
Bayesian Decision Theory in Biostatistics

David Draper (joint work with **Dimitris Fouskakis,**
Ioannis Ntzoufras and Ken Pietz)

Department of Applied Mathematics and Statistics
University of California, Santa Cruz, USA

draper@ams.ucsc.edu

www.ams.ucsc.edu/~draper

ISBA 2010 World Meeting

Benidorm, Spain

3 June 2010

What Biostatisticians Do

The practice of **statistics** in general (and **biostatistics** in particular) can be roughly divided into **four activities**:

- **Description** of **available information** (e.g., one or more **data sets**) relevant to answering a **question** of interest, without an attempt to **generalize outward** from the available data;
- **Inference** about aspects of the **underlying process** that gave rise to the data;
- **Prediction** of **future data values** under **interesting scenarios**; and
- **Decision-making** (choosing an **action** from among the **available possibilities**, in spite of the **current uncertainty** about **relevant unknowns**), e.g., **experimental or sampling design**.

Description is largely **non-probabilistic** and relatively **uncontroversial**.

Two probability paradigms are in **widespread use** today in biostatistics:

- **Frequentist** probability: Restrict attention to phenomena that are **inherently repeatable** under (essentially) **identical conditions**;

Use of Frequentist and Bayesian Probability in Biostatistics

then, for an event A of interest, $P_F(A)$ is the limiting **relative frequency** with which A occurs in the n (hypothetical) repetitions, as $n \rightarrow \infty$.

- **Bayesian** probability: numerical **weight of evidence** in favor of an uncertain proposition, obeying a series of **reasonable axioms** to ensure that Bayesian probabilities are **coherent (internally logically consistent)**.

Two facts about these paradigms:

- With the **frequentist** approach, **inference is much easier** than (good) **prediction** and **decision-making**.
- For several reasons (e.g., **computing technology**), the **frequentist paradigm dominated work in biostatistics** in the **20th century**.

An **unpleasant by-product** of these two facts is that

In **biostatistical work** it's a **common practice** to use **frequentist inferential tools**, such as **hypothesis testing** and **benefit-only variable selection methods**, for **decision-theoretic purposes** for which they **may not be optimal**.

Four Examples

In this talk, time permitting, I'll describe **four examples** of how **thinking decision-theoretically** can lead to **better results**.

- **Variable selection in generalized linear models** is a familiar task that is usually accomplished in what may be termed a **benefit-only** manner: we try, **using inferential tools**, (e.g.) to find a **subset** of the available predictors that **maximizes predictive accuracy on future data**.

This ignores the **cost of data collection of the predictors**, which may **vary considerably** from one variable to another; **Bayesian decision theory** with an **appropriate utility structure** can improve on this.

References:

- Fouskakis D, Draper D (2008). Comparing stochastic optimization methods for variable selection in binary outcome prediction, with application to health policy. *Journal of the American Statistical Association*, **103**, 1367–1381.
- Fouskakis D, Ntzoufras I, Draper D (2009a). Bayesian variable selection using cost-adjusted BIC, with application to cost-effective measurement of quality of health care. *Annals of Applied Statistics*, **3**, 663–690.

Four Examples (continued)

— Fouskakis D, Ntzoufras I, Draper D (2009b). Population-based reversible jump MCMC for Bayesian variable selection and evaluation under cost constraints. *Journal of the Royal Statistical Society, Series C*, **58**, 383–403.

- When a **clinical trial** has been **adequately planned** (“**appropriately powered**”), as far as **sample size** is concerned, to bring the notions of **clinical** and **statistical significance** into **good agreement** with respect to its **primary objectives**, it may well still be true that it is “**underpowered**” for **secondary subgroup analyses**.
The use of **frequentist multiple comparisons** (inferential) methods in such situations — e.g., to make **choices** about whether to run **new trials** on the **promising subgroups** — is a **bad idea** that can nevertheless be seen in the **literature** (e.g., in a published trial I’m now reanalyzing, assessing the **efficacy of an HIV vaccine**).

Four Examples (continued)

The **problem** (of course) is that (in frequentist language) **multiple comparisons methods are terrified of making type I errors without any concern about type II mistakes.**

The use of **Bayesian decision theory**, to make the **trade-off** explicit in **cost-benefit** terms, can again come to **more sensible conclusions.**

- In **phase II clinical trials**, where the **sample sizes** are typically **fairly small, good frequentist statisticians** (such as the ones with whom I've worked at **Roche** in Switzerland) know that it may well be a **bad idea** to conduct **hypothesis tests** at the usual **0.05 level**, because this does not strike a **sensible balance** between **type I** and **type II error**; the **usual thing** to do (when this problem is realized at all) is to **informally choose a higher type I error rate**, such as **0.2** or **0.25** or **0.3**.

Setting the problem up **decision-theoretically** instead of **inferentially** offers **explicit and non-ad-hoc guidance** on where the **optimal balance** between **type I** and **type II errors** may be found.

Four Examples (continued)

- (joint work with **Ken Pietz**: VA Health Services Research Center of Excellence, Houston, and **Baylor** College of Medicine)

We have as a client one **VA network** interested in **improving quality of care**; here's an example of the data: **101 providers** (physicians), treating **18,763 patients** (this factor fully **nested** in provider factor) with a history of **hypertension**.

Once per quarter each patient is scored as **{in compliance}** or not, where **compliance** = {either at baseline for this quarter, blood pressure **under control**, or if not under control or not measured at baseline then **appropriate action** taken during this quarter}; **87% overall compliance rate**.

Each provider thus has an **observed patient compliance rate**; these range from about **70%** up to nearly **100%**; **70% = bad care?**

Not necessarily: some patients **harder** to bring in compliance (e.g., **sicker**); need to compare **observed** and **expected** compliance rates for each provider, where expected rate based on **relevant patient-level covariates**.

How balance **false positive** and **false negative** rates in **classifying** providers?

Measuring Sickness at Admission

Variable selection (choosing the “best” subset of predictors) in generalized linear models is an old problem, dating back at least to the 1960s, and many methods have been proposed to try to solve it; but virtually all of them ignore an aspect of the problem that can be important: the **cost of data collection of the predictors**.

Example 1. (Fouskakis and Draper, *JASA*, 2008; Fouskakis, Ntzoufras and Draper (FND), *AoAS*, 2009; *JRSS-C*, 2009). In the field of **quality of health care measurement**, patient **sickness at admission** is often assessed by using **logistic regression of mortality within 30 days of admission** on a fairly large number of **sickness indicators** (on the order of **100**) to construct a sickness scale, employing standard **variable selection methods** (e.g., **backward selection** from a model with all predictors) to find an “**optimal**” subset of **10–20** indicators.

Such “**benefit-only**” methods ignore the considerable **differences** among the sickness indicators in **cost of data collection**, an issue that’s **crucial** when admission sickness is used to drive programs (now implemented or

Choosing Utility Function (continued)

under consideration in several countries, including the U.S. and U.K.) that attempt to **identify substandard hospitals** by comparing **observed and expected mortality rates (given admission sickness)**.

When both **data-collection cost** and **accuracy of prediction** of 30-day mortality are considered, a large **variable-selection problem** arises in which **costly variables that do not predict well enough** should be **omitted** from the final scale.

There are **two main ways** to solve this problem — you can (a) put **cost** and **predictive accuracy** on the **same scale** and **optimize**, or (b) **maximize** the latter subject to a **bound** on the former — leading to **three methods**:

- (1) a **decision-theoretic cost-benefit approach** based on **maximizing expected utility** (Fouskakis and Draper, 2008),
- (2) an **alternative cost-benefit approach** based on **posterior model odds** (FND, 2009a), and
- (3) a **cost-restriction-benefit analysis** that **maximizes predictive accuracy** subject to a **bound on cost** (FND, 2009b).

The Data

Data (Kahn et al., *JAMA*, 1990): $p = 83$ **sickness indicators** gathered on **representative sample** of $n = 2,532$ elderly American patients hospitalized in the period 1980–86 with **pneumonia**; original RAND **benefit-only scale** based on **subset** of 14 predictors:

Variable	Cost (U.S.\$)	Correlation	Good?
Total APACHE II score (36-point scale)	3.33	0.39	
Age	0.50	0.17	*
Systolic blood pressure score (2-point scale)	0.17	0.29	**
Chest X-ray congestive heart failure score (3-point scale)	0.83	0.10	
Blood urea nitrogen	0.50	0.32	**
APACHE II coma score (3-point scale)	0.83	0.35	**
Serum albumin (3-point scale)	0.50	0.20	*
Shortness of breath (yes, no)	0.33	0.13	**
Respiratory distress (yes, no)	0.33	0.18	*
Septic complications (yes, no)	1.00	0.06	
Prior respiratory failure (yes, no)	0.67	0.08	
Recently hospitalized (yes, no)	0.67	0.14	
Ambulatory score (3-point scale)	0.83	0.22	
Temperature	0.17	-0.16	*

Decision-Theoretic Cost-Benefit Approach

Approach (1) (decision-theoretic cost-benefit). **Problem formulation:**

Suppose (a) the 30-day **mortality outcome** y_i and data on p **sickness indicators** (x_{i1}, \dots, X_{ip}) have been collected on n individuals sampled exchangeably from a **population** \mathcal{P} of patients with a given disease, and (b) the goal is to **predict** the death outcome for n^* **new patients** who will in the future be sampled exchangeably from \mathcal{P} , (c) on the basis of some or all of the predictors $X_{.j}$, when (d) the **marginal costs of data collection** per patient c_1, \dots, c_p for the $X_{.j}$ **vary considerably**.

