific significance of the departure can be determined.
 - Suggests that the p -value threshold should decrease with increasing n .
- What is the prescription for finding a suitable T , particularly when we have nuisance parameters? Intuitively we want a statistic that takes large values (has high power) for realistic alternatives.
- Another aspect of p -values is that they use a tail area and

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

388

Critique of Neyman-Pearson approach

- How do we decide on a size α ?
- The emphasis is on fixing α , but this should also decrease with increasing n .
- Just stating that a p -value less than α was achieved is throwing away information — but if we give an observed p -value, how do we interpret.
- What if H_0 and H_1 are both wrong (always true in practice!).

Under both approaches, wouldn't we rather know $\Pr(H_0|\mathbf{y})$? This requires a Bayesian approach.

389

Bayesian decision theory approach to hypothesis testing

Under a Bayes decision theory approach to hypothesis testing the “decision” δ is taken that minimizes the posterior expected loss.

$L(\delta, H)$		Decision	
		$\delta = 0$	$\delta = 1$
Truth H	H_0	0	L_α
	H_1	L_β	0

Table 14: Losses corresponding to the decision δ , when the truth is H , L_α is the loss associated with a type I error and L_β with a type II error.

For hypothesis testing we have two possible data generating mechanisms:

$$\begin{aligned} H_0 &\rightarrow \beta_0 | H_0 \rightarrow \mathbf{y} | \beta_0 \\ H_1 &\rightarrow \beta_1 | H_1 \rightarrow \mathbf{y} | \beta_1 \end{aligned}$$

As always in the Bayesian approach in which unknowns are treated as random, the true hypothesis H is viewed as an unknown parameter for which the posterior may be derived.

We have the posterior probability of H_j :

$$\Pr(H_j | \mathbf{y}) = \frac{p(\mathbf{y} | H_j) \times \pi_j}{p(\mathbf{y})}$$

with π_j the prior probability of hypothesis H_j and

$$p(\mathbf{y} | H_j) = \int p(\mathbf{y} | \beta_j) f(\beta_j | H_j) d\beta_j$$

where $f(\beta_j | H_j)$ is the prior distribution over the parameters associated with hypothesis H_j , $j = 1, 2$.

391

With respect to Table 14, the posterior expected loss associated with the decision δ is

$$E[L(\delta, H)] = L(\delta, H_0) \Pr(H_0 | \mathbf{y}) + L(\delta, H_1) \Pr(H_1 | \mathbf{y})$$

so that for the two possible decisions (accept/reject H_0) the expected losses are:

$$\begin{aligned} E[L(\delta = 0, H)] &= 0 \times \Pr(H_0 | \mathbf{y}) + L_\beta \Pr(H_1 | \mathbf{y}) \\ E[L(\delta = 1, H)] &= L_\alpha \Pr(H_0 | \mathbf{y}) + 0 \times \Pr(H_1 | \mathbf{y}) \end{aligned}$$

To find the decision that minimizes posterior expected loss, we therefore need to take the smaller of:

$$\begin{aligned} E[L(\delta = 0, H)] &= L_\beta \times \Pr(H_1|\mathbf{y}) \\ E[L(\delta = 1, H)] &= L_\alpha \times \Pr(H_0|\mathbf{y}). \end{aligned}$$

We should choose $\delta = 1$ if

$$L_\beta \times \Pr(H_1|\mathbf{y}) \geq L_\alpha \Pr(H_0|\mathbf{y})$$

i.e. if

$$\frac{\Pr(H_1|\mathbf{y})}{1 - \Pr(H_1|\mathbf{y})} \geq \frac{L_\alpha}{L_\beta}$$

or

$$\Pr(H_1|\mathbf{y}) \geq \frac{L_\alpha/L_\beta}{1 + L_\alpha/L_\beta}.$$

Hence we only need to specify the ratio of losses.

If equal errors choose H_1 if $\Pr(H_1|\mathbf{y}) > \Pr(H_0|\mathbf{y})$.

393

Rearranging we can say that we should decide on H_1 if

$$\frac{\Pr(H_0|\mathbf{y})}{\Pr(H_1|\mathbf{y})} = \text{Bayes factor} \times \frac{\pi_0}{\pi_1} < \frac{L_\beta}{L_\alpha}$$

where

$$\text{Bayes factor} = \frac{p(\mathbf{y}|H_0)}{p(\mathbf{y}|H_1)} = \frac{\text{Posterior Odds}}{\text{Prior Odds}}.$$

So if a type I error is 4 times as bad as a type II error we should choose H_1 if the posterior odds on H_0 drop below 0.25.

394

Critique of Bayes Approach

- Need to specify prior distributions on all parameters under null and alternative, and on the hypotheses. In general cannot get away with improper priors when hypothesis testing is considered (unlike estimation).
- All of the calculations above should be conditioned on $H_0 \cup H_1$ — we are really obtaining the posterior probability of the null given one of the null or alternatives is true, and under our assumed data and prior models.
- The calculation of the Bayes factor requires integrals, which will usually be intractable.
- How to decide on the ratio of losses?

395

Calibrating p -values

Let $\text{data} = T > t_{\text{obs}}$ — we want $\Pr(H_0 | \text{data})$ but to obtain this we must specify alternatives – consider a simple alternative, say $H_1 : \beta = \beta_1$.

Then

$$\Pr(H_0 | \text{data}) = \frac{\Pr(\text{data} | H_0)\pi_0}{\Pr(\text{data} | H_0)\pi_0 + \Pr(\text{data} | H_1)\pi_1}$$

where $\pi_j = \Pr(H_j)$, $j=0, 1$. Dividing by $\Pr(H_1 | \text{data})$:

$$\begin{aligned} \text{Posterior Odds of } H_0 &= \frac{\Pr(\text{data} | H_0)}{\Pr(\text{data} | H_1)} \times \text{Prior Odds of } H_0 \\ &= \frac{p\text{-value}}{\text{power at } H_1} \times \text{Prior Odds of } H_0 \end{aligned}$$

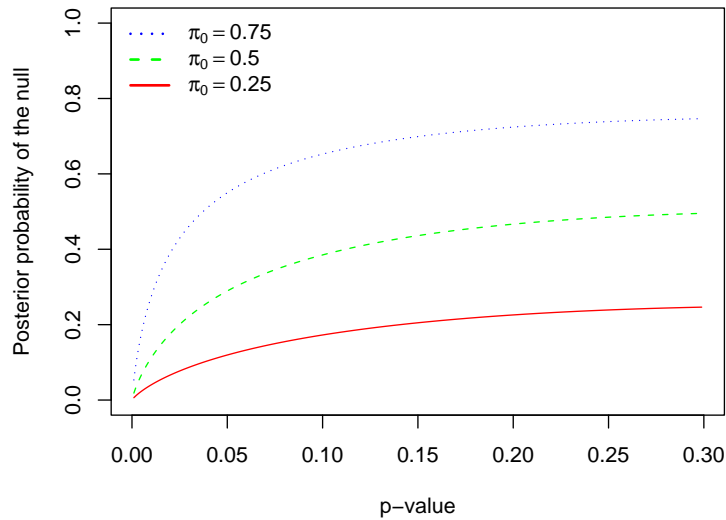
which depends on:

- The **prior** on H_0 , π_0 .
- The **power**, $\Pr(\text{data} | H_1)$ — greater power, more “evidence” in p -value.

396

Sellke, Bayarri and Berger (2001) show that for a p -value $p < 1/e$:

$$\Pr(H_0 | \text{data}) \geq \left\{ 1 - \frac{1}{e p \log p} \times \frac{\pi_1}{\pi_0} \right\}^{-1}$$



397

The interpretation of p -values

Sanity check! Why does anyone use p -values?

Historically it was usual to carry out single experiments and the prior on the alternative was not tiny.

With $\pi_0 = 0.5$:

- p -value = 0.05 gives $\Pr(H_0 | \text{data}) > 0.29$.
- p -value = 0.01 gives $\Pr(H_0 | \text{data}) > 0.11$.

But perhaps this is one of the reasons so many “findings” are not reproducible (along with confounding, multiple testing, errors-in-variables,...).

