

## Approximate Bayesian methods for the analysis of epidemiological data

Jim Young  
Basel Institute for Clinical  
Epidemiology and Biostatistics

## The Swiss HIV Cohort Study

- SCHS enrolls HIV infected adults.
- Visits scheduled every 6 months.
- Measurement of biomarkers (CD4, RNA), cardiovascular risk assessed.
- Drug (start, stop) and disease dates.

## Motivating examples

- Which drugs cause metabolic syndrome – 250 cases, 16 drugs, 7 confounders?
- Why do patients fail salvage therapy – 29 failures in 115 patients?
- Is drug D associated with liver disease – 15 cases and 75 matched controls?
- Statisticians too successful? Questions concern modest effects and little data.

## So why is this a problem?

- Insufficient confounder control by design (restriction or matching).
- Exchangeability within strata defined by a sufficient set of covariates.
- ‘Exposure effectively randomised by natural circumstances.’
- ‘Small sample bias’ – inflated MLEs (even in conditional logistic regression).

## Variable selection is no solution

- 10 to 15 events required per predictor.
- Invariably some variables omitted in an attempt to better estimate others.
- Automatic variable selection, pre-testing, repeated fitting of simplistic models.
- Models and estimates do not replicate.

## How can Bayesian methods help?

- Provide additional information - uninformative priors are pointless.
- Vaguely informative priors – ‘at least reasonable if not liberally inclusive.’
- Sensitivity analysis = alternative priors.
- ‘Shrinkage’ versus data must pull estimate away from the prior.

## Why approximate methods?

- Because MCMC is unrealistic precision.
- ‘Semi-quantitative inference about an adjusted risk comparison.’
- Better to think hard about available background information.
- Use standard software to get a rough answer – anything else is just fantasy.

JY 28/04/10

## Method 1: Hierarchical models

- The ‘multiple exposure problem’: many possible correlated causes (exposures).
- 2nd level (prior) model for correlation between exposures and residual effects.
- Prior estimate of the residual variance.
- ‘Semi-Bayes’ without other priors.
- Fit using GLIMMIX macro.

JY 28/04/10

## Metabolic syndrome

- Marker for heart disease and diabetes
- Some antiviral drugs worse than others?
- 1249 patients starting HAART, 251 develop MS, 16 drugs & 7 covariates.
- Interval censoring - MS only known to occur within an interval between visits.

JY 28/04/10

## Approximate discrete Cox model

$$\log(-\log(p_{ij})) = \alpha_j + X_{ij}^T \beta + W_{ij}^T \theta + \log(\Delta t_{ij})$$

- Risk sets of patients  $i$  at risk of MS at visit  $j$ , given no MS at previous visit.
- Likelihood exact for regular visits.
- Cumulative exposures  $\beta$  correlated if they belong to the same drug class  $\pi$ .

JY 28/04/10

## Second level model for exposures

$$\beta = Z\pi + \delta$$

- Drug class mean is prior estimate of  $\beta$ .
- $\delta$  residual effect with variance  $\sigma^2$ .
- $\sigma^2 = 1/8$ , 95% prior probability hazard ratios for residual effects ( $\delta$ ) within the range 0.50 to 2.0 (ie  $[\ln(2)/1.96]^2 = 1/8$ ).
- Better to over- than under-estimate  $\sigma^2$  ( $\sigma^2 = \infty$  is conventional model).

JY 28/04/10

## GLIMMIX

- GLIMMIX options: ERROR=binomial, LINK=clogl, OFFSET=name.
- Repeated calls to PROC MIX with ...
- MODEL statement: outcome =  $\alpha$  indicators,  $W$  covariate matrix,  $XZ$  class mean matrix.
- RANDOM statement: X exposure matrix / GDATA=residual variance matrix.

JY 28/04/10

## You will need

- %GLIMMIX for version 8 or later (<http://support.sas.com/kb/25/030.html>).
- Process output with macro from Witte (<http://darwin.cwru.edu/~witte/glimmix>).
- Make minor changes to this macro so it works with SAS V8 or later (see Young).

JY 28/04/10

## Results for two boosted PIs

### Associations between MS and drugs: HR (95% CI) / 6 months exposure

Class / Drug	Conventional Full (ie $\sigma^2 = \infty$ )	Stepwise* $\alpha = 0.2$	Hierarchical $\sigma^2 = 1/8$
PI+RTV	-	-	1.0 (0.7-1.5)
atazanavir	0.8 (0.5-1.3)	0.7 (0.5-1.2)	0.9 (0.6-1.3)
indinavir	1.4 (1.1-1.9)	1.4 (1.0-1.8)	1.3 (1.0-1.8)

\* Backwards with high  $\alpha$  better in simulation, forwards with default  $\alpha = 0.05$  gives indinavir HR 1.5 (1.1-1.9)

JY 28/04/10

## Method 2: Data augmented priors

- Create prior that is approximately log-normal for a risk, hazard or odds ratio.
- 'Prior' data representing this distribution added to the real data.
- Use standard software with separate stratum for real data and each prior.
- Approximate versus semi-Bayes.