What is the **best subset** of the $X_{.j}$ to choose, if a **fixed amount of money** is available for this task and you're **rewarded** based on the **quality** of your predictions?

Since data on **future patients** are **not available**, we use a **cross-validation** approach in which (i) a random subset of n_M observations is drawn for creation of the mortality predictions (the **modeling** subsample) and (ii) the quality of those predictions is assessed on the remaining $n_V = (n - n_M)$ observations (the **validation** subsample, which serves as a proxy for future patients).

Utility Elicitation

Here **utility** is quantified in **monetary terms**, so that **data collection** part of **utility function** is simply **negative of total amount of money** required to gather data on specified predictor subset (**manual data abstraction** from hardcopy patient charts will gradually be replaced by **electronic medical records**, but still widely used in **quality of care studies**).

Letting $I_j = 1$ if $X_{.j}$ is included in a given model (and 0 otherwise), the **data-collection utility** associated with subset $I = (I_1, \dots, I_p)$ for patients in the **validation subsample** is

$$U_D(I) = -n_V \sum_{j=1}^p c_j I_j, \quad (1)$$

where c_j is the **marginal cost per patient of data abstraction** for variable j (the second column in the table above gave examples of these marginal costs).

To measure the **accuracy** of a model's predictions, a metric is needed that quantifies the **discrepancy** between the actual and predicted values, and in this problem **the metric must come out in monetary terms** on a scale comparable to that employed with the data-collection utility.

Utility Elicitation (continued)

In the setting of this problem the outcomes Y_i are **binary death indicators** and the **predicted values** \hat{p}_i , based on statistical modeling, take the form of **estimated death probabilities**.

We use an approach to the comparison of **actual** and **predicted** values that involves **dichotomizing** the \hat{p}_i with respect to a **cutoff**, to mimic the decision-making reality that **actions** taken on the basis of observed-versus-expected quality assessment will have an **all-or-nothing character** at the hospital level (for example, regulators must decide either to subject or not subject a given hospital to a more detailed, more expensive quality audit based on **process criteria**).

In the first step of our approach, given a particular **predictor subset** I , we fit a **logistic regression model** to the **modeling** subsample M and apply this model to **validation** subsample V to create predicted death probabilities \hat{p}_i^I .

In more detail, letting $Y_i = 1$ if patient i dies and 0 otherwise, and taking X_{i1}, \dots, X_{ik} to be the k **sickness predictors** for this patient under model I , the usual **sampling model** which underlies logistic regression in this case is

Utility Elicitation (continued)

$$\begin{aligned} (Y_i | p_i^I) &\stackrel{\text{indep}}{\sim} \text{Bernoulli}(p_i^I), \\ \log\left(\frac{p_i^I}{1-p_i^I}\right) &= \beta_0 + \beta_1 X_{i1} + \dots + \beta_k X_{ik}. \end{aligned} \quad (2)$$

We use **maximum likelihood** to fit this model (as a computationally efficient approximation to Bayesian fitting with relatively diffuse priors), obtaining a vector $\hat{\beta}$ of estimated logistic regression coefficients, from which the **predicted death probabilities** for the patients in subsample V are as usual given by

$$\hat{p}_i^I = \left[1 + \exp\left(-\sum_{j=0}^k \hat{\beta}_j X_{ij}\right) \right]^{-1}, \quad (3)$$

where $X_{i0} = 1$ (\hat{p}_i^I may be thought of as the **sickness score** for patient i under model I).

In the second step of our approach we **classify** patient i in the validation subsample as **predicted dead or alive** according to whether \hat{p}_i^I exceeds or falls short of a **cutoff** p^* , which is chosen — by searching on a discrete grid from 0.01 to 0.99 by steps of 0.01 — to **maximize the predictive accuracy** of model I .

Utility Elicitation (continued)

We then cross-tabulate actual versus predicted death status in a 2×2 **contingency table**, **rewarding** and **penalizing** model I according to the numbers of patients in the **validation sample** which fall into the cells of the right-hand part of the following table.

		Rewards and Penalties		Counts	
		Predicted		Predicted	
		Died	Lived	Died	Lived
Actual	Died	C_{11}	C_{12}	n_{11}	n_{12}
	Lived	C_{21}	C_{22}	n_{21}	n_{22}

The left-hand part of this table records the **rewards and penalties** in US\$.

The **predictive utility** of model I is then

$$U_P(I) = \sum_{l=1}^2 \sum_{m=1}^2 C_{lm} n_{lm}. \quad (4)$$

To **elicit** the **utility values** C_{lm} we reason as follows.

Utility Elicitation (continued)

- (1) Clearly C_{11} (the **reward** for correctly predicting death at 30 days) and C_{22} (the **reward** for correctly predicting living at 30 days) should be **positive**, and C_{12} (the **penalty** for a false prediction of living) and C_{21} (the **penalty** for a false prediction of death) should be **negative**.
- (2) Since it's **easier** to correctly predict that a person lives than dies with these data (the overall pneumonia 30-day death rate in the RAND sample was 16%, so a prediction that every patient lives would be right about **84%** of the time), it's natural to specify that $C_{11} > C_{22}$.
- (3) Since it's arguably **worse** to label a “bad” hospital as “good” than the other way around, one should take $|C_{12}| > |C_{21}|$, and furthermore it's natural that the magnitudes of the **penalties** should exceed those of the **rewards**.
- (4) We completed the utility specification by **eliciting** information from **health experts** in the U.S. and U.K, first to **anchor** C_{21} to the cost of subjecting a “good” hospital to an unnecessary process audit and then to obtain **ratios** relating the other C_{lm} to C_{21} .

Utility Elicitation (continued)

Since the **utility structure** we use is based on the idea that hospitals have to be treated in an **all-or-nothing** way in acting on the basis of their apparent quality, the approach taken was (i) to quantify the **monetary loss** L of incorrectly subjecting a “good” hospital to a detailed but unnecessary process audit and then (ii) to **translate** this from the hospital to the patient level.

Rough **correspondence** may be made between left-hand part of contingency table above at **patient level** and **hospital-level** table with rows representing **truth** (“bad” in row 1, “good” in row 2) and columns representing **decision taken** (“process audit” in column 1, “no process audit” in column 2).

Unnecessary process audits then correspond to cell (2, 1) in these tables (hospitals where a process audit is **not needed** will typically have an **excess** of patients who are predicted to die but actually live).

Discussions with health experts in the U.S. and U.K. suggested that **detailed process audits** cost on the order of $L = \$5,000$ per hospital (in late 1980s U.S. dollars), and RAND data indicated that the mean number of pneumonia patients per hospital per year in the U.S. at the time of the RAND quality of care study was **71.8**.

Utility Elicitation (continued)

This **fixed** C_{21} at approximately $\frac{-\$5,000}{71.8} = -\69.6 .

Our **health experts** judged that C_{12} should be the **largest** in absolute value of the C_{lm} , and averaging across the expert opinions, expressed as orders of magnitude base 2, the elicitation results were $\left| \frac{C_{12}}{C_{21}} \right| = 2$, $\left| \frac{C_{11}}{C_{21}} \right| = \frac{1}{2}$, and $\left| \frac{C_{22}}{C_{21}} \right| = \frac{1}{8}$, finally yielding $(C_{11}, C_{12}, C_{21}, C_{22}) = \mathbf{\$(34.8, -139.2, -69.6, 8.7)}$.

The results in Fouskakis and Draper (2008) use these values; Draper and Fouskakis (2000) present a **sensitivity analysis** on the choice of the C_{lm} which demonstrates **broad stability** of the findings when the utility values mentioned above are **perturbed** in reasonable ways.

With the C_{lm} in hand, the **overall expected utility function** to be maximized over I is then simply

$$E[U(I)] = E[U_D(I) + U_P(I)], \quad (5)$$

where this expectation is over **all possible cross-validation splits** of the data.

