

Value of Information: Sensitivity Analysis and Research Design in Bayesian Evidence Synthesis

Christopher Jackson, Anne Presanis, Stefano Conti, Daniela De Angelis*
MRC Biostatistics Unit, University of Cambridge; NHS England

November 20, 2018

Abstract

Suppose we have a Bayesian model that combines evidence from several different sources. We want to know which model parameters most affect the estimate or decision from the model, or which of the parameter uncertainties drive the decision uncertainty. Furthermore we want to prioritise what further data should be collected. These questions can be addressed by Value of Information (VoI) analysis, in which we estimate expected reductions in loss from learning specific parameters or collecting data of a given design. We describe the theory and practice of VoI for Bayesian evidence synthesis, using and extending ideas from health economics, computer modelling and Bayesian design. The methods are general to a range of decision problems including point estimation and choices between discrete actions. We apply them to a model for estimating prevalence of HIV infection, combining indirect information from surveys, registers and expert beliefs. This analysis shows which parameters contribute most of the uncertainty about each prevalence estimate, and the expected improvements in precision from specific amounts of additional data. These benefits can be traded with the costs of sampling to determine an optimal sample size.

Keywords: decision theory, research prioritisation, uncertainty

*This work was funded by the Medical Research Council, grant code U105260566, and from Public Health England (funding DDA and SC). The authors are grateful to the HIV department of Public Health England for providing the data and permission to use the example, to the NATSAL team for providing data, and Louise Logan for advice on costs of data collection in GUM Anon.

1 Introduction

Bayesian modelling is a natural paradigm for decision making, in the presence of uncertainty, based on multiple sources of evidence. However, as more data sources, parameters and assumptions are built into a model, it becomes harder to see the influence of each input or assumption. The modelling process should involve an investigation of where the weak parts of the model are, to identify which uncertainties in the model inputs contribute most to the uncertainty in the final result or decision (*sensitivity analysis*). We might then want to assess and compare the potential value of obtaining datasets of specific designs or sizes to strengthen different parts of the model. Furthermore, we may want to formally trade off the costs of sampling with the resulting expected improvement to decision making.

Annual estimation of HIV prevalence in the United Kingdom has, for several years, been based on a Bayesian synthesis of evidence from various surveillance systems and other surveys (Goubar et al. 2008, Presanis et al. 2010, De Angelis et al. 2014, Kirwan et al. 2016). This is an example of a class of problems called *multiparameter evidence synthesis* (MPES) (e.g. Ades & Sutton 2006), where the quantities of interest cannot be estimated directly, but can be inferred from multiple indirect data sources linked through a network of model assumptions that can be expressed as a directed acyclic graph. Markov Chain Monte Carlo (MCMC) is typically required to estimate the posterior. The HIV MPES model is used to inform health policies, thus it is crucial to be able to assess sensitivity to uncertain inputs and to indicate how the model could be strengthened with further data.

These dual aims can be achieved with *value of information* (VoI) analysis, a decision-theoretic framework based on expected reductions in loss from future information. The concepts of VoI were first set out in detail by Raiffa & Schlaifer (1961), while Parmigiani & Inoue (2009) give a more recent overview. The expected value of *partial perfect information*

(EVPPI) is the expected reduction in loss if the exact value of a particular parameter or parameters θ_0 were learnt, also interpreted as the amount of decision uncertainty that is due to θ_0 . The expected value of *sample information* (EVSI) is the expected reduction in loss from a study of a specific design. The EVSI can be traded off with the costs of data collection to give the *expected net benefit of sampling* (ENBS). Therefore, as well as recommending a policy based on minimising expected loss under the *current* model and data, the decision-maker may also recommend collecting *further* data according to a design which minimises the ENBS.