398

Variable Selection

We define the *null* model as that which contains an intercept only, the *minimal* model as the smallest model which is consistent with prior information.

So for example, in an epidemiological investigation we would almost always want to include terms for age and gender.

The minimal model may also be a function of the design so in matched case-control studies we include a term for each of the matching sets. Similarly in clinical trials in which treatments are randomized within a priori chosen strata, we again will include a term for strata.

399

Selection of Regressors

Trade-off: as we include more covariates, bias is reduced, but variability may be increased, dependent on how strong a predictor the covariate is (and its association with other covariates) — this is why we don't fit the full model.

We now describe some of the approaches to subset selection that have been proposed in the literature:

Forward selection. Begins with the simplest model. At each stage the 'best' unselected variable that satisfies the selection criterion is added. Best here is defined to be that variable whose deviance (or Wald or score statistic) is largest. This variable is added to the regression if its statistic is greater than a threshold of a specified significance level. This value, is contentious. Note that a maximum of p models will be considered in this procedure (out of 2^p).

Backward elimination. Begins with the full model. At each stage the covariate with the smallest deviance value that is less than a specified value is removed.

Stepwise regression (Efroymson's algorithm). Follows forward selection with backward elimination.

Difficulties

There are a number of problems with selection methods (Miller, 1990). In the conventional use of a hypothesis test, for the correct interpretation of significance levels the hypotheses must be specified before the data are examined. Note: same problem with transformations of y and/or x , choice of variance models and error distributions.

As a Bayesian all of the possible models that will be fitted should be assigned, priors with model averaging producing the summaries, averaging over all models.

Similarly for interval estimates to be valid the model must be specified *a priori*.

There now exists great potential for over-fitting in which models become too dataset-specific as they are refined on the basis of the examination of diagnostics.

In practice, if refinement is carried out through the fitting of alternative models (e.g. transformation of covariates, choice of distribution for the responses), then interval estimates will often be too narrow since they are produced by conditioning on the final model, and hence do not reflect the mechanism by which the model was selected.

401

We illustrate some of the difficulties of model selection with two simple examples.

Example 1: If we carry out a *single* hypothesis test and report the estimate of β_1 in a simple linear regression *only if* the null hypothesis of $\beta_1 = 0$ is rejected.

Figure 53 results – the bias is clear.

There are close links with publication bias in meta-analysis.

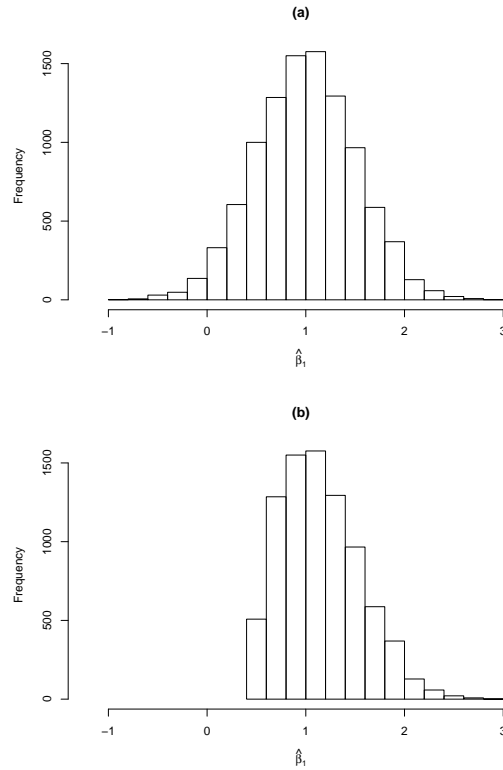


Figure 53: $E[\hat{\beta}_1] = 1.00$, while $E[\hat{\beta}_1 | \text{rejection of } H_0] = 1.27$.

403

Example 2: Suppose we are interested in β_1 but we wish to “control” for β_2 by testing whether the latter is significant.

In the following simulation, a multiple linear regression in X_1 and X_2 was carried out. The true values were $\beta_1 = \beta_2 = 1$ and X_1, X_2 were simulated from a bivariate normal with means zero, variances one, and correlation 0.7.

Figure 54 shows the results, in (a) we display the sampling distributions of $\hat{\beta}_1$ from the adjusted model. The mean and standard deviation of the distribution of $\hat{\beta}_1$ are 1.00 and 1.23.

Panel (b) displays the sampling distribution of the *reported* estimator.

The mean and standard deviation of the distribution of the reported estimate of β_1 are 1.23 and 1.01, respectively, showing positive bias and a reduced variance.

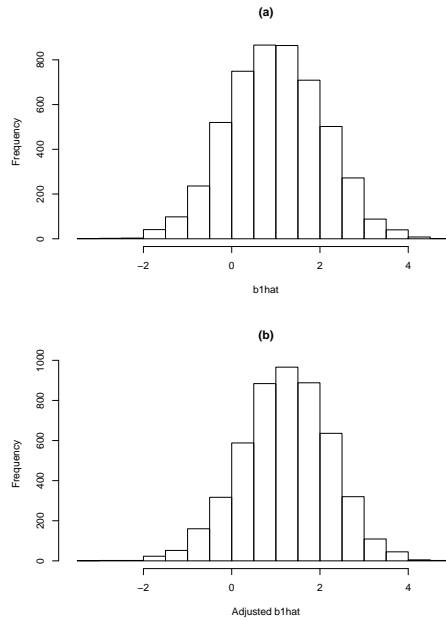


Figure 54: (a) Sampling distribution of $\widehat{\beta}_1$, (b) sampling distribution of $\widehat{\beta}_1$ given “control” for the possibility that $\beta_2 \neq 0$.

405

Frequentist model selection difficulties

From a frequentist standpoint estimators and test statistics should be examined via their long-run behaviour *given* the model-fitting process, including refinement. To be more explicit, let P denote the procedure by which a final model M is decided upon. Then suppose it is of interest to examine the bias of a statistic T ,

$$E[T|P] = E_{M|P}\{E[T|M]\}. \quad (50)$$

In general it will be incorrect to report $T | \widehat{M}$ where \widehat{M} is the final model chosen, since this does not reflect the procedure by which \widehat{M} was chosen, but rather acts as if the final model is the “truth”.

We know that

$$\text{var}(T|P) = E_{M|P}[\text{var}(T|M)] + \text{var}_{M|P}(E[T|M]).$$

but $\text{var}(T|\widehat{M})$ is reported (which approximates the first term only).

Under a frequentist approach inference follows from the behaviour of an estimator under repeated sampling.

406

Bayesian model selection difficulties

From a Bayesian standpoint the same problem of dredging exists because the posterior distribution should reflect all sources of uncertainty and *a priori* all possible models that may be entertained should be explicitly stated, with prior distributions being placed upon different likelihoods and the parameters of these likelihoods.

Model averaging (see later) should then be carried out across the different possibilities, a process which is fraught with difficulties not least in placing “comparable” priors over what may be fundamentally different objects.

(One solution is to place prior on “model-free” quantities.)

407

Model Averaging

Suppose the action concerns a parameter of interest T (which for simplicity we assume is univariate) that is well-defined for all models.

We have

$$E[T|\mathbf{y}] = \sum_{j=1}^J E[T|\mathbf{y}, M_j] \times \Pr(M_j|\mathbf{y}),$$

and

$$\begin{aligned} \text{var}(T|\mathbf{y}) &= \sum_{j=1}^J \text{var}(T|\mathbf{y}, M_j) \times \Pr(M_j|\mathbf{y}) \\ &+ \sum_{j=1}^J \{E[T|\mathbf{y}, M_j] - E[T|\mathbf{y}]\}^2 \times \Pr(M_j|\mathbf{y}). \end{aligned}$$

This latter term shows how not only parameter uncertainty but *model* uncertainty is accounted for.

- Specification of priors is not trivial.
- Interpretation.
- Continuous model expansion.

408

A possible compromise

One solution is to never refine the model for a given data set. This approach is operationally pure but pragmatically dubious (unless one is in the context of a randomized experiment) since we may obtain appropriate inference for a model that is a very poor description of the phenomenon under study.