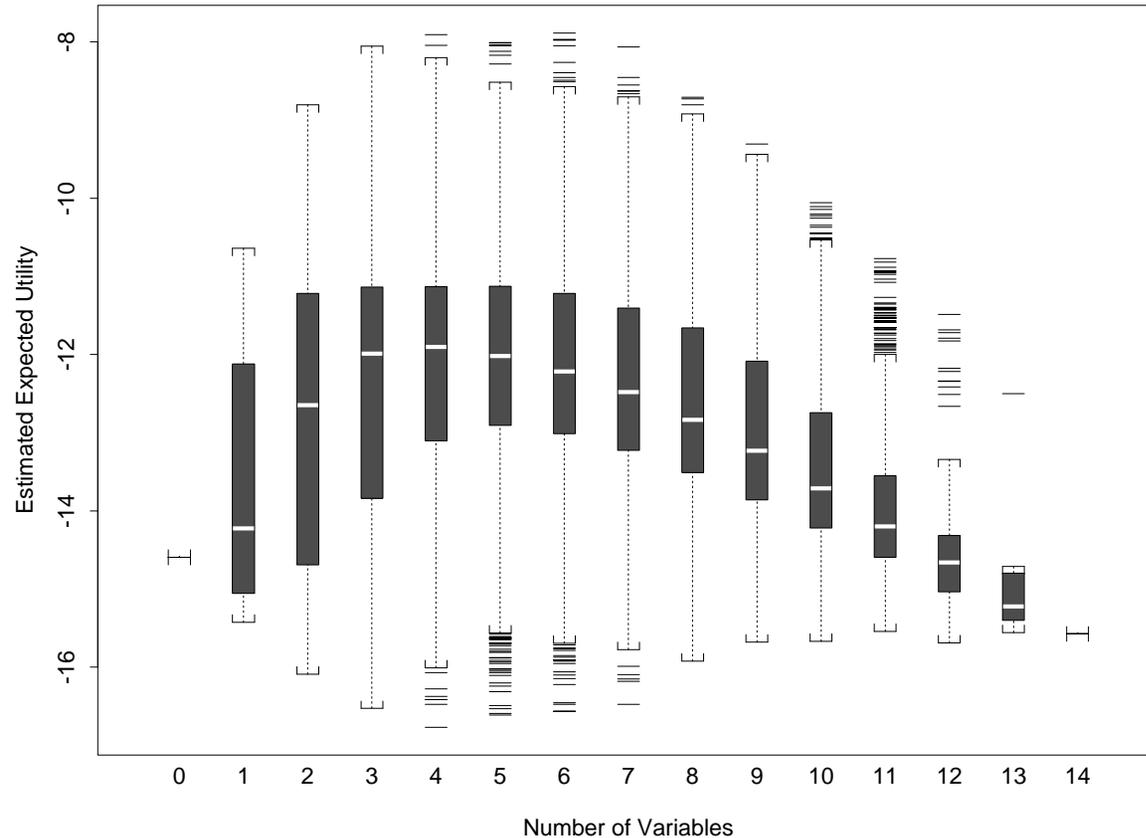
Results

The number of possible cross-validation splits is **far too large** to evaluate the expectation in (5) directly; in practice we therefore use **Monte Carlo methods** to evaluate it, **averaging** over N random modeling and validation **splits**.

Results. We explored this approach in **two settings**:

- a **Small World** created by focusing only on the $p = 14$ variables in the **original RAND scale** ($2^{14} = 16,384$ is a **small enough number of possible models** to do **brute-force enumeration** of the estimated expected utility of all models), and
- the **Big World** defined by all $p = 83$ available predictors ($2^{83} \doteq 10^{25}$ is **far too large** for brute-force enumeration; we compared a variety of **stochastic optimization methods** — including **simulated annealing, genetic algorithms, and tabu search** — on their ability to find **good variable subsets**).

Results: Small World



The **20 best models** included the **same three variables** 18 or more times out of 20, and never included six other variables; the **five best models** were minor variations on each other, and included **4–6 variables** (last column in table on page 10).

Approach (2)

The best models **save almost \$8 per patient** over the full 14-variable model; this would amount to **significant savings** if the observed-versus-expected assessment method were **applied widely**.

Approach (2) (alternative cost-benefit) **Maximizing expected utility**, as in Approach (1) above, is a natural Bayesian way forward in this problem, but (a) the elicitation process was **complicated** and (b) the **utility structure** we examine is only one of a number of plausible alternatives, with utility framed from **only one point of view**; the broader question for a decision-theoretic approach is **whose utility should drive the problem formulation**.

It's well known (e.g., Arrow, 1963; Weerahandi and Zidek, 1981) that **Bayesian decision theory** can be **problematic** when used **normatively** for **group decision-making**, because of **conflicts in preferences** among members of the group; in the context of the problem addressed here, it can be **difficult** to identify a **utility structure acceptable to all stakeholders** (including **patients, doctors, hospitals, citizen watchdog groups, and state and federal regulatory agencies**) in the quality-of-care-assessment process.

Approach (2) (continued)

As an **alternative**, in Approach (2) we propose a **prior distribution** that accounts for the **cost** of each variable and results in a set of **posterior model probabilities** which correspond to a **generalized cost-adjusted version of the Bayesian information criterion** (BIC).

This provides a **principled approach** to performing a **cost-benefit trade-off** that **avoids ambiguities** in identification of an **appropriate utility structure**.

Details. Bayesian **parametric model comparison** and **variable selection** are based on specifying a model m , its likelihood $f(\mathbf{y}|\boldsymbol{\theta}_m, m)$, the prior distribution of model parameters $f(\boldsymbol{\theta}_m|m)$ and the corresponding prior model weight (or probability) $f(m)$, where $\boldsymbol{\theta}_m$ is a parameter vector under model m and \mathbf{y} is the data vector.

Parametric inference is based on the posterior distribution $f(\boldsymbol{\theta}_m|\mathbf{y}, m)$, and quantifying **model uncertainty** by estimating the posterior model probability $f(m|\mathbf{y})$ is also an important issue.

Parametric Model Comparison

Hence, when we consider a set of competing models $\mathcal{M} = \{m_1, m_2, \dots, m_{|\mathcal{M}|}\}$, we focus on the **posterior probability** of model $m \in \mathcal{M}$, defined as

$$\begin{aligned} f(m|\mathbf{y}) &= \frac{f(\mathbf{y}|m)f(m)}{\sum_{m_l \in \mathcal{M}} f(\mathbf{y}|m_l)f(m_l)} = \left(\sum_{m_l \in \mathcal{M}} PO_{m_l, m} \right)^{-1} \\ &= \left[\sum_{m_l \in \mathcal{M}} B_{m_l, m} \frac{f(m_l)}{f(m)} \right]^{-1}, \end{aligned} \quad (6)$$

where $PO_{m_i, m_j} = \frac{f(m_i|\mathbf{y})}{f(m_j|\mathbf{y})}$ is the **posterior model odds** and B_{m_i, m_j} is the **Bayes factor** for comparing models m_i and m_j .

When we limit ourselves in the comparison of only **two models** we typically focus on PO_{m_i, m_j} and B_{m_i, m_j} , which have the desirable property of **insensitivity** to the selection of the model space \mathcal{M} .

By definition the **Bayes factor** is the ratio of the **posterior model odds** over the **prior model odds**; thus **large values** of B_{m_i, m_j} (usually greater than **12**, say) indicate **strong posterior support** of model m_i against model m_j .

Variable Selection in Logistic Regression

The **posterior model probabilities** and **integrated likelihoods** $f(\mathbf{y}|m_i)$ in (6) are **rarely analytically tractable**; we use a combination of **Laplace approximations** and **Markov Chain Monte Carlo** (MCMC) methodology to approximate posterior odds and Bayes factors.

In the sickness-at-admission problem at issue here, we use a simple **logistic regression** model with response $Y_i = 1$ if patient i dies and 0 otherwise.

We further denote by X_{ij} the **sickness predictor variable** j for patient i and by γ_j an **indicator**, often used in Bayesian variable selection problems, taking the value 1 if variable j is included in the model and 0 otherwise; thus in this case $\mathcal{M} = \{0, 1\}^p$, where p is the total number of variables.

In order to map the set of **binary model indicators** γ onto a model m we can use a **representation** of the form $m(\gamma) = \sum_{i=1}^p 2^{i-1} \gamma_i$.