These concepts have been applied in various forms in three distinct areas: health economics, computer modelling and Bayesian design. In health economic modelling, there is a large literature on calculation and application of VoI, see, e.g. Felli & Hazen (1998), Willan & Pinto (2005), Claxton & Sculpher (2006), Welton et al. (2008). The model output is then the expected net benefit of each alternative policy, a known deterministic function $g(\theta)$ of uncertain inputs θ , and the decision problem is the choice of policy that minimises $g(\theta)$. In computer modelling, see, e.g. Oakley & O’Hagan (2004) and Saltelli et al. (2004), the influence of a particular element θ_0 of θ is calculated as the expected reduction in $\text{var}(g(\theta))$, if we were to learn θ_0 exactly. This is equivalent to the EVPPI for θ_0 under a decision problem defined as point estimation of $g(\theta)$ with quadratic loss (Oakley & O’Hagan 2004). The decision-theoretic view of Bayesian experimental design also has a long history, see, e.g. Lindley (1956), Bernardo & Smith (1994), Chaloner & Verdinelli (1995), Berger (2013), and a recent review of the computational challenges by Ryan et al. (2016).

However, the current tools in any one of these three areas cannot be applied directly to MPES. First of all, it is not always feasible or desirable to make a discrete decision with a quantifiable loss, as in health economic modelling. Instead, the aim of evidence synthesis is often to estimate one or more quantities. For a scalar quantity of interest, we might then

define the “loss” as the posterior variance of this quantity, as Oakley & O’Hagan (2004) described in the computer modelling context. In computer modelling, however, tools to estimate the expected value of a proposed study to learn a particular θ_0 more precisely have not been developed, and it is not clear what an appropriate loss for a vector of model outputs would be. Challenges also arise with computation. Current methods for computing the expected variance reduction in the computer modelling field (Sobol’ 2001, Saltelli et al. 2004) assume the output is an explicit function $g(\theta)$ of the inputs, therefore do not apply in MPES, where this function is unknown and the outputs must be estimated by MCMC. For Bayesian design, Ryan et al. (2016) reviewed methods where evaluating the expected utility of a design (equivalent to the EVSI) is relatively inexpensive, so that maximising the utility over a complex design space is feasible. However, this can again be difficult with MCMC. Given a sample from the posterior $p(\theta|\mathbf{x})$, potential future datasets \mathbf{y} under a specific design can be simulated cheaply from the posterior predictive distribution, but then to obtain the expected utility, the posterior $p(\theta|\mathbf{x}, \mathbf{y})$ needs to be repeatedly updated for different \mathbf{y} , which is feasible with Monte Carlo only for smaller problems (e.g. Han & Chaloner 2004).

Here we use and extend methods from health economics, computer modelling and Bayesian design to devise a new VoI framework for sensitivity analysis and research design in evidence syntheses based on graphical models fitted by MCMC. This is a broader class of models than those typically used in health economics or computer modelling, since the model “output” is not necessarily a known function of the inputs, but depends on the model parameters θ and observed data \mathbf{x} through a network of statistical models or deterministic functions, potentially with hierarchical relationships. We apply this new VoI framework to the part of the HIV prevalence estimation model that estimates prevalence in men who have sex with men (MSM), in London. Here the decision problem is point estimation of a

single scalar or a vector of parameters, followed by the choice of what extra data should be collected in the future. We use ideas from Bayesian design to choose appropriate loss functions in this context. We also generalize methods of computing EVPPI (Strong et al. 2014) and EVSI (Strong et al. 2015), developed for finite choices in health economics, to a broader class of decision problems, including point estimation. The method for computing EVSI enables the expected utility over all potential \mathbf{y} to be estimated cheaply without an additional level of simulation, assuming only that the information provided by \mathbf{y} can be represented as a low-dimensional sufficient statistic $T(\mathbf{y})$.

In § 2 we describe the general MPES model, and define the expected value of information under different decision problems and loss functions, and in § 3 we present methods to compute them. In § 4 we describe the model for HIV prevalence estimation, and in § 5 we use VoI to identify the areas of greatest uncertainty in this model and determine what specific data should be collected to improve the precision of the estimates of various subgroup-specific prevalences. Finally we discuss potential extensions to the methods and application and the associated challenges.