The philosophy suggested here is to think as carefully as possible about the initial model class before the analysis proceeds, but after fitting to carry out model checking and refine the model in the face of *clear* model misspecification, with refinement ideally being carried out within distinct *a priori* known classes.

With reference to (50), if a model is chosen because it is clearly superior to the alternatives, then it may be reasonable to assume that $E[T | P] \approx E[T | \widehat{M}]$, because \widehat{M} would be consistently chosen in repeated sampling under these circumstances.

So, for example, examining quantile-quantile plots for different t distributions and picking the one that produces the straightest line would not be a good idea.

Inference then proceeds as if the final model were the one that were chosen initially. This is clearly a subjective procedure but can be informally justified via either philosophical approaches.

409

In a similar vein, under a Bayesian approach the above procedure is consistent with model-averaging but with the posterior model weight being concentrated upon the chosen model (since alternative models are only rejected on the basis of clear inadequacy).

The aim is to provide probability statements, from either philosophical standpoints that are “honest” representations of uncertainty. The above approach is relevant to analyses that are more confirmatory in their outlook, as opposed to being used for prediction, or for more exploratory purposes (for example, to gain clues to models that may be appropriate for future data analyses).

Conclusions

- For Confirmatory studies – try to avoid any model selection. Use background context to specify model.
- Exploratory studies – stepwise and all subsets may point to important variables, but attaching a p-value is difficult. Model averaging is another possibility.
- Prediction – some form of shrinkage should be used. Cross-validation may be used to choose a model. Bayesian model averaging also useful.

411

Lindley's Paradox

Lindley (1957) gave a much-quoted example in which Bayesian and frequentist inference differ.

We assume that $\bar{Y}|\theta \sim N(\theta, \sigma^2/n)$ with σ^2 known and θ unknown.

We first state Lindley's paradox, without loss of generality suppose the null is $H_0 : \theta = 0$, with alternative $H_1 : \theta \neq 0$.

Let

$$\bar{y}_n = Z_p \sigma / \sqrt{n}$$

where p is the p -value and Z_p is the corresponding quantile of the normal distribution, i.e. that Z_p such that $\Pr(Z > Z_p) = p/2 = \Pr(Z < -Z_p)$.

We define \bar{y}_n in this way so that as n increases the p -value remains constant.

For a Bayesian analysis assume that we assign π_0 to H_0 , and under the alternative $\theta \sim N(0, \tau^2)$. Then we have

$$\Pr(H_0|\bar{y}_n) = \frac{\text{BF} \times \text{PO}}{1 + \text{BF} \times \text{PO}}$$

where the Bayes factor is given by

$$\text{BF} = \frac{p(\bar{y}_n|H_0)}{p(\bar{y}_n|H_1)}$$

and $\text{PO} = \pi_0/(1 - \pi_0)$ is the prior odds.

We have

$$\begin{aligned}\bar{y}_n|H_0 &\sim N(0, \sigma^2/n) \\ \bar{y}_n|H_1 &\sim N(0, \sigma^2/n + \tau^2)\end{aligned}$$

so that

$$\text{BF} = \frac{(2\pi\sigma^2/n)^{-1/2} \exp\left[-\frac{\bar{y}_n^2}{2\sigma^2/n}\right]}{(2\pi[\sigma^2/n + \tau^2])^{-1/2} \exp\left[-\frac{\bar{y}_n^2}{2(\sigma^2/n + \tau^2)}\right]}$$

413

For $\sigma^2 = 1, \tau^2 = 0.2^2, \pi_0 = 0.5, p = 0.05$ and $n = 1, \dots, 100,000$, Figure 55 shows the posterior probability of the null as a function of n .

From the starting position of $\Pr(H_0|\bar{y}_n) = 0.5$ the curve initially falls, reaching a minimum at around $n = 100$, and then increases towards 1, illustrating the “paradox”.

For large values of n , \bar{y}_n is very close to the null value of 0, but there is high power to detect any difference from 0, and so a p -value of 0.05 is not difficult to achieve.

The Bayes factor calculates the density under the alternative also, and values close to 0 are more likely under the null, Figure 56 illustrates for $n = 1000$, the green vertical line denotes \bar{y}_n .

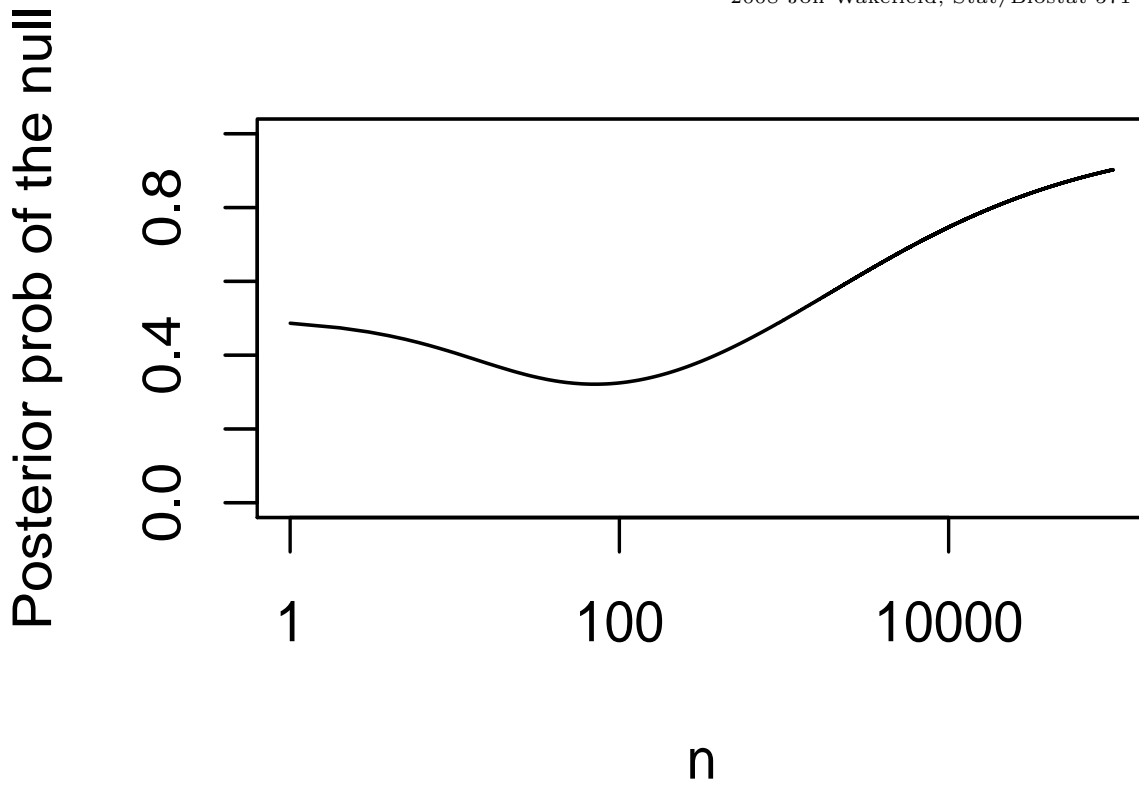


Figure 55: Posterior probability of the null for a fixed p -value of 0.05.