Hence the **model formulation** can be summarized as

$$(Y_i | \gamma) \stackrel{\text{indep}}{\sim} \text{Bernoulli}[p_i(\gamma)],$$
$$\eta_i(\gamma) = \log \left[\frac{p_i(\gamma)}{1 - p_i(\gamma)} \right] = \sum_{j=0}^p \beta_j \gamma_j X_{ij}, \quad (7)$$

Prior on Model Parameters

$$\eta(\gamma) = \mathbf{X} \text{diag}(\gamma) \boldsymbol{\beta} = \mathbf{X}_{\gamma} \boldsymbol{\beta}_{\gamma},$$

defining $X_{i0} = 1$ for all $i = 1, \dots, n$ and $\gamma_0 = 1$ with **prior probability one** since here the intercept is always included in all models.

Here $p_i(\gamma)$ is the **death probability** (which may be thought of as the **sickness score**) for patient i under model γ , $\boldsymbol{\eta}(\gamma) = [\eta_1(\gamma), \dots, \eta_n(\gamma)]^T$,

$$\boldsymbol{\gamma} = (\gamma_0, \gamma_1, \dots, \gamma_p)^T, \boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_p)^T, \text{ and}$$

$\mathbf{X} = (X_{ij}, i = 1, \dots, n; j = 0, 1, \dots, p)$; the vector $\boldsymbol{\beta}_{\gamma}$ stands for the subvector of $\boldsymbol{\beta}$ which is included in the model specified by γ , i.e.,

$\boldsymbol{\beta}_{\gamma} = (\beta_i : \gamma_i = 1, i = 0, 1, \dots, p)$, and is equivalent to the $\boldsymbol{\theta}_m$ vector defined above; similarly \mathbf{X}_{γ} is the submatrix of \mathbf{X} with columns corresponding to variables included in the model specified by γ .

Prior on model parameters. We proceed in **two steps**:

- (1) First we build a **prior** on $\boldsymbol{\beta}$ that is a modified version of the **unit information prior** for this problem (to avoid **Lindley's paradox**); then
- (2) We **adjust** this prior for **differences in marginal costs** of variables.

Sensitivity to Prior Variance

Step (1). One important problem in **Bayesian model evaluation** using **posterior model probabilities** is their **sensitivity** to the **prior variance** of the model parameters: large variance of the β_γ (used to represent prior ignorance) will **increase** the posterior probabilities of the **simpler** models considered in the model space \mathcal{M} (**Lindley's paradox**).

We address this issue by using ideas proposed by Ntzoufras *et al.* (2003): we use a **prior distribution** of the form

$$f(\beta_\gamma|\gamma) = N(\mu_\gamma, \Sigma_\gamma) \quad (8)$$

with **prior covariance matrix** given by $\Sigma_\gamma = n \left[\mathcal{J}(\beta_\gamma) \right]^{-1}$, where n is the total sample size and $\mathcal{J}(\beta_\gamma)$ is the information matrix

$$\mathcal{J}(\beta_\gamma) = \mathbf{X}_\gamma^T \mathbf{W}_\gamma \mathbf{X}_\gamma;$$

here \mathbf{W}_γ is a diagonal matrix which in the Bernoulli case takes the form

$$\mathbf{W}_\gamma = \text{diag} \{p_i(\gamma)[1 - p_i(\gamma)]\}.$$

Unit Information Prior

This is the **unit information prior** of Kass and Wasserman (1996), which corresponds to adding **one data point** to the data.

Here we use this prior as a **base**, but we specify $p_i(\gamma)$ in the information matrix according to our prior information; in this manner we **avoid (even minimal) reuse** of the data in the prior.

When little prior information is available, a reasonable prior mean for β_γ is $\mu_\gamma = \mathbf{0}$.

This corresponds to a prior mean on the log-odds scale of zero, from which a sensible prior estimate for all model probabilities is $p_i(\gamma) = 1/2$; with this choice (8) becomes

$$f(\beta_\gamma|\gamma) = N\left[\mathbf{0}, 4n \left(\mathbf{X}_\gamma^T \mathbf{X}_\gamma\right)^{-1}\right]. \quad (*) \quad (9)$$

This prior distribution can also be motivated by combining the idea of **imaginary data** with the **power prior** approach of Chen *et al.* (2000); it turns out that (9) introduces additional information to the posterior equivalent to adding **one data point** to the likelihood and therefore we support **a priori** the simplest model with a weight of one data point.

Laplace Approximation

Step (2). To introduce **costs** we again proceed in **two sub-steps**:

(2a) First we specify a **Laplace approximation** (and the **BIC approximation** that corresponds to it) for the **posterior model odds** in our problem, using the **prior** in Step (1), and

(2b) Then we see how to **adjust** the approximations in Step (2a) to account for **cost differences** among the variables.

Step (2a). We denote by $PO_{k\ell}$ the **posterior odds** of model $\gamma^{(k)}$ versus model $\gamma^{(\ell)}$; then we have

$$-2 \log PO_{k\ell} = -2 \left[\log f(\gamma^{(k)} | \mathbf{y}) - \log f(\gamma^{(\ell)} | \mathbf{y}) \right]. \quad (10)$$

Following the approach of Raftery (1996), we can approximate the posterior distribution of a model γ using the following **Laplace approximation**:

$$\begin{aligned} -2 \log f(\gamma | \mathbf{y}) &= -2 \log f(\mathbf{y} | \tilde{\boldsymbol{\beta}}_{\gamma}, \gamma) - 2 \log f(\tilde{\boldsymbol{\beta}}_{\gamma} | \gamma) - d_{\gamma} \log(2\pi) \\ &\quad - \log |\boldsymbol{\Psi}_{\gamma}| - 2 \log f(\gamma) + O(n^{-1}), \end{aligned} \quad (11)$$

Details

where $\tilde{\beta}_\gamma$ is the posterior mode of $f(\beta_\gamma|\mathbf{y}, \gamma)$, $d_\gamma = \sum_{j=0}^p \gamma_j$ is the dimension of the model γ , and Ψ_γ is minus the inverse of the Hessian matrix of $h(\beta_\gamma) = \log f(\mathbf{y}|\beta_\gamma, \gamma) + \log f(\beta_\gamma|\gamma)$ evaluated at the posterior mode $\tilde{\beta}_\gamma$.

Under the **model formulation** given by equation (7) and the **prior distribution** (9) we have that

$$\begin{aligned} \Psi_\gamma &= \left[- \frac{\partial^2 \log f(\mathbf{y}|\beta_\gamma, \gamma)}{\partial \beta_\gamma^2} \Big|_{\beta_\gamma = \tilde{\beta}_\gamma} - \frac{\partial^2 \log f(\beta_\gamma|\gamma)}{\partial \beta_\gamma^2} \Big|_{\beta_\gamma = \tilde{\beta}_\gamma} \right]^{-1} \\ &= \left(\mathbf{X}_\gamma^T \text{diag} \left\{ \frac{\exp(\mathbf{X}_{\gamma,i} \tilde{\beta}_\gamma)}{[1 + \exp(\mathbf{X}_{\gamma,i} \tilde{\beta}_\gamma)]^2} + \frac{1}{4n} \right\} \mathbf{X}_\gamma \right)^{-1}, \end{aligned} \quad (12)$$

where $\mathbf{X}_{\gamma,i}$ is row i of the matrix \mathbf{X}_γ for $i = 1, \dots, n$.

By substituting the **prior** (9) in expression (11) we get

$$-2 \log f(\gamma|\mathbf{y}) = -2 \log f(\mathbf{y}|\tilde{\beta}_\gamma, \gamma) + \phi(\gamma) - 2 \log f(\gamma) + O(n^{-1}), \quad (13)$$

Penalized Log Likelihood Ratio

$$\text{where } \phi(\gamma) = \frac{1}{4n} \tilde{\beta}_\gamma^T \mathbf{X}_\gamma^T \mathbf{X}_\gamma \tilde{\beta}_\gamma + d\gamma \log(4n) + \log \frac{|\Psi_\gamma^{-1}|}{|\mathbf{X}_\gamma^T \mathbf{X}_\gamma|}. \quad (14)$$

From the above expression it's clear that the logarithm of a posterior model probability can be regarded as a **penalized log-likelihood** evaluated at the posterior mode of the model, in which the term $\phi(\gamma) - 2 \log f(\gamma)$ can be interpreted as the **penalty** imposed upon the log-likelihood.