2 Theory and methods

2.1 Bayesian graphical modelling for evidence synthesis

In our motivating applications, the general model can be represented as a directed acyclic graph (Figure 1) in the standard way, see, e.g. Lauritzen (1996). Nodes in the graph may represent scalar or vector quantities. A set of datasets $\mathbf{x} = \{x_1, \dots, x_n\}$ is observed, most generally from n different sources. These data are assumed to arise from statistical models with parameters μ_1, \dots, μ_n respectively, collectively denoted $\boldsymbol{\mu}$. The “founder nodes” of the

graph are denoted $\boldsymbol{\phi} = (\phi_1, \dots, \phi_p)$ and given a joint prior distribution $\boldsymbol{\phi} \sim p(\cdot)$ which may also include substantive information. The full set of unknowns is denoted $\boldsymbol{\theta}$. Most simply, the $\boldsymbol{\mu}$ could equal the $\boldsymbol{\phi}$ or be related to the $\boldsymbol{\phi}$ through deterministic functions, so that $\boldsymbol{\theta} = \boldsymbol{\phi}$. More generally, some of the relationships in the graph could be stochastic, defining a hierarchical model, where the $\boldsymbol{\mu}$ themselves arise from a distribution with parameters given by the $\boldsymbol{\phi}$ or descendants of $\boldsymbol{\phi}$. The vector of unknowns $\boldsymbol{\theta}$ would then comprise $\boldsymbol{\phi}$ and the stochastic descendants of $\boldsymbol{\phi}$ such as random effects.

We further denote $\boldsymbol{\alpha}$ as an intermediate node in the graph, the model “output”, which is used for decision-making. This could be any unknown quantity, including one of the $\boldsymbol{\mu}$ or $\boldsymbol{\phi}$, a function of these, or a prediction of new data. We may also plan to collect additional data, either from the same source as one of the existing datasets (e.g. y_1 in Figure 1), or from a new source informing a parameter μ_{n+1} on which no direct data (y_2) were available.

This DAG (Figure 1) is a generalisation of the typical structure (Figure 2) used in computer modelling (Oakley & O’Hagan 2004) where the output $\boldsymbol{\alpha}$ is a known (usually complicated) deterministic function of uncertain model inputs $\boldsymbol{\phi}$, which are given substantive priors that may be derived separately from data.

2.2 Expected value of information: definitions

In a general decision-theoretic framework, the purpose of the model is to choose a decision or action d from a space of possible decisions \mathcal{D} , to minimise an expected loss $E_{\boldsymbol{\theta}}(L(d, \boldsymbol{\theta}))$, with the expectation taken with respect to the posterior distribution of $\boldsymbol{\theta}$. Let $\boldsymbol{\alpha} = \boldsymbol{\alpha}(\boldsymbol{\theta})$ be the minimal subset or function of $\boldsymbol{\theta}$ necessary to make the decision, so that $E_{\boldsymbol{\theta}}(L(d, \boldsymbol{\theta})) = E_{\boldsymbol{\alpha}}(L(d, \boldsymbol{\theta}))$, $\forall d \in \mathcal{D}$. For example, the purpose could be the choice of decision d among a finite set $\mathcal{D} = \{1, \dots, D\}$ expected to minimise a loss defined as a function of the parameters, so that $\boldsymbol{\alpha}$ would be a vector with D components $\alpha_d = f_d(\boldsymbol{\theta}) = L(d, \boldsymbol{\theta})$.