415

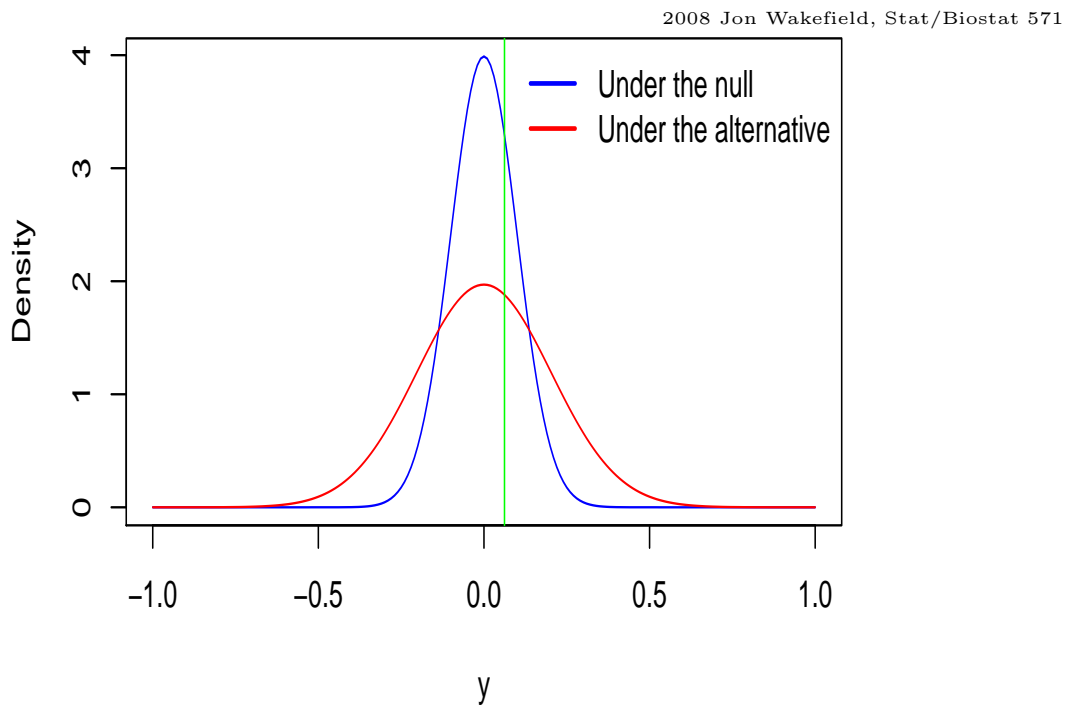


Figure 56: Numerator (blue) and denominator (red) of the Bayes factor for $n = 1000$. The green line represents \bar{y}_n for this n .

416

Bayesian analysis with a tail-area

We now consider a Bayesian analysis in which the data appear in the form of knowing that $|\bar{Y}_n| \geq \bar{y}_n$, a censored observation. This is clearly not realistic since a Bayesian would condition on the *actual* value observed, but it does allow considerations of statements about power made by several authors.

We have

$$\text{BF} = \frac{\Pr(|\bar{Y}_n| \geq \bar{y}_n | H_0)}{\Pr(|\bar{Y}_n| \geq \bar{y}_n | H_1)}$$

so that the numerator is the p -value, and the denominator is given by

$$\begin{aligned} \Pr(|\bar{Y}_n| \geq \bar{y}_n | H_1) &= \int \Pr(|\bar{Y}_n| \geq \bar{y}_n | \theta) \pi(\theta) d\theta \\ &= \int \{ \Pr(\bar{Y}_n \geq \bar{y}_n | \theta) + \Pr(\bar{Y}_n \leq -\bar{y}_n | \theta) \} \pi(\theta) d\theta \\ &= \int \left\{ \Phi \left(Z \geq \frac{\sqrt{n}(\bar{y} - \theta)}{\sigma} \right) + \Phi \left(Z \leq \frac{\sqrt{n}(-\bar{y} - \theta)}{\sigma} \right) \right\} \pi(\theta) d\theta \end{aligned}$$

where $Z \sim N(0, 1)$, so that we are evaluating the average of the power over the prior $\pi(\theta)$.

We emphasize that it is the post-data power that is being evaluated, i.e. it is based on the observed value of the statistic.

417

Figure 57 gives the average power as a function of n , and we see a monotonic increase with sample size towards the value 1.

Hence, as seen in Figure 58 the Bayes factor based on the tail-area information is monotonic decreasing towards the p -value as n increases (with $\pi_0 = 0.5$ this gives the posterior probability of the null also). Hence this justifies the claims of a number of authors that greater credence should be given to p -values based on large sample sizes/power.

The difference in behavior between a genuine Bayesian analysis that conditions on the actual statistic and that based on the tail area is apparent.

As noted by Lindley (1957, p. 189–190), “...the paradox arises because the significance level argument is based on the area under a curve and the Bayesian argument is based on the ordinate of the curve”.

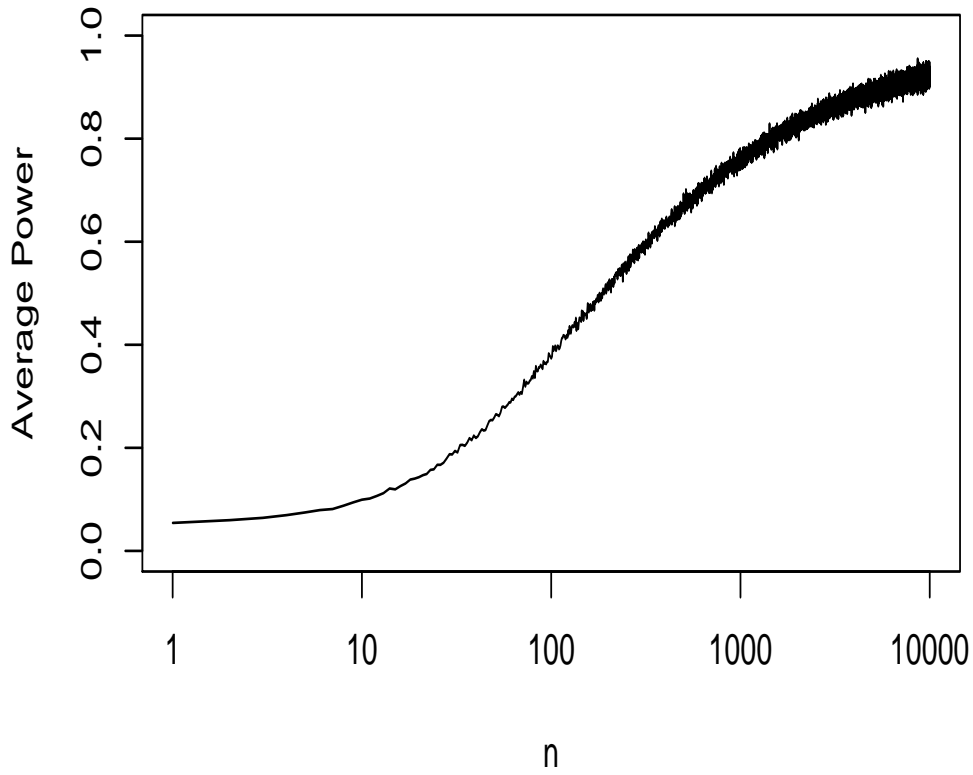


Figure 57: Average power under the $N(0, 0.2^2)$ prior for a fixed p -value of 0.05.

419

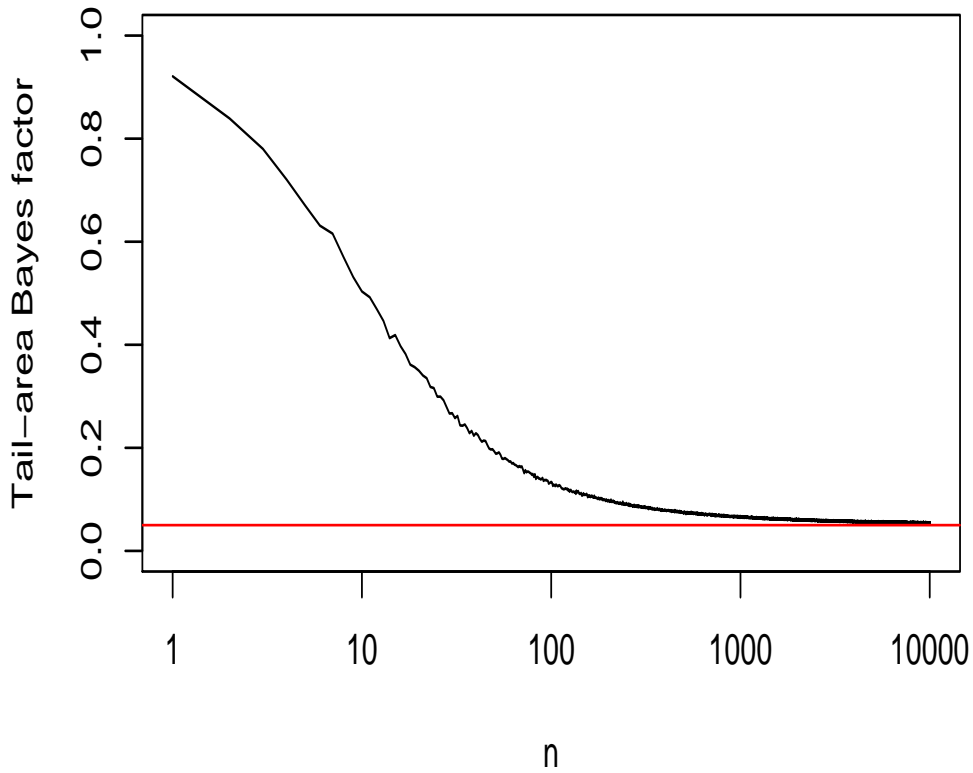


Figure 58: Bayes factor based on the observed tail area.