In **pairwise model comparisons**, we can directly use the **posterior model odds** (10), which can now be written as

$$\begin{aligned} -2 \log PO_{k\ell} &= -2 \log \left\{ \frac{f(\mathbf{y} | \tilde{\beta}_{\gamma^{(k)}}, \gamma^{(k)})}{f(\mathbf{y} | \tilde{\beta}_{\gamma^{(\ell)}}, \gamma^{(\ell)})} \right\} + \phi(\gamma^{(k)}) - \phi(\gamma^{(\ell)}) \\ &\quad - 2 \log \frac{f(\gamma^{(k)})}{f(\gamma^{(\ell)})} + O(n^{-1}). \end{aligned} \quad (15)$$

Therefore, the comparison of the two models is based on a **penalized log-likelihood ratio**, where the penalty is now given by

$$\psi(\gamma^{(k)}, \gamma^{(\ell)}) = \phi(\gamma^{(k)}) - \phi(\gamma^{(\ell)}) - 2 \log \frac{f(\gamma^{(k)})}{f(\gamma^{(\ell)})}.$$

Decomposing the Penalty Term

Each **penalty term** is divided into **two parts**: $\phi(\gamma)$ and $-2 \log f(\gamma)$.

The first term, $\phi(\gamma)$, has its source in the **marginal likelihood** $f(\mathbf{y}|\gamma)$ of model γ and can be thought of as a measure of **discrepancy** between the **data** and the **prior information** for the model parameters; the second part comes from the **prior model probabilities** $f(\gamma)$.

Indifference on the space of all models, usually expressed by the **uniform distribution** (i.e., $f(\gamma) \propto 1$), eliminates the second term from the model comparison procedure, since the penalty term in (15) will then be based only on the difference of the first penalty terms $\phi(\gamma^{(k)}) - \phi(\gamma^{(\ell)})$.

For this reason the penalty term $\phi(\gamma)$ is the **imposed penalty** which appears in the penalized log-likelihood expression of the **Bayes factor** $BF_{k\ell}$ with a uniform prior on model space.

A **simpler but less accurate** approximation of $\log PO_{k\ell}$ can be obtained following the arguments of Schwartz (1978):

BIC Approximation

$$\begin{aligned} -2 \log PO_{k\ell} &= -2 \log \left[\frac{f(\mathbf{y} | \hat{\beta}_{\gamma^{(k)}}, \gamma^{(k)})}{f(\mathbf{y} | \hat{\beta}_{\gamma^{(\ell)}}, \gamma^{(\ell)})} \right] + \left(d_{\gamma^{(k)}} - d_{\gamma^{(\ell)}} \right) \log n \\ &\quad - 2 \log \frac{f(\gamma^{(k)})}{f(\gamma^{(\ell)})} + O(1) \\ &= BIC_{k\ell} - 2 \log \frac{f(\gamma^{(k)})}{f(\gamma^{(\ell)})} + O(1), \end{aligned} \tag{16}$$

where $BIC_{k\ell}$ is the **Bayesian Information Criterion** for choosing between models $\gamma^{(k)}$ and $\gamma^{(\ell)}$ and $\hat{\beta}_{\gamma}$ is vector of maximum likelihood estimates of β_{γ} .

Since $BIC_{k\ell}$ is an $O(1)$ approximation, it might **diverge** from the exact value of the logarithm of the Bayes factor even for large samples; even so, it has often been shown to provide a **reasonable measure of evidence** (for finite n) and its straightforward calculation has encouraged its **widespread use** in practice.

Step (2b). From the above argument and equations (13) and (15), it's clear that an **additional penalty** can be directly imposed on the posterior model probabilities and odds via the **prior model probabilities** $f(\gamma)$.

Cost Adjustment

Therefore we may use **prior model probabilities** to induce **prior preferences** for specific variables depending on their **costs**.

For this reason we propose to use **prior model probabilities** of the form

$$\boxed{(*)} \quad f(\gamma_j) \propto \exp \left[-\frac{\gamma_j}{2} \left(\frac{c_j - c_0}{c_0} \right) \log n \right] \quad \text{for } j = 1, \dots, p, \quad (17)$$

where c_j is the **marginal cost per observation** for variable X_j and (as will be seen below) the desire for our approach to yield a **cost-adjusted generalization of BIC** compels the definition $c_0 = \min\{c_j, j = 1, \dots, p\}$.

We further assume that the **constant term** is included in all models by specifying $f(\gamma_0 = 1) = 1$, resulting in

$$-2 \log f(\gamma) = \sum_{j=1}^p \gamma_j \frac{c_j}{c_0} \log n - d\gamma \log n + 2 \sum_{j=1}^p \log \left[1 + n^{-\frac{1}{2}} \left(1 - \frac{c_j}{c_0} \right) \right]. \quad (18)$$

If all variables have the **same cost** or we're indifferent concerning the cost then we can set $c_j = c_0$ for $j = 1, \dots, p$, which reduces to the **uniform prior** on model space ($f(\gamma) \propto 1$) and posterior odds equal to the **usual Bayes factor**.

Cost Adjustment (continued)

When comparing two models $\gamma^{(k)}$ and $\gamma^{(\ell)}$, the **additional penalty** imposed on the log-likelihood ratio due to the **cost-adjusted prior model probabilities** is given by

$$\begin{aligned} -2 \log \left[\frac{f(\gamma^{(k)})}{f(\gamma^{(\ell)})} \right] &= \sum_{j=1}^p \left(\gamma_j^{(k)} - \gamma_j^{(\ell)} \right) \frac{c_j}{c_0} \log n - \left(d_{\gamma^{(k)}} - d_{\gamma^{(\ell)}} \right) \log n \\ &= \left[\frac{C_{\gamma^{(k)}} - C_{\gamma^{(\ell)}}}{c_0} - \left(d_{\gamma^{(k)}} - d_{\gamma^{(\ell)}} \right) \right] \log n, \end{aligned} \quad (19)$$

where $C_{\gamma} = \sum_{j=1}^p \gamma_j c_j$ is the **total cost** of model γ ; thus two models of the **same dimension and cost** will have the **same prior weight**.

In the simpler case where we compare **two nested models** that differ only on the status of variable j , the prior model ratio simplifies to

$$-2 \log \left[\frac{f(\gamma_j = 1, \gamma_{\setminus j})}{f(\gamma_j = 0, \gamma_{\setminus j})} \right] = \left(\frac{c_j}{c_0} - 1 \right) \log n, \quad (20)$$

where $\gamma_{\setminus j}$ is the vector of γ **excluding** element γ_j .

Cost-Adjusted Laplace Approximation

The above expression can be viewed as a **prior penalty** for including the variable j in the model, while the term $\left(\frac{c_j}{c_0} - 1\right)$ can be interpreted as the **proportional additional penalty** imposed upon $(-2 \log BF)$ if the variable X_j is included in the model due to its **increased cost**.

Using the **prior model odds** (19) in the **approximate posterior model odds** (15) we obtain

$$-2 \log PO_{k\ell} = -2 \log \left[\frac{f(\mathbf{y} | \tilde{\beta}_{\gamma^{(k)}}, \gamma^{(k)})}{f(\mathbf{y} | \tilde{\beta}_{\gamma^{(\ell)}}, \gamma^{(\ell)})} \right] + \psi(\gamma^{(k)}, \gamma^{(\ell)}) + O(n^{-1}), \quad (21)$$

where the **penalty term** is given by

$$\begin{aligned} \psi(\gamma^{(k)}, \gamma^{(\ell)}) &= \frac{1}{4n} \left(\tilde{\beta}_{\gamma^{(k)}}^T \mathbf{X}_{\gamma^{(k)}}^T \mathbf{X}_{\gamma^{(k)}} \tilde{\beta}_{\gamma^{(k)}} - \tilde{\beta}_{\gamma^{(\ell)}}^T \mathbf{X}_{\gamma^{(\ell)}}^T \mathbf{X}_{\gamma^{(\ell)}} \tilde{\beta}_{\gamma^{(\ell)}} \right) \\ &\quad + \left(d_{\gamma^{(k)}} - d_{\gamma^{(\ell)}} \right) \log(4) + \log \frac{|\Psi_{\gamma^{(k)}}^{-1}|}{|\mathbf{X}_{\gamma^{(k)}}^T \mathbf{X}_{\gamma^{(k)}}|} \\ &\quad - \log \frac{|\Psi_{\gamma^{(\ell)}}^{-1}|}{|\mathbf{X}_{\gamma^{(\ell)}}^T \mathbf{X}_{\gamma^{(\ell)}}|} + \frac{C_{\gamma^{(k)}} - C_{\gamma^{(\ell)}}}{c_0} \log n. \end{aligned} \quad (22)$$

Cost-Adjusted BIC

Finally we consider the **BIC-based approximation** (16) to the logarithm of the posterior model odds with the prior model odds (19), yielding (*)

$$-2 \log PO_{k\ell} = -2 \log \left[\frac{f(\mathbf{y} | \hat{\beta}_{\gamma^{(k)}}, \gamma^{(k)})}{f(\mathbf{y} | \hat{\beta}_{\gamma^{(\ell)}}, \gamma^{(\ell)})} \right] + \frac{C_{\gamma^{(k)}} - C_{\gamma^{(\ell)}}}{c_0} \log n + O(1). \quad (23)$$

The **penalty term** $d\gamma \log n$ of model γ used in (16) has been replaced in the above expression by the **cost-dependent penalty** $c_0^{-1} C_\gamma \log n$; **ignoring costs** is equivalent to taking $c_j = c_0$ for all j , yielding $c_0^{-1} C_\gamma = d\gamma$, the **original BIC expression**.