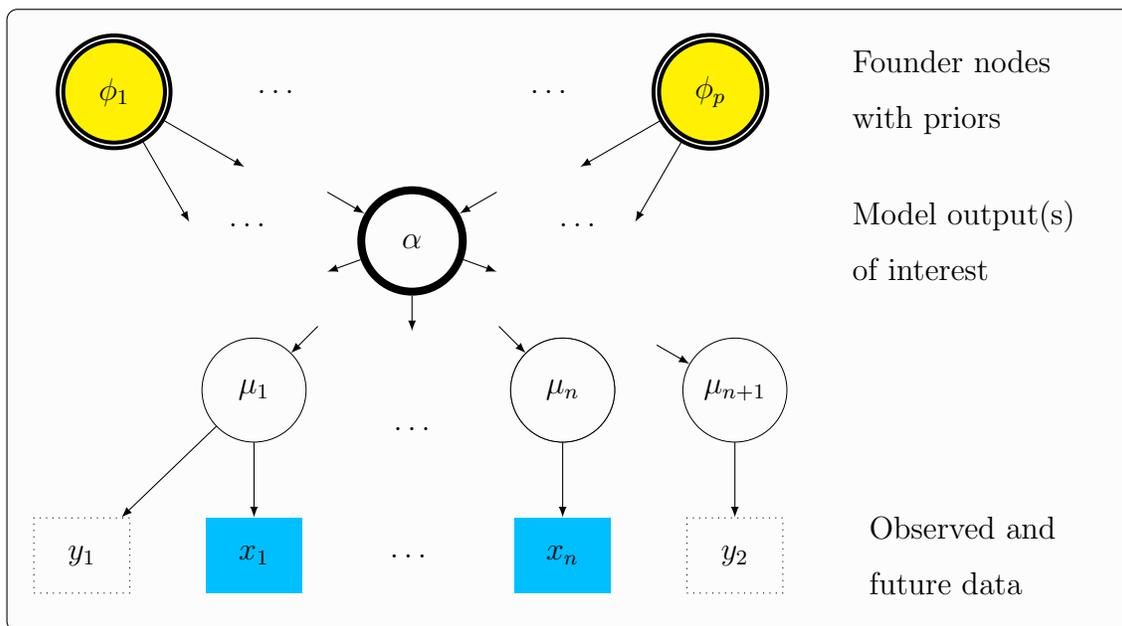


Figure 1: Directed acyclic graph for Bayesian evidence synthesis

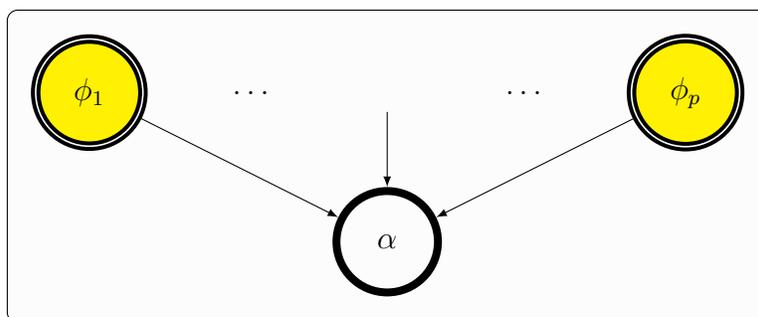


Figure 2: Graph representing a known deterministic model

This is the typical situation in health policy decisions (e.g. Claxton & Sculpher 2006), where a treatment d is chosen to maximise a measure of utility such as expected quality-adjusted survival. Alternatively, as in our examples, the decision could simply be the choice of a point estimate $\hat{\alpha}$ of some parameter α , in which case the decision space \mathcal{D} is the support of α (see §2.3). Alongside making a decision, we wish to also determine where further research should be prioritised to reduce the uncertainty about the decision, and given the costs of data collection, to determine the optimal design of further research (see §2.4).

For general decision problems, let $d^* = \arg \min_d E_{\theta}(L(d, \theta))$ be the optimal decision under *current* knowledge about θ , represented by the posterior distribution $p(\theta|\mathbf{x})$. Suppose now we are in a position to collect new information. Let d_y^* be the optimal decision given further knowledge of a quantity \mathbf{y} (either parameters or potential data) that informs α , so that the updated posterior would be $p(\theta|\mathbf{x}, \mathbf{y})$. We define the following quantities.