420

Multiple Hypothesis Testing

In subgroup analyses we want to examine the data for effects in different sub-groups.

Possibilities when m tests are performed and 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**.

How do we select a rule that will determine K ?

421

The Statistical Set-Up

The p -value was introduced in a single test situation to control the Type I error rate – historically the emphasis was placed on avoiding a Type I error.

The **family-wise error rate** (FWER) is the probability of making **at least** one Type I error, i.e. $\Pr(B \geq 1 | \text{all } H_0 \text{ true})$.

More recently there has been more interest in the **false discovery rate** — the expected proportion of rejected nulls that are actually true, see Benjamini and Hochberg (1995) and Storey (2002).

422

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 FWER is given by:

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

where α^* is the level for each test.

Bonferroni takes $\alpha^* = \alpha/m$ to give $\text{FWER} \leq \alpha$.

For control at $\alpha = 0.05$ with $m = 10$ tests take $\alpha^* = 0.05/10 = 0.005$.

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

423

Bayesian Bonferroni

If we have prior probability of the null $\pi_{0i} = \pi_0$ for $i = 1, \dots, m$ then

$$\Pi_0 = \Pr(\text{prior probability that all nulls are true}) = \pi_0^m$$

For example, if $\pi_0 = 0.5$ and $m = 10$, $\Pi_0 = 0.00098$, which may be deemed too small.

Westfall et al. (1997) show that if we take $\pi_{0i} = \Pi_0^{1/m}$ (so that the probability that all nulls are true is Π_0) — call this prior 2 — then for independent tests we have approximately

$$\alpha_B = \Pr(H_{0i} | \mathbf{y}_i, \text{prior 2}) \approx m \times \Pr(H_{0i} | \mathbf{y}_i, \text{prior 1}) = m \times \alpha_B^*$$

where prior 1 is $\pi_{0i} = \Pi_0$.

CHAPTER 14: MISSING DATA

A serious problem in data analysis is the existence of missing data. We concentrate on missing responses in a dependent data situation.

Implications of missing data:

1. Data are unbalanced – not a problem given modern regression techniques.
2. Information loss.
3. Depending on the mechanism of missingness, bias in estimation may result.

Missing data can arise in numerous ways, and understanding the mechanism is crucial to appropriate modeling assumptions.

In a longitudinal study, if *drop-out* occurs at a certain time then no additional data are observed after that point.

425

Examples:

1. In a health-air pollution study an individual may be unavailable for measurement because he/she took a job in another area.
2. In a clinical trial, patients may be removed from the study if their longitudinal measurements are below/above some limit.
3. Censoring – measurement instruments may be inaccurate below a lower limit of detection, this limit is then reported.
4. The value of the outcome may itself determine the missingness, but the outcome is unobserved.

In 1, the missingness will not be a problem unless the person moved area because of health problems. In 2, the missingness will be a function of the responses on previous occasions, while in 3 and 4 it depends on the actual measurement that would have been recorded.

426

Example: Simulated Data

Data were simulated in which the data ($m = 200, n_i = 10, i = 1, \dots, m$) were generated from a linear mixed model in which intercepts and slopes are random (and independent), with measurement error and $\beta_0 = 100, \beta_1 = -5$.

Figure 59 shows the resultant data.

We then simulated drop-out by a mechanism in which if the outcome falls below 65, the subsequent observations are lost (but we retain the initial one below 65).

Figure 60 shows the resultant data (509 data points were lost).

427

2008 Jon Wakefield, Stat/Biostat 571

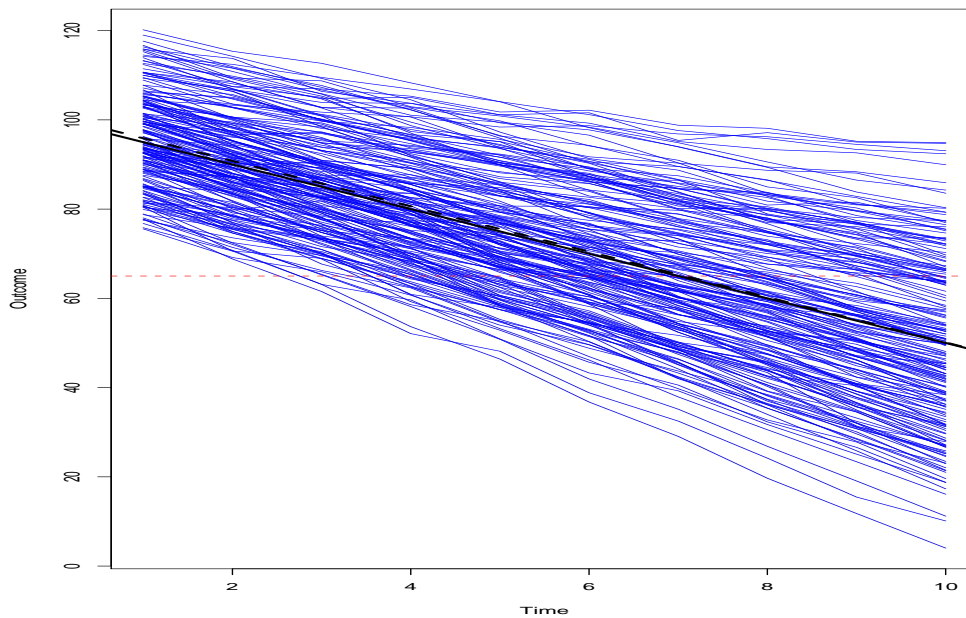


Figure 59: Full simulated data set: solid line is truth and dashed the LS line.

428

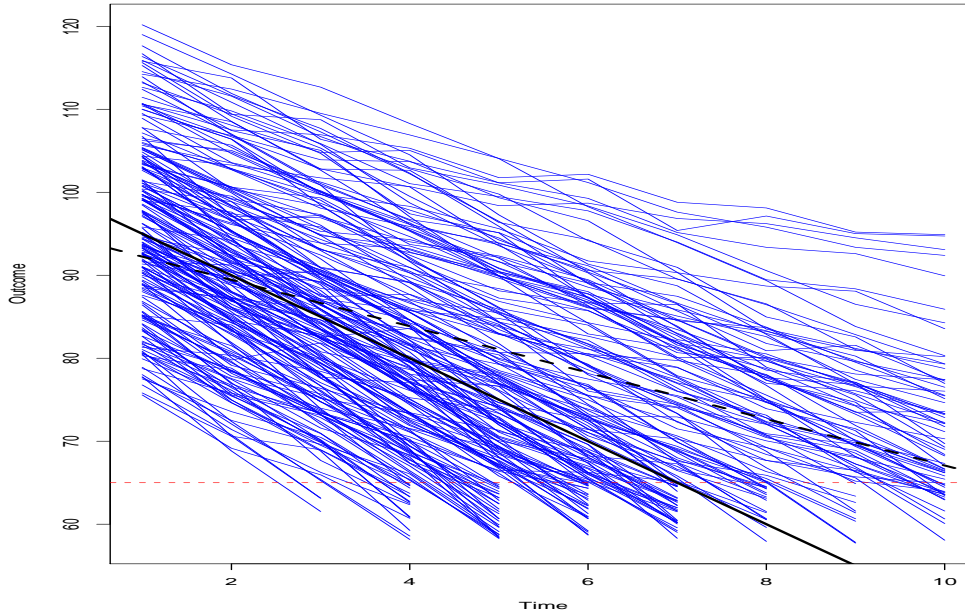


Figure 60: Simulated data set with drop-out: solid line is truth and dashed the LS line.