Therefore, we may interpret the quantity $\log n$ as the **imposed penalty** for each variable included in the **model γ when no costs are considered** (or when costs are equal).

Moreover, this baseline penalty term is inflated **proportionally** to the cost ratio $\frac{c_j}{c_0}$ for each variable X_j ; for example, if the cost of a variable X_j is **twice** the minimum cost ($c_j = 2c_0$) then the imposed penalty is equivalent to **adding two variables with the minimum cost**.

MCMC Implementation

For this reason, (23) can be considered as a **cost-adjusted generalization of BIC** when prior model probabilities of type (17) are adopted.

MCMC implementation. As noted earlier, in our quality of care study with $p = 83$ predictors there are on the order of 10^{25} possible models.

In such situations, **sampling algorithms** will not be able to estimate posterior model probabilities with **high accuracy** in a reasonable amount of CPU time due to the **large model space**.

For this reason, we implemented the following **two-step method**:

(1) First we use a **model search tool** to identify variables with **high marginal posterior inclusion probabilities** $f(\gamma_j | \mathbf{y})$, and we create a **reduced model space** consisting only of those variables whose marginal probabilities are above a **threshold value**.

According to Barbieri and Berger (2004) this method of selecting variables based on their **marginal probabilities** may lead to the identification of models with **better predictive abilities** than approaches based on maximizing posterior model probabilities.

MCMC Implementation (continued)

Although Barbieri and Berger proposed **0.5** as a threshold value for $f(\gamma_j = 1|\mathbf{y})$, we used the lower value of **0.3**, since our aim was only to **identify and eliminate** variables not contributing to models with high posterior probabilities.

(2) Then we use a **model search tool** in the **reduced model space** to estimate **posterior model probabilities** (and the corresponding odds).

To ensure **stability** of our findings we explored the use of **two model search tools** in step (1):

- A **reversible-jump MCMC algorithm** (RJMCMC), as implemented for variable selection in generalized linear models by Dellaportas *et al.* (2002) and Ntzoufras *et al.* (2003); and
- the **MCMC model composition** (MC^3) algorithm (Madigan and York, 1995).

More specifically, we implemented **reversible-jump moves within Gibbs** for the model indicators γ_j , by proposing the new model to differ from the current one in each step by a **single term** j with probability one.

MCMC Implementation (continued)

The **algorithm** can be summarized as follows:

- (1) For $j = 1, \dots, p$, use **RJMCMC** to compare the current model γ with the proposed one γ' with components $\gamma'_j = 1 - \gamma_j$ and $\gamma'_k = \gamma_k$ for $k \neq j$ with probability one; the **updating sequence** of γ_j is randomly determined in each step.
- (2) For $j = 0, \dots, p$, if $\gamma_j = 1$ then **generate** model parameters β_j from the corresponding posterior distribution $f(\beta_j | \beta_{\setminus j}, \gamma, \mathbf{y})$, otherwise set $\beta_j = 0$.

In our context the **MC³ algorithm** may be summarized by the following steps:

- (1) For $j = 1, \dots, p$, propose a **move** from the current model γ to a new one γ' with components $\gamma'_j = 1 - \gamma_j$ and $\gamma'_k = \gamma_k$ for $k \neq j$ with probability one; the **updating sequence** of γ_j is randomly determined in each step.
- (2) **Accept** the proposed model γ' with probability

$$\alpha = \min \left[1, \frac{f(\gamma' | \mathbf{y})}{f(\gamma | \mathbf{y})} \right] = \min (1, PO_{\gamma, \gamma'}).$$

MCMC Implementation (continued)

Since the **posterior model odds** $PO_{\gamma, \gamma'}$ used in MC^3 are **not analytically available** here, we also explored **two methods** for calculating them — approximating the acceptance probabilities with **cost-adjusted Laplace** (equation 21) and **cost-adjusted BIC** (equation 23) — and in addition we further explored one additional form of **sensitivity analysis**: initializing the MCMC runs at the **null model** (with no predictors) and the **full model** (with all predictors).

All of this was done both for the **benefit-only analysis** (specified by **setting all variable costs equal**) and the **cost-benefit approach**.

In moving from the **full** to the **reduced** model space to implement step (1) of our two-step method, for both the benefit-only and cost-benefit analyses we found a **striking level of agreement** — across (a) the two model search tools, (b) the two methods to approximate the acceptance probabilities in MC^3 , and (c) the two choices for initializing the MCMC runs — in the **subset of variables** defining the reduced model space; this made it **unnecessary** to perform a similar sensitivity analysis in step (2).

Results

Results are therefore presented below **only for RJMCMC** (starting from the full model).

Convergence of the RJMCMC algorithm was checked using **ergodic mean plots** of the **marginal inclusion probabilities** for the full model space and the **posterior model probabilities** for the reduced space.

In what follows we refer to the **cost-benefit results** as “**RJMCMC**,” but we could equally well have used the term “ **MC^3 with cost-adjusted BIC**” (or just “**cost-adjusted BIC**” for short), because the results from the two methods were in such **close agreement**.

Results. The table below presents the **marginal posterior probabilities** of the variables that exceeded the threshold value of 0.30, in each of the **benefit-only** and **cost-benefit** analyses, together with their data collection costs (in minutes of abstraction time rather than US\$), in the **Big World** of all 83 predictors.

In both the **benefit-only** and **cost-benefit** situations our methods reduced the initial list of $p = 83$ available candidates down to **13** predictors.

Results (continued)

Index	Variable Name	Cost	Marginal Posterior Probabilities	
			Analysis	
			Benefit-Only	Cost-Benefit
1	SBP Score	0.50	0.99	0.99
2	Age	0.50	0.99	0.99
3	Blood Urea Nitrogen	1.50	1.00	0.99
4	Apache II Coma Score	2.50	1.00	
5	Shortness of Breath Day 1?	1.00	0.97	0.79
8	Septic Complications?	3.00	0.88	
12	Initial Temperature	0.50	0.98	0.96
13	Heart Rate Day 1	0.50		0.34
14	Chest Pain Day 1?	0.50		0.39
15	Cardiomegaly Score	1.50	0.71	
27	Hematologic History Score	1.50	0.45	
37	Apache Respiratory Rate Score	1.00	0.95	0.32
46	Admission SBP	0.50	0.68	0.90
49	Respiratory Rate Day 1	0.50		0.81
51	Confusion Day 1?	0.50		0.95
70	Apache pH Score	1.00	0.98	0.98
73	Morbid + Comorbid Score	7.50	0.96	
78	Musculoskeletal Score	1.00		0.54

Note that the **most expensive** variables with high marginal posterior probabilities in the **benefit-only** analysis were **absent** from the set of promising variables in the **cost-benefit** analysis (e.g., **Apache II Coma Score**).

Results (continued)

Common variables in both analyses: $X_1 + X_2 + X_3 + X_5 + X_{12} + X_{70}$

Benefit-Only Analysis

k	Common Variables Within Each Analysis	Additional Variables	Model Cost	Posterior Probabilities	PO_{1k}
1	$X_4 + X_{15} + X_{37} + X_{73}$	$+X_8 + X_{27} + X_{46}$	22.5	0.3066	1.00
2		$+X_8 + X_{27}$	22.0	0.1969	1.56
3		$+X_8$	20.5	0.1833	1.67
4		$+X_{27} + X_{46}$	19.5	0.0763	4.02
5				17.5	0.0383

Cost-Benefit Analysis

k	Common Variables Within Each Analysis	Additional Variables	Model Cost	Posterior Probabilities	PO_{1k}
1	$X_{46} + X_{51}$	$+X_{49} + X_{78}$	7.5	0.1460	1.00
2		$+X_{14} + X_{49} + X_{78}$	7.5	0.1168	1.27
3		$+X_{13} + X_{49} + X_{78}$	7.5	0.0866	1.69
4		$+X_{13} + X_{14} + X_{49} + X_{78}$	8.0	0.0665	2.20
5		$+X_{14} + X_{49}$	7.0	0.0461	3.17
6		$+X_{49}$	6.5	0.0409	3.57
7		$+X_{37} + X_{78}$	7.5	0.0382	3.82
8		$+X_{13} + X_{14} + X_{49}$	7.5	0.0369	3.96
9		$+X_{13}$	6.5	0.0344	4.25

Results (continued)

	Analysis		Percentage
	Benefit-Only	Cost-Benefit	Difference
Minimum Deviance	1553.2	1635.8	+5.3
Median Deviance	1564.5	1644.8	+5.1
Cost	22.5	7.5	-66.7
Dimension	13	10	-23.1

The table above presents a comparison of **measures of fit, cost and dimensionality** between the best models in the reduced model space of the benefit-only and cost-benefit analyses (percentage difference is in relation to benefit-only).