1. The *expected value of perfect information* (EVPI) is the expected loss of the decision d^* under current information, minus the expected loss for the decision d_{α}^* we would make if we knew the true α (Raiffa & Schlaifer 1961).

$$EVPI = E_{\theta}(L(d^*, \theta)) - E_{\theta}(L(d_{\alpha}^*, \theta))$$

Since additional information is always expected to reduce the expected loss of the optimal decision (Parmigiani & Inoue 2009), the EVPI is an upper bound on the expected gains from any new information.

2. The *expected value of partial perfect information* (EVPPI) for a particular parameter ϕ is the expected reduction in loss if ϕ were to be learnt precisely. Since this precise value is not yet known, an expectation must be taken over all possible values.

$$EVPPI(\phi) = E_{\theta}(L(d^*, \theta)) - E_{\phi}[E_{\theta|\phi}(L(d_{\phi}^*, \theta))] \quad (1)$$

where d_ϕ^* is the optimal decision if ϕ were known. This is an upper bound on the potential value of data \mathbf{y} which inform only ϕ . In a graphical model, this means data \mathbf{y} that are conditionally independent of θ given ϕ , for example $\mathbf{y} = y_1$ and $\phi = \mu_1$ in Figure 1.

3. The *expected value of sample information* $EVSI(\mathbf{y})$ is the reduction in loss we would expect from collecting an additional dataset \mathbf{y} of a specific design.

$$EVSI(\mathbf{y}) = E_\theta(L(d^*, \theta)) - E_{\mathbf{y}} [E_{\theta|\mathbf{y}}(L(d_\mathbf{y}^*, \theta))] \quad (2)$$

The inner expectation is now with respect to the updated posterior distribution of $\theta|\mathbf{y}$, after learning \mathbf{y} as well as the existing data \mathbf{x} , or “preposterior” (Berger 2013).

2.3 Value of information in different decision problems

Finite-action decisions For a choice of d among a finite set $\{1, \dots, D\}$ with loss $L(d, \theta) = \alpha_d$ and $\alpha = \{\alpha_1, \dots, \alpha_d, \dots, \alpha_D\}$, the expected loss with current information is $\min_d \{E_\alpha(\alpha_d)\}$, so (Raiffa & Schlaifer 1961)

$$\begin{aligned} EVPI &= \min_d \{E_\alpha(\alpha_d)\} - E_\alpha \min_d \{\alpha_d\} \\ EVPPI(\phi) &= \min_d \{E_\alpha(\alpha_d)\} - E_\phi \min_d \{E_{\theta|\phi}(\alpha_d)\} \\ EVSI(\mathbf{y}) &= \min_d \{E_\alpha(\alpha_d)\} - E_{\mathbf{y}} \min_d \{E_{\alpha|\mathbf{y}}(\alpha_d)\} \end{aligned} \quad (3)$$

Point estimation When the decision is the choice of a point estimate $\hat{\alpha}$ of a vector of parameters α , with quadratic loss

$$L(\hat{\alpha}, \alpha) = (\hat{\alpha} - \alpha)^T H (\hat{\alpha} - \alpha) \quad (4)$$

for a symmetric, positive-definite H , the optimal estimate with current information is the posterior mean, $\hat{\alpha} = E_\alpha(\alpha)$.

For a scalar $\boldsymbol{\alpha} = \alpha$ and $H = 1$, the expected loss is $\text{var}(\alpha)$ under current information and zero under perfect information, so that $EVPI = \text{var}(\alpha)$ and

$$EVPPI(\boldsymbol{\phi}) = \text{var}(\alpha) - E_{\boldsymbol{\phi}} [\text{var}_{\alpha|\boldsymbol{\phi}}(\alpha|\boldsymbol{\phi})] \quad (5)$$

$$EVSI(\mathbf{y}) = \text{var}(\alpha) - E_{\mathbf{y}} [\text{var}_{\alpha|\mathbf{y}}(\alpha|\mathbf{y})] \quad (6)$$

the expected reduction in variance given new information. Expression (5) is used by Oakley & O'Hagan (2004) and Saltelli et al. (2004) as a measure of sensitivity of the output of a deterministic model $\alpha = g(\phi, \dots)$ to an uncertain input ϕ , termed the *main effect* of ϕ , but this has not been extended to the EVSI of potential data \mathbf{y} in a point estimation context.