429

Mechanisms of Missingness

The impact of missing data depends crucially on the mechanism of missingness, that is the probability model for missingness.

We let \mathbf{R}_i be a vector of response indicators for the i -th units so that

$$R_{ij} = \begin{cases} 1 & \text{if } Y_{ij} \text{ is observed} \\ 0 & \text{if } Y_{ij} \text{ is missing} \end{cases}$$

We partition the complete data vector $\mathbf{Y}_i = (\mathbf{Y}_i^O, \mathbf{Y}_i^M)$ into those components that are observed, \mathbf{Y}_i^O , and those that are missing \mathbf{Y}_i^M .

There are two ways of factoring the data:

$$\begin{aligned} p(\mathbf{Y}, \mathbf{R} \mid \mathbf{x}) &= p(\mathbf{Y} \mid \mathbf{x}) \times p(\mathbf{R} \mid \mathbf{Y}, \mathbf{x}) \\ p(\mathbf{Y}, \mathbf{R} \mid \mathbf{x}) &= p(\mathbf{Y} \mid \mathbf{R}, \mathbf{x}) \times p(\mathbf{R} \mid \mathbf{x}) \end{aligned}$$

The first is known as a **selection model** (individuals are selected according to their outcome), and the second as a **pattern mixture model** (we “mix” pattern specific models). We concentrate on the former.

Three situations are distinguished:

1. Missing completely at random (MCAR).
2. Missing at random (MAR).
3. Not missing at random (NMAR).

each of which we now discuss in detail.

Unfortunately the terminology is confusing!

431

Missing Completely at Random (MCAR)

Data are MCAR if

$$\Pr(R_{ij} = 1 \mid \mathbf{Y}^O, \mathbf{Y}^M, \mathbf{x}) = \Pr(R_{ij} \mid \mathbf{x}),$$

so that the missingness does not depend on the response data, observed or unobserved.

This implies that

$$E[Y_{ij} \mid R_{ij} = 1, \mathbf{X}_i] = E[Y_{ij} \mid \mathbf{X}_i]$$

No selection bias.

Missing at Random (MAR)

Data are MCAR if

$$\Pr(R_{ij} = 1 \mid \mathbf{Y}^O, \mathbf{Y}^M, \mathbf{x}) = \Pr(R_{ij} \mid \mathbf{Y}^O, \mathbf{x}),$$

so that the missingness may depend on observed values.

This implies that

$$E[Y_{ij} \mid R_{ij} = 1, \mathbf{X}_i] \neq E[Y_{ij} \mid \mathbf{X}_i]$$

which suggests that the GEE approach might be in trouble in terms of biased parameter estimates.

432

Not Missing at Random (NMAR)

If the missingness depends on \mathbf{Y}^M , i.e.

$$\Pr(R_{ij} = 1 \mid \mathbf{Y}^O, \mathbf{Y}^M, \mathbf{x}) = \Pr(R_{ij} \mid \mathbf{Y}^O, \mathbf{Y}^M, \mathbf{x}).$$

In this case the mechanism is also sometimes referred to as **non-ignorable**.

This selection bias is not fixable, since we don't know the outcomes that caused the problems. Models can be postulated, but are not checkable from the observed data alone.

In general it is obviously best if we know why the data are missing.

433

Approaches**Complete-case analysis**

A simple approach is to exclude units that did not provide data at all intended occasions. Clearly there is a loss of information in this process, and bias will result unless the data are MCAR. Not to be recommended.

Available-case analysis

This approach uses the largest set of available data for estimating parameters. Will provide biased estimates unless the data are MCAR.

Last observation carried forward

In a longitudinal setting we could simply “fill-in” the missing values, extrapolating from the last observed value. As a general method not to be recommended.

Imputation

An appealing approach is to “fill-in”, or impute, the missing values and then carry out a conventional analysis. Complex models for the missingness can be incorporated (closely related to *data augmentation* which we describe later).

434

Likelihood-based approach

Let $\boldsymbol{\theta}$ be the parameters of the model for \mathbf{Y} , and $\boldsymbol{\phi}$ the parameters for \mathbf{R} .

In general, a natural way to decompose the data as

$$\begin{aligned} p(\mathbf{Y}^O, \mathbf{Y}^M, \mathbf{R} \mid \boldsymbol{\theta}, \boldsymbol{\phi}) &= p(\mathbf{Y}^O, \mathbf{Y}^M \mid \boldsymbol{\theta}, \boldsymbol{\phi}) \times \Pr(\mathbf{R} \mid \mathbf{Y}^O, \mathbf{Y}^M, \boldsymbol{\theta}, \boldsymbol{\phi}) \\ &= p(\mathbf{Y}^O, \mathbf{Y}^M \mid \boldsymbol{\theta}) \times \Pr(\mathbf{R} \mid \mathbf{Y}^O, \mathbf{Y}^M, \boldsymbol{\phi}) \end{aligned}$$

where we have also assumed that the data and missingness models have distinct parameters.

We require a distribution for the observed data, \mathbf{Y}^O, \mathbf{R} :

$$p(\mathbf{Y}^O, \mathbf{R} \mid \boldsymbol{\theta}, \boldsymbol{\phi}) = \int p(\mathbf{Y}^O, \mathbf{Y}^M \mid \boldsymbol{\theta}) \times \Pr(\mathbf{R} \mid \mathbf{Y}^O, \mathbf{Y}^M, \boldsymbol{\phi}) d\mathbf{Y}^M.$$

This is an example of a [selection model](#).

435

Suppose we are in the MAR situation so that

$$\Pr(\mathbf{R} \mid \mathbf{Y}^O, \mathbf{Y}^M, \boldsymbol{\phi}) = \Pr(\mathbf{R} \mid \mathbf{Y}^O, \boldsymbol{\phi}).$$

In this situation the likelihood is given by

$$\begin{aligned} p(\mathbf{Y}^O, \mathbf{R} \mid \boldsymbol{\theta}, \boldsymbol{\phi}) &= \int p(\mathbf{Y}^O, \mathbf{Y}^M \mid \boldsymbol{\theta}) d\mathbf{Y}^M \times \Pr(\mathbf{R} \mid \mathbf{Y}^O, \boldsymbol{\phi}) \\ &= p(\mathbf{Y}^O \mid \boldsymbol{\theta}) \times \Pr(\mathbf{R} \mid \mathbf{Y}^O, \boldsymbol{\phi}) \end{aligned}$$

Hence we have the log-likelihood

$$\log p(\mathbf{Y}^O \mid \boldsymbol{\theta}) + \log \Pr(\mathbf{R} \mid \mathbf{Y}^O, \boldsymbol{\phi})$$

and can ignore the second term and don't have to model the missingness mechanism.

Important Point: We need to get the model right!!!

436

Simulation Study

Model	$\widehat{\beta}_1$	(s.e.)	$\widehat{\beta}_2$	(s.e.)
GEE ind	101.0	0.715	-5.084	0.109
GEE exch	101.0	0.715	-5.084	0.109
LMEM1	101.0	0.981	-5.084	0.037
LMEM2	101.0	0.720	-5.084	0.109
GEEiDO ind	95.0	0.894	-2.796	0.134
GEEDO exch	98.8	0.787	-4.304	0.114
LMEM1DO	98.8	0.837	-4.282	0.041
LMEM2DO	100.1	0.722	-5.097	0.112

Table 15: Results of GEE and LMEM analyses of full and drop-out simulated data, LMEM1 is random intercepts only, LMEM2 is random intercepts and slopes.