- The **deviance statistic** for the benefit-only RAND model summarized in Table 1 turned out to be **1587.3** (achieved with **14** predictors), **substantially worse** than the median deviance (**1564.5**, achieved with **13** predictors) of the best model visited by the **benefit-only** approach we investigate; in other words, in this case study, **frequentist backward selection** from the model with all predictors (the RAND approach) was **substantially out-performed** by Bayesian RJMCMC.

Results (continued)

	Analysis		Percentage
	Benefit-Only	Cost-Benefit	Difference
Minimum Deviance	1553.2	1635.8	+5.3
Median Deviance	1564.5	1644.8	+5.1
Cost	22.5	7.5	-66.7
Dimension	13	10	-23.1

- The minimum and median values of the posterior distribution of the **deviance** statistic for the benefit-only analysis were **lower** by a **relatively modest 5.3% and 5.1%** compared to the corresponding values of the cost-benefit analysis, but the **cost** of the best model in the cost-benefit analysis was almost **67% lower** than that for the benefit-only analysis; similarly, the dimensionality of the best model in the cost-benefit analysis was about **23% lower** than that for the benefit-only analysis.

These values indicate that the **loss of predictive accuracy** with the **cost-benefit analysis** is **small** compared to the **substantial gains** achieved in **cost** and **reduced model complexity**.

Utility Versus Cost-Adjusted BIC

Index	Variable Name	Cost (Minutes)	Method		
			Utility Good?	RJMCMC Good? Posterior Probability	
1	Systolic Blood Pressure Score (2-point scale)	0.5	**	**	0.99
2	Age	0.5	*	**	0.99
3	Blood Urea Nitrogen	1.5	**	**	1.00
4	APACHE II Coma Score (3-point scale)	2.5	**	**	1.00
5	Shortness of Breath Day 1 (yes, no)	1.0	**	**	0.99
6	Serum Albumin (3-point scale)	1.5	*	**	0.55
7	Respiratory Distress (yes, no)	1.0	*	**	0.92
8	Septic Complications (yes, no)	3.0			0.00
9	Prior Respiratory Failure (yes, no)	2.0			0.00
10	Recently Hospitalized (yes, no)	2.0			0.00
12	Initial Temperature	0.5	*	**	0.95
17	Chest X-ray Congestive Heart Failure Score (3-point scale)	2.5			0.00
18	Ambulatory Score (3-point scale)	2.5			0.00
48	Total APACHE II Score (36-point scale)	10.0			0.00

It's clear that the **utility** and **cost-adjusted BIC** approaches have reached **nearly identical conclusions** in the **Small World** of $p = 14$ predictors.

Utility Versus Cost-Adjusted BIC (continued)

With $p = 83$ the **agreement** between the two methods is also **strong** (although not as strong as with $p = 14$): using a **star system** for variable importance given in FND (2007a), **60** variables were **ignored** by both methods, **8** variables had **identical** star patterns, **3** variables were chosen as **important by both methods** but with different star patterns, **10** variables were marked as important by the utility approach and not by RJMCMC, and **2** variables were singled out by RJMCMC and not by utility: thus the two methods **substantially** agreed on the importance of **71 (86%) of the 83 variables**.

p	Method	Model	Median		
			Cost	Deviance	LS_{CV}
14	RJMCMC	$X_1 + X_2 + X_3 + X_4 + X_5 + X_6 + X_7 + X_{12}$	9.0	1654	-0.329
		$X_1 + X_2 + X_3 + X_4 + X_5 + X_7 + X_{12}$	7.5	1676	-0.333
	Utility	$X_1 + X_3 + X_4 + X_5$	5.5	1726	-0.342
83	RJMCMC	$X_1 + X_2 + X_3 + X_5 + X_{12}$ $+ X_{46} + X_{49} + X_{51} + X_{70} + X_{78}$	7.5	1645	-0.327
	Utility	$X_1 + X_3 + X_4 + X_{12}$ $+ X_{46} + X_{49} + X_{57}$	6.5	1693	-0.336

To the extent that the two methods **differ**, the **utility** method favors models that **cost somewhat less** but also **predict somewhat less well**.

Utility Versus Cost-Adjusted BIC (continued)

The fact that the **two methods** may yield **somewhat different results** in **high-dimensional problems** does not mean that either is **wrong**; they are both **valid solutions** to **similar but not identical problems**.

Both methods lead to **noticeably better models** (in a **cost-benefit** sense) than frequentist or Bayesian **benefit-only** approaches, when — as is often the case — **cost** is an issue that must be included in the **problem formulation** to arrive at a **policy-relevant solution**.

Summary. In **comparing two or more models**, to say whether one is **better** than another I have to face the question: **better for what purpose?**

- This makes **model specification** a **decision problem**: I need to either
- (a) elicit a **utility structure** that's specific to the **goals** of the current study and **maximize expected utility** to find the best models, or
 - (b) (if (a) is too hard, e.g., because the problem has a **group decision** character) I can look for a **principled alternative** (like the **cost-adjusted Laplace and BIC methods** described here) that **approximates** the utility approach while **avoiding ambiguities in utility specification**.

HIV-1 Vaccine Efficacy

Recall two of the main points in this talk: (1) **Inference** and **decision-making** are **not the same thing**. (2) People sometimes use **inferential tools** to make an **implied decision** when **decision-making methods** lead to a **better choice**.

Example 2: A randomized controlled trial of an **rgp120 vaccine** against **HIV** (rgp120 HIV Vaccine Study Group (2005). Placebo-controlled phase 3 trial of a recombinant glycoprotein 120 vaccine to prevent HIV-1 infection. *Journal of Infectious Diseases*, **191**, 654–663).

5403 healthy HIV-negative volunteers at high risk of getting HIV were **randomized**, **3598** to the **vaccine** and **1805** to **placebo** (in both cases, 7 injections over 30 months), and followed for **36 months**; the **main outcome** was presence or absence of **HIV infection** at the end of the trial, with

Vaccine Efficacy (VE) defined as

$$VE = 100(1 - \text{relative risk of infection}) = 100 \left[1 - \frac{P(\text{infection}|\text{vaccine})}{P(\text{infection}|\text{placebo})} \right].$$

Secondary frequentist analyses examined **differences in VE** by **gender**, **ethnicity**, **age**, and **education** and **behavioral risk score** at baseline.

Vaccine Efficacy

Group	Rate (%) of HIV-1 Infection		VE (95% CI)	P Value		
	Vaccine	Placebo		Unadj	Adj	D-M
All	241/3598 (6.7)	127/1805 (7.0)	6 (-17 to 24)	.59	> .5	
Black (Non-Hisp)	6/233 (2.6)	9/116 (7.8)	67 (6 to 88)	.028	.24	
Black Women	1/112 (0.9)	4/57 (7.0)	87 (19 to 98)	.033		
Nonwhite	30/604 (5.0)	29/310 (9.4)	47 (12 to 68)	.012	.13	
Nonwhite Men	27/461 (6.1)	25/236 (10.6)	43 (3 to 67)	.036		

The trial found a **small decline** in infection overall (**6.7% vaccine, 7.0% placebo**) that was **neither practically nor statistically significant**; large preventive **effects** of the **vaccine** were found for some **subgroups** (e.g., **nonwhites**), but **statistical significance vanished** after adjustment for **multiple comparisons**.

Frequentist Multiple Comparisons Adjustment

Group	Rate (%) of		VE (95% CI)	<i>P</i> Value		
	HIV-1 Infection Vaccine	Placebo		Unadj	Adj	D-M
Nonwhite	30/604 (5.0)	29/310 (9.4)	47 (12 to 68)	.012	.13	

Note that the *P* value for the **nonwhite** subgroup was **0.012** before, but **0.13** after, (frequentist) multiple comparisons adjustment.

However, **frequentist multiple comparisons methods** are an **inferential approach** to what should really be a **decision problem** (Should this **vaccine** be given to **nonwhite** people at high risk of getting HIV? Should **another trial** focusing on **nonwhites** be run?), and when **multiple comparison methods** are viewed as “**solutions**” to a **Bayesian decision problem** they **do not have a sensible implied utility structure**: they’re **terrified of announcing that an effect is real when it’s not** (a **type I error**), and have **no built-in penalty for failing to announce an effect is real when it is** (a **type II error**).