When $\boldsymbol{\alpha}$ is a vector, the typical situation where a MPES of the form in Figure 1 is carried out, we could conduct independent value of information analyses for each component of $\boldsymbol{\alpha}$. In more formal decision analyses we may want a scalar loss for the overall vector $\boldsymbol{\alpha}$. There are various alternatives based on generalisations $v(\boldsymbol{\alpha})$ of the variance, which can be used instead of the scalar variance $\text{var}(\alpha)$ in equations (5)–(6) to define the expected value of information. These have been applied in the context of Bayesian study design, and we explain two examples that can be adopted for EVPPI and EVSI in our context as follows.

1. If $H = \mathbf{c}\mathbf{c}^T$ in the quadratic loss (4), for some vector of weights \mathbf{c} , then the expected loss is $v(\boldsymbol{\alpha}) = \mathbf{c}^T \text{cov}(\boldsymbol{\alpha}) \mathbf{c} = \text{var}(\mathbf{c}^T \boldsymbol{\alpha})$, corresponding to optimal (under squared error loss) estimation of the weighted sum of the parameters, $\mathbf{c}^T \boldsymbol{\alpha}$. For example, when the elements α_s of $\boldsymbol{\alpha}$ are weighted equally, the goal is to minimise the sum of all elements (r, s) of the covariance matrix, $v(\boldsymbol{\alpha}) = \sum_{r,s} \text{cov}(\boldsymbol{\alpha})_{r,s}$, or, if the α_s are also independent of each other, $v(\boldsymbol{\alpha}) = \text{tr}(\text{cov}(\boldsymbol{\alpha})) = \sum_s \text{var}(\alpha_s)$. The same *absolute* reductions in variance for different components of $\boldsymbol{\alpha}$ would then be valued equally. More generally, if \mathbf{c} is given a prior, then loss (4) also arises (see Chaloner & Verdinelli (1995) and references therein). Designs that minimise (4) are Bayesian analogues of

classical *A-optimal* designs. See also Lamboni et al. (2011) for similar measures of sensitivity for multivariate outputs in deterministic computer models.

2. A Bayesian *D-optimal* design, on the other hand, minimises the *determinant* $v(\boldsymbol{\alpha}) = \det(\text{cov}(\boldsymbol{\alpha}))$ (Chaloner & Verdinelli 1995, Ryan et al. 2016). This simplifies to the *product* of the $\text{var}(\alpha_s)$ when the α_s are independent and equally-weighted. Equivalently, a standardised version $\det(\text{cov}(\boldsymbol{\alpha}))^{1/S}$, where S is the number of components of $\boldsymbol{\alpha}$, represents a geometric average variance of the α_s , adjusted for their covariance. Here the same *relative* reductions in variance for different components of $\boldsymbol{\alpha}$ would then be valued equally, which would be more appropriate when the output of interest $\boldsymbol{\alpha}$ comprises quantities on very different scales and/or with different interpretations.

2.4 Maximising the expected net benefit of sampling

The EVSI measures the expected *benefits* from sampling. The *costs* of sampling should also be considered. The decision-maker can then choose the design and sample size for data \mathbf{y} to maximise the *expected net benefit of sampling* $E_{\mathbf{y}}(b(\mathbf{y}) - c(\mathbf{y}))$, where $b(\mathbf{y}) = EVSI(\mathbf{y})$ is the benefit and $c(\mathbf{y})$ is the cost of obtaining data \mathbf{y} (Parmigiani & Inoue 2009). This requires benefits and costs to be measured on the same scale, which can be achieved in different ways. Improved precision of point estimates might be valued in monetary terms, as described below and illustrated in §5.4. Alternatively, the better knowledge given by the new data could lead to indirect benefits which could be valued, for example, improved health from better-informed health-related decision making, as discussed in §6. We will assume $c(\mathbf{y})$ depends only on the design and sample size, thus is known in advance of observing \mathbf{y} , so that $E_{\mathbf{y}}(c(\mathbf{y})) = c(\mathbf{y})$.

To directly translate improved precision to a monetary benefit, the decision-maker

should specify the amount they are willing to pay to reduce the posterior variance by a certain amount. This willingness to pay may depend on the original posterior variance. Formally, the benefit function $b(\mathbf{y}) = f(v_0, v_{\mathbf{y}})$, specified by the decision-maker, places a value on a reduction in variance (or its multivariate analogue as in §2.3) from $v_0 = \text{var}(\boldsymbol{\alpha})$ to $v_{\mathbf{y}}$, the variance after collecting new data \mathbf{y} . For example, if any *absolute* variance reduction is valued the same way (as in A-optimal design, see §2.3), $f(v_0, v_{\mathbf{y}}) = \lambda(v_0 - v_{\mathbf{y}})$, where λ is the constant willingness-to-pay for one unit of variance reduction. The expected benefit is then $E_{\mathbf{y}}(b(\mathbf{y})) = \lambda(v_0 - E_{\mathbf{y}}(\text{var}_{\boldsymbol{\alpha}|\mathbf{y}}(\boldsymbol{\alpha}|\mathbf{y})))$, which equals $EVSI(\mathbf{y})$ using the quadratic loss function (4) multiplied by a constant λ . Alternatively, if the same *relative* gains are valued equally (as in D-optimal design), the decision-maker could specify λ as the amount they are willing to pay to (e.g.) halve the variance, so that $f(v_0, v_{\mathbf{y}}) = f(v_0, v_0/2^k) = k\lambda$, for $k = \log(v_0/v_{\mathbf{y}})/\log(2)$.

3 Computation of value of information

3.1 Partial perfect information

Computation of the EVPPI in general is not straightforward. Given a sample from the posterior distribution, the first term in (1) can be calculated by a Monte Carlo mean. The double expectation in the second term is more challenging. While it can be evaluated using nested Monte Carlo, this is expensive. Strong et al. (2014) proposed a method for estimating the EVPPI in the special case of finite choice decisions (equation 3) which uses only a single Monte Carlo loop. To estimate EVPPI (1) in a broader class of decision problems, which also includes point estimation, the method needs to be generalised.

Strong et al. (2014) estimated formula (3) by expressing

$$\alpha_d = E_{\alpha_d|\phi}(\alpha_d|\phi) + \epsilon = g_d(\phi) + \epsilon \quad (7)$$

for each $d = 1, \dots, D$, where ϵ is an error term with mean zero. Then $g_d(\phi)$ is estimated by regression of α_d on ϕ , fitted to a Monte Carlo sample of $(\alpha_d^{(k)}, \phi^{(k)}) : k = 1, \dots, K$. If ϕ comprises p parameters that could be learnt simultaneously, the regression will have p predictors. Since the functional form of $g_d()$ will not be known in general, nonparametric regression methods are used. This produces a fitted value $\hat{g}_d(\phi^{(k)})$ for each k , which allows the second term in (3) to be estimated by a Monte Carlo mean

$$E_\phi[\min_d(E_{\theta|\phi}(\alpha_d))] = E_\phi[\min_d(E_{\alpha_d|\phi}(\alpha_d|\phi))] \approx \frac{1}{K} \sum_{k=1}^K \min_d(\hat{g}_d(\phi^{(k)})). \quad (8)$$

Our generalisation of this approach computes EVPPI (equation 1) in a broader class of decision problems defined as follows. Given a state of knowledge about the decision-relevant quantities α represented by a distribution $\psi(\cdot)$, the expected loss under the optimal decision should be a known function h of the mean of α under that distribution:

$$E_\psi(L(d_\psi^*, \theta)) = h(E_\psi(\alpha)). \quad (9)$$

If $\psi(\cdot)$ is the current posterior, this is $h(E_\alpha(\alpha))$, and if we were to learn the value