- Bias for GEE is bad (particularly working independence).
- Bias for LMEM if we only assume random intercepts (terrible se on $\widehat{\beta}_2$).
- LMEM with random intercepts and slopes recovers the truth.

437

Models for Drop-out

If the missingness is monotone, in the sense that if $R_{ij} = 0$ then $R_{ik} = 0$ for all $k > j$, then we define the drop-out time as

$$D_i = \min_k \{R_{ik} = 0\}.$$

Hence $2 \leq D_i \leq n_i + 1$, with $D_i = n_i + 1$ for an individual that does not drop out.

The reason for drop-out may be that the individual was not responding well, and their outcomes reflect this.

To examine this possibility we could fit logistic models of the form:

$$\log \left(\frac{\Pr(D_i = k | D_i \geq k, Y_{i1}, \dots, Y_{ik})}{\Pr(D_i > k | D_i \geq k, Y_{i1}, \dots, Y_{ik})} \right) = \phi_0 + \phi_1 Y_{ik-1}$$

and look for evidence that $\phi_1 \neq 0$.

Bayesian Inference via Data Augmentation

Data augmentation is a auxiliary variable method that treats the missing observations as unknown parameters – this can lead to simple MCMC schemes.

General formulation: we have posterior

$$\begin{aligned} p(\boldsymbol{\theta}, \mathbf{Y}^M \mid \mathbf{Y}^O) &= p(\boldsymbol{\theta} \mid \mathbf{Y}^M, \mathbf{Y}^O)p(\mathbf{Y}^M \mid \mathbf{Y}^O) \\ &= p(\mathbf{Y}^M \mid \boldsymbol{\theta}, \mathbf{Y}^O)p(\boldsymbol{\theta} \mid \mathbf{Y}^O) \end{aligned}$$

MCMC scheme:

1. Auxiliary variables:

$$\mathbf{Y}^M \sim p(\mathbf{Y}^M \mid \mathbf{Y}^O, \boldsymbol{\theta}).$$

2. Model parameters:

$$\boldsymbol{\theta} \sim p(\boldsymbol{\theta} \mid \mathbf{Y}^O, \mathbf{Y}^M).$$

The auxiliary variable scheme may be modified to $p(\mathbf{Y}^M \mid \mathbf{Y}^O, \boldsymbol{\theta}) \sim p(\mathbf{Y}^M \mid \boldsymbol{\theta})$, depending on the missing data model, as we now illustrate.

439

Example: Censoring Model

Suppose we have data Y_i measured at times t_i , $j = 1, \dots, n$, but measurements *below the lower limit of detection*, D (assumed known) are not recorded. Also suppose that the data generating model (likelihood) is:

$$Y \mid \boldsymbol{\beta}, \sigma \sim_{ind} N(\eta(\boldsymbol{\beta}, t), \sigma^2).$$

Clearly setting such measurements to zero or ignoring the measurements will lead to bias in estimation.

Figure 61 illustrates for a set of simulated data in which the true slope was -0.01 ; the slope estimates are -0.0099 , -0.0095 and -0.0087 for the full data, set equal to D and ignored schemes, respectively.

440

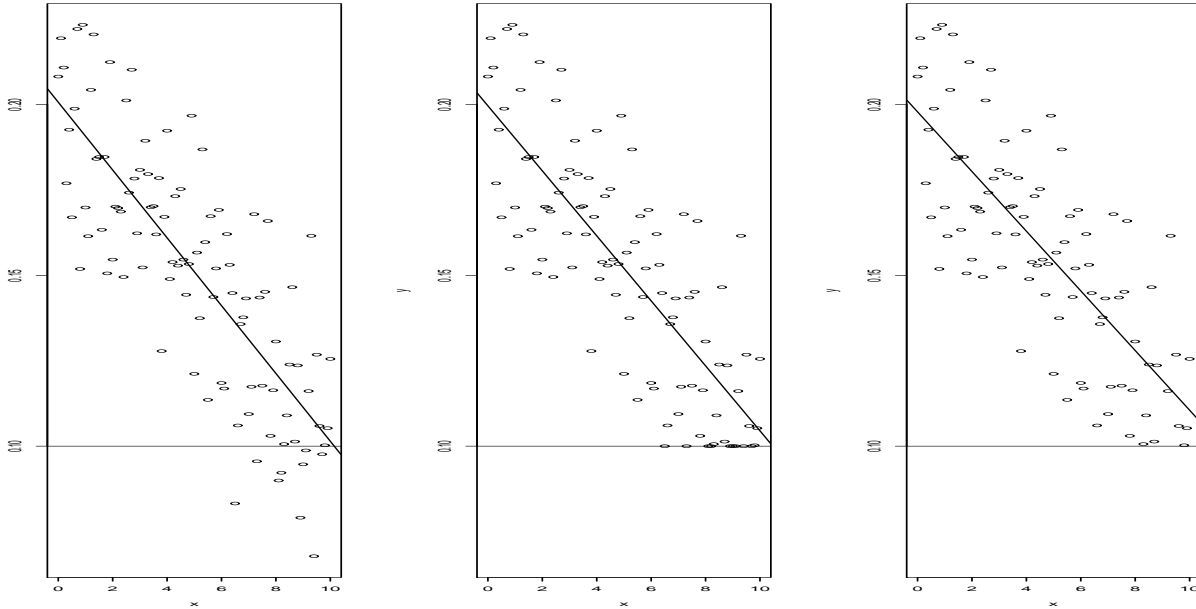


Figure 61: All data (left), assigned to lower limit (middle), ignored (right). Horizontal line is the lower limit of detection.

441

Suppose that the last c measurements are censored, the remaining $n - c$ being uncensored. Then

$$\begin{aligned}
 p(\mathbf{y} | \theta) &= \prod_{i=1}^{n-c} p(y_i | \boldsymbol{\beta}, \sigma^2) \prod_{i=c+1}^n \Pr(Y_i < D | \boldsymbol{\beta}, \sigma^2) \\
 &= \prod_{i=1}^{n-c} \phi\left(\frac{y_i - \eta(\boldsymbol{\beta}, t_i)}{\sigma}\right) \prod_{i=c+1}^n \Phi\left(\frac{D - \eta(\boldsymbol{\beta}, t_i)}{\sigma}\right)
 \end{aligned}$$

where

$$\phi(z) = (2\pi)^{-1/2} \exp(-z^2/2)$$

and

$$\Phi(z_0) = \Pr(Z < z_0) = \int_{-\infty}^{z_0} \phi(z) dz$$

where Z is an $N(0, 1)$ random variable.

To perform likelihood or Bayesian inference we need to numerically evaluate the distribution function of a normal distribution for each likelihood calculation.

442

Data Augmentation Scheme

Letting $\mathbf{Y}^O = \{Y_i, i = 1, \dots, n - c\}$ and $\mathbf{Y}^M = \{Y_i, i = n - c + 1, \dots, n\}$, we iterate between

1. $y_i \mid \boldsymbol{\beta}, \sigma \sim \text{TruncNorm}(\eta(\boldsymbol{\beta}, t_i), \sigma^2)$, on $(-\infty, D)$, $i = n - c + 1, \dots, n$.
2. $\boldsymbol{\beta} \mid y_1, \dots, y_n, \sigma^2 \propto \prod_{i=1}^n p(y_i \mid \boldsymbol{\beta}, \sigma^2) \pi(\boldsymbol{\beta})$. Usual (uncensored) posterior.
3. $\sigma^2 \mid y_1, \dots, y_n, \boldsymbol{\beta} \propto \prod_{i=1}^n p(y_i \mid \boldsymbol{\beta}, \sigma^2) \pi(\sigma^2)$. Usual (uncensored) posterior.

We give an example in a different context — survival analysis with censored data.

443

Related Example: Survival Analysis with Censored Data

From the WinBUGS manual: *Mice: Weibull regression*

Dellaportas and Smith (1993) analyse data from Grieve (1987) on photocarcinogenicity in four groups, each containing 20 mice, who have recorded a survival time and whether they died or were censored at that time.

A portion of the data, giving survival times in weeks, are shown below.

A * indicates censoring.