Decision-Making

In the **frequentist** approach, **type II errors** are supposed to be **taken care of** by having done a **power calculation** at the time the **experiment** was **designed**, but this **begs the question** of **what decision** should be **taken**, **now that this study has been run**, about whether to **run a new trial** and/or **give the vaccine to nonwhite people now**.

When the problem is **reformulated** as a **decision** that properly **weighs all of the real-world costs and benefits**, the **result** (interpreted in **frequentist** language) would be a **third P value column** in the table on page 4 (a column called **“Implied P from a decision-making perspective”**, or **D-M** for short) that would look a lot more like the **first (unadjusted) P value column** than the **second (multiple-comparisons adjusted) column**, leading to the **decision** that a **new trial for nonwhites for this vaccine** is a **good clinical and health policy choice**.

The point is that when the **problem** is really to **make a decision**, **decision-theoretic methods** typically lead to **better choices** than **inferential methods** that were **not intended to be used** in this way.

Decision-Theoretic Re-Analysis

Group	Rate (%) of		VE (95% CI)	P Value		
	Vaccine	Placebo		Unadj	Adj	D-M
All	241/3598	127/1805	6 (-17	.59	> .5	A
Volunteers	(6.7)	(7.0)	to 24)			
Black (Non-Hisp)	6/233	9/116	67 (6	.028	.24	More Like
	(2.6)	(7.8)	to 88)			
Black Women	1/112	4/57	87 (19	.033		The
	(0.9)	(7.0)	to 98)			
Nonwhite	30/604	29/310	47 (12	.012	.13	Unadj
	(5.0)	(9.4)	to 68)			
Nonwhite Men	27/461	25/236	43 (3	.036		Col
	(6.1)	(10.6)	to 67)			

When both **type I** and **type II losses** are properly **traded off** against each other (and **gains** are correctly factored in as well), the **right choice** is (at a minimum) to **run a new trial** in which **Nonwhites** (principally **Blacks** and **Asians**, both **men and women**) are the **primary study group**.

Details

Example 3: This can be seen in an **even simpler setting**: consider a **randomized controlled Phase 3 clinical trial** with **no subgroup analysis**, and define Δ to be the **population mean health improvement** from the **treatment T** as compared with the **control condition C** .

There will typically be **some point c along the number line** (a kind of **practical significance threshold**), which may not be **0**, such that if $\Delta \geq c$ the **treatment** should be **implemented** (note that this is really a **decision problem**, with action space $a_1 = \{\text{implement } T\}$ and $a_2 = \{\text{don't}\}$).

The **frequentist hypothesis-testing inferential approach** to this problem would test $H_0: \Delta < c$ against $H_A: \Delta \geq c$, with **(reject H_0)** corresponding to action a_1 .

In the **frequentist inferential approach** H_0 would be rejected if $\hat{\Delta} \geq \Delta^*$, where $\hat{\Delta}$ is a **good estimator** of Δ based on **clinical trial data D** with **sample size n** and Δ^* is chosen so that the corresponding P value is no greater than α , the **type I error probability** (the chance of **rejecting H_0** when H_0 is **true**).

Details (continued)

As noted above, α is usually chosen to be a **conventional value** such as **0.05**, in conjunction with choosing n large enough (if you can do this at **design time**) so that the **type II error probability** β is no more than **another conventional value** such as **0.2** (the **real-world consequences of type I and type II errors** are **rarely contemplated** in choosing α and β , and in practice you won't necessarily have a **large enough** n for, e.g., **subgroup analyses** to correctly control the **type II error probability**).

The **Bayesian decision-theoretic** approach to this **decision problem** requires me to specify a **utility function** that addresses these **real-world consequences** (and others as well); a **realistic utility structure** here would depend **continuously** on Δ , but I can look at an **oversimplified utility structure** that permits **comparison with hypothesis-testing**: for $u_{ij} \geq 0$,

Action	Truth	
	$\Delta \geq c$	$\Delta < c$
a_1	u_{11}	$-u_{12}$
a_2	$-u_{21}$	u_{22}

Details (continued)

Action	Truth	
	$\Delta \geq c$	$\Delta < c$
a_1	u_{11}	$-u_{12}$
a_2	$-u_{21}$	u_{22}

The **utilities** may be considered from the point of view of several different **actors** in the drama; in the context of the **HIV vaccine study**, for instance, considering the situation from the viewpoint of a **non-HIV+ person at high risk of becoming HIV+**,

- u_{11} is the **gain** from **using** a vaccine that is **thought** to be **effective** and **really is effective**;
- $-u_{12}$ is the **loss** from **using** a vaccine that is **thought** to be **effective** and **really is not effective**;
- $-u_{21}$ is the **loss** from **not using** a vaccine that is **thought** to be **not effective** but **really is effective**; and
- u_{22} is the **gain** from **not using** a vaccine that is **thought** to be **not effective** and **really is not effective** (i.e., $u_{22} = 0$).

Details (continued)

Note that the **frequentist inferential approach** at **analysis time** only requires you to think about something (α) corresponding to **one** of these **four ingredients** ($-u_{12}$), and even then α is on the **wrong (probability) scale** (the u_{ij} will be on a **real-world-relevant scale** such as **quality-adjusted life years (QALYs)**).

The **optimal Bayesian decision** turns out to be

$$\text{choose } a_1 \text{ (implement } T) \leftrightarrow P(\Delta \geq c|D) \geq \frac{u_{12} + u_{22}}{u_{11} + u_{12} + u_{21} + u_{22}} = u^*.$$

The **frequentist inferential approach** is **equivalent** to this **only if**

$$\alpha = 1 - u^* = \frac{u_{11} + u_{21}}{u_{11} + u_{12} + u_{21} + u_{22}}.$$

In the context of the **HIV vaccine**, with realistic values of the u_{ij} that **appropriately weigh** both the **loss from taking the vaccine when it doesn't work** and **failing to take the vaccine when it does work**, the analogous **frequentist inferential “action”** would be to **reject H_0** for P **values** that are **much larger** than the usual threshold

(e.g., **0.3** instead of **0.05**).

VA Quality of Care Example (continued)

Need to compare **observed** and **expected** compliance rates for each provider, where expected rate based on **relevant patient-level covariates**.

We use a **method** for making this **comparison** developed in **Draper and Gittoes** (2004: Statistical analysis of performance indicators in UK higher education (with discussion). *Journal of the Royal Statistical Society, Series A*, **167**, 449–474).

In this method all **covariates** are treated as **categorical**; in the VA case study we have **age** (7 categories); **race** (3); **number of medicines taken** (4); and **relative risk score** [overall illness burden] (6), for a total of **504** possible cells, of which **457** were non-empty.

The method produces a **standardized z score** for each provider, who is then **provisionally classified** as providing **H** (unusually **good** care), **N** (**normal** care) or **L** (unusually **bad** care); how choose **cutpoints** on z scale?

Ad-hoc **inferential** approach: use conventional **sensitivity** and **specificity** goals to **implicitly trade off false positive and false negative errors**.

VA Quality of Care Example (continued)

Bayesian decision-theory approach: elicit 3×3 **utility matrix** (cross-tabulating **classification** against **truth**) from client and choose cutpoints to **maximize expected utility**; do **sensitivity analysis on utility matrix** to assess **stability of findings**; current **client utility matrix** (H (high) = unusually good care, N (normal) = normal care, L (low) = unusually bad care):

		Truth		
		H	N	L
What Our Method Says	H	+2	-1	-6
	N	-1	0	-3
	L	-5	-2	+2

Performing **MEU calculation** requires **truth** to compare with **provisional classification**; to obtain this, we created **simulation environment** closely matching **VA client network** except that **we control truth**; this involves **provider-level heterogeneity parameter** σ : $\sigma = 0$ means **no H, no L, all N**; as σ increases, get more **truly good** and **truly bad** simulated providers.

VA Quality of Care Example (continued)

Preliminary results: $\sigma = 0.3$ is closest to client reality; for this σ , optimal cut-points on z scale are at ± 3.5 (note that Bayesian decision theory has led to a natural, and rather strong, adjustment for multiplicity); simulated success rates (%) with $\sigma = 0.3$ and $z = \pm 3.5$ cutpoints:

		Truth		
		H	N	L
What Our Method Says	H	0.01	0.03	0.00
	N	4.50	92.44	3.02
	L	0.00	0.00	0.00

Comments: (1) So far have only explored $+z$ and $-z$ cutoffs **equal in absolute value**; need to explore **asymmetric cutoffs**; (2) client notes that method at present (with this utility matrix) is **somewhat conservative** (its mistakes are almost all in saying N when truth was H or L); this suggests utility matrix may need more **refining**, e.g., to give **more positive reward** to the (H, H) and (L, L) cells.