JY 28/04/10

## Information in data and prior

- Assumes posterior approximately normal with MLE weighting prior and real data by their information (inverse variance).
- So amount of 'shrinkage' (influence of prior) set by spread of prior.
- Aim to shrink unstable estimates from data towards (sensible) prior values.

JY 28/04/10

## Salvage therapy with darunavir

- 'Triple class failure' not so common.
- Salvage with DRV remarkably successful.
- Only 29 patients fail out of 115 – why?
- Relevant factors: patient health, adherence and potency of therapy.
- Assert vaguely informative priors...

JY 28/04/10

## Vaguely informative priors

- Generic risk factors age and sex: 'uncertain direction' median hazard ratio (HR) of 1.0 (95% limits 0.25 to 4).
- Viral load: 'possible harmful' HR 1.5 (0.38 to 6), P(HR>1)=0.72.
- CD4 cells: 'possible beneficial' HR 0.67 (0.17 to 2.7), P(HR>1)=0.28.
- Poor adherence: 'probably harmful' HR 2.0 (0.5 to 8), P(HR>1)=0.84.
- GSS: 'probably beneficial' HR 0.5 (0.13 to 2), P(HR>1)=0.16.
- These priors correspond to normal distributions for the log hazard ratio with variance 0.5.
- GSS: Sum susceptibility to each drug from resistance tests.

JY 28/04/10

## Discrete Cox model (again)

- Visits: 0 - 24, >24 - 48, >48 - 72 weeks.
- Assess VF using variants of FDA time to loss of virological response algorithm.

```
proc logistic data=real descending;
class visit;
model fail / trials = visit age female
rna cd4 pooradhere trt_gss /
clparm=pl link=cloglog offset=LogDays;
run;
```

JY 28/04/10

## Add prior data

- Data are three risk sets, one per visit.
- Add an additional 'visit' per prior.

```
data prior;
input Fail Trials Visit Female Age CD4
RNA PoorAdhere TRT_GSS Logdays; cards;
4 1.E5 4 0 0 0 0 0 0 1
4 1.E5 4 1 0 0 0 0 0 1
...
4 1.E5 9 0 0 0 0 0 0 1
4 2.E5 9 0 0 0 0 0 1 1
;
```

JY 28/04/10

## Check prior and re-run analysis

- Check each prior.

```
proc logistic data=prior descending;
where visit=4;
model fail / trials = female /
clparm=both link=cloglog offset=LogDays;
run;
```

- Append prior data to real data.
- Re-run model.

JY 28/04/10

## What have we learnt?

### Associations between VF and risk factors: HR (95% CI)

Risk Factor	Prior median	MLE	Posterior median
Age	1.0 (0.25-4) Uncertain direction	0.6 (0.4-0.9)	0.6 (0.4-0.9) Certainly beneficial
Female	1.0 (0.25-4) Uncertain direction	2.1 (0.8-5)	1.7 (0.8-4) Probably harmful
TRT_GSS	0.5 (0.13-2) Probably beneficial	1.0 (0.6-1.5)	0.9 (0.6-1.4) Uncertain direction

JY 28/04/10

## Method 3: Propensity scores

- 15 cases of NCPH (liver disease), 75 matched controls.
- Is disease associated with cumulative exposure (per year) to didanosine (DDI)?
- A case series? Is analysis warranted?
- Authors: 'full multivariate model not possible' so they fit bivariate models.

JY 28/04/10

## To adjust or not to adjust?

- Univariate OR for DDI 3.4 (1.5-8.1)
- Matching on confounders limited to only a few factors; otherwise no matches.
- But adding further covariates, additional sparse data bias overwhelms any reduction in bias from confounding.
- Median adjusted OR for DDI among 10 bivariate models 4.2 (1.3-14).

JY 28/04/10

## How to achieve exchangeability?

- Matching ensures diseased and non-diseased in each stratum. Does not ensure exchangeability in each stratum.
- ‘Exposure effectively randomised by natural circumstances.’
- Further adjust but to avoid sparse data bias: (1) impose constraints via DAP (2) adjust by just one additional variable.



JY 28/04/10

## Adjust for a single variable

- Propensity score = probability of exposure given covariates in combined (exposed and unexposed) population.
- Build model for DDI exposure in cohort, plus DAP on DDI and fit as a bivariate model to case control strata.
- Is outcome associated with DDI if both case and control as likely to be exposed?



JY 28/04/10

## Summary – approximate Bayes

- Hierarchical models
  - Data augmented priors
  - Propensity scores (not only Bayesian)
  - PLUS sensitivity to plausible alternatives
  - AND appropriate precision in estimates
- = **Meaningful epidemiological analyses.**



JY 28/04/10

## Acknowledgements

- I have been hugely influenced by the work of Sander Greenland; however...
- The mistakes I've made are my own.
- Martin Rickenbach & clinicians - the quality of the SHCS is outstanding.
- Heiner Bucher inexplicably tolerates my increasing eccentricity as a statistician.



JY 28/04/10