Mouse	Irradia control	Vehicle control	Test substan	Positive control
1	12	32	22	27
.....				
18	*40	30	24	12
19	31	37	37	17
20	36	27	29	26

444

The survival distribution is assumed to be Weibull. That is

$$p(t_i, z_i) = r \exp(\beta z_i) t_i^{r-1} \exp\{-\exp(\beta z_i) t_i^r\}$$

where t_i is the failure time of an individual with covariate vector z_i and β is a vector of unknown regression coefficients. The baseline hazard is given by

$$\lambda_0(t_i) = r t_i^{r-1}.$$

Setting $\mu_i = \exp(\beta z_i)$ gives the parameterization

$$t_i \sim \text{Weibull}(r, \mu_i)$$

For censored observations, the survival distribution is a truncated Weibull, with lower bound corresponding to the censoring time. The regression coefficients β are assumed a priori to follow independent Normal distributions with zero mean and “vague” precision 0.0001. The shape parameter r for the survival distribution was given a Gamma(1, 0.0001) prior, which is slowly decreasing on the positive real line.

Median survival for individuals with covariate vector z_i is given by

$$m_i = (\log 2 \exp(-\beta z_i))^{1/r}.$$

445

WinBUGS code

```

model
{
  for(i in 1 : M) {
    for(j in 1 : N) {
      t[i, j] ~ dweib(r, mu[i])I(t.cen[i, j],)
    }
    mu[i] <- exp(beta[i])
    beta[i] ~ dnorm(0.0, 0.001)
    median[i] <- pow(log(2) * exp(-beta[i]), 1/r)
  }
  r ~ dexp(0.001)
  veh.control <- beta[2] - beta[1]
  test.sub <- beta[3] - beta[1]
  pos.control <- beta[4] - beta[1]
}

```

446

```
list(t = structure(.Data = c(12, 1, 21, 25, 11, 26, 27, 30, 13, 12,
21, 20, 23, 25, 23, 29, 35, NA, 31, 36, 32, 27, 23, 12, 18, NA, NA,
38, 29, 30, NA, 32, NA, NA, NA, NA, 25, 30, 37, 27, 22, 26, NA, 28,
19, 15, 12, 35, 35, 10, 22, 18, NA, 12, NA, NA, 31, 24, 37, 29, 27,
18, 22, 13, 18, 29, 28, NA, 16, 22, 26, 19, NA, NA, 17, 28, 26, 12,
17, 26), .Dim = c(4, 20)), t.cen = structure(.Data = c( 0, 0, 0, 0, 0,
0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 40, 0, 0, 0, 0, 0, 0, 0, 0, 40, 40,
0, 0, 0, 40, 0, 40, 40, 40, 40, 0, 0, 0, 0, 0, 0, 10, 0, 0, 0, 0, 0,
0, 0, 0, 0, 24, 0, 40, 40, 0, 0, 0, 0, 0, 0, 0, 0, 0, 0, 20, 0, 0,
0, 0, 29, 10, 0, 0, 0, 0, 0, 0), .Dim = c(4, 20)), M = 4, N = 20)
```

We note a number of tricks in setting up this model.

First, individuals who are censored are given a missing value in the vector of failure times t , whilst individuals who fail are given a zero in the censoring time vector $t.cen$.

447

The truncated Weibull is modelled using $I(t.cen[i],)$ to set a lower bound. Second, we set a parameter $\beta[j]$ for each treatment group j . The contrasts $\beta[j]$ with group 1 (the irradiated control) are calculated at the end. Alternatively, we could have included a grand mean term in the relative risk model and constrained $\beta[1]$ to be zero.

node	mean	sd	MC error	2.5%	median	97.5%	start	sample
median[1]	23.9	1.967	0.05889	20.3	23.89	28.09	1001	10000
median[2]	35.2	3.359	0.04757	29.46	34.93	42.64	1001	10000
median[3]	26.9	2.383	0.0582	22.62	26.79	31.91	1001	10000
median[4]	21.4	1.799	0.03362	18.2	21.32	25.36	1001	10000
pos.control	0.3409	.3457	0.00723	-0.327	0.3429	1.009	1001	10000
r	3.03	0.3182	0.02749	2.388	3.045	3.64	1001	10000
test.sub	-0.351	0.3459	0.004433	-1.035	-0.3541	0.3303	1001	10000
veh.control	-1.16	0.3679	0.005974	-1.893	-1.156	-0.444	1001	10000

448

GEE Approaches

Suppose that if the full data had been observed there would have been n_i observations on each individual, $i = 1, \dots, m$.

We write the usual estimating equation as

$$\mathbf{G}(\boldsymbol{\beta}) = \sum_{i=1}^m \mathbf{D}_i^T \mathbf{W}_i^{-1} \mathbf{R}_i (\mathbf{Y}_i - \boldsymbol{\mu}_i)$$

where \mathbf{R}_i is the diagonal matrix with elements R_{ij} , $j = 1, \dots, n_i$.

For the estimator, $\hat{\boldsymbol{\beta}}$ to be consistent we require \mathbf{G} to be unbiased. The random variables are now \mathbf{Y}, \mathbf{R} and so we have

$$\begin{aligned} \mathbb{E}_{Y,R}[\mathbf{G}(\boldsymbol{\beta})] &= \mathbb{E}_R\{\mathbb{E}_{Y|R}[\mathbf{G}(\boldsymbol{\beta})]\} \\ &= \sum_{i=1}^m \mathbb{E}_{R_i}\{\mathbb{E}_{Y_i|R_i}[D_i^T W_i^{-1} R_i (Y_i - \boldsymbol{\mu}_i)]\} \\ &= \sum_{i=1}^m \mathbb{E}_{R_i}\{D_i^T W_i^{-1} R_i \mathbb{E}_{Y_i|R_i}[Y_i - \boldsymbol{\mu}_i]\} \\ &= \sum_{i=1}^m \mathbb{E}_{R_i}\{D_i^T W_i^{-1} R_i (\mathbb{E}_{Y_i|R_i}[Y_i] - \boldsymbol{\mu}_i)\} \end{aligned}$$

449

Hence, to obtain an unbiased estimating equation we require

$$\mathbb{E}[\mathbf{Y}_i | \mathbf{R}_i, \mathbf{x}_i] = \mathbb{E}[\mathbf{Y}_i | \mathbf{x}_i] = \boldsymbol{\mu}_i$$

so that we are fine under MCAR but not under MAR, since the distribution of $\mathbf{Y}_i | \mathbf{x}_i, \mathbf{R}_i$ is different from that of $\mathbf{Y}_i | \mathbf{x}_i$ under MAR.

To rectify the situation we need to modify the usual estimating equation.

Let

$$\pi_{ij} = E[R_{ij} \mid \mathbf{x}_i, \mathbf{H}_{i,j-1}]$$

where $\mathbf{H}_{i,j-1} = (Y_{i1}, \dots, Y_{i,j-1})$ contains the “history” of responses.

Consider the estimating equation:

$$\sum_{i=1}^m \mathbf{D}_i^T \mathbf{W}_i^{-1} \mathbf{P}_i (\mathbf{Y}_i - \boldsymbol{\mu}_i)$$

where \mathbf{P}_i is a diagonal matrix which contains terms R_{ij}/π_{ij} , for $j = 1, \dots, n_i$.

We have

$$E_Y \left\{ \sum_{i=1}^m E_{R|Y} \left[\mathbf{D}_i^T \mathbf{W}_i^{-1} \mathbf{P}_i (\mathbf{Y}_i - \boldsymbol{\mu}_i) \right] \right\} = E_Y \left\{ \sum_{i=1}^m \mathbf{D}_i^T \mathbf{W}_i^{-1} E_{R|Y} [\mathbf{P}_i] (\mathbf{Y}_i - \boldsymbol{\mu}_i) \right\} = 0$$

since $E[\mathbf{P}_i] = \mathbf{I}$ if π_{ij} is correctly specified.

In both GEE and likelihood we are basically accounting for the biased sampling scheme of MAR; likelihood does this by assuming a model, while GEE adjusts by modeling the probabilities of seeing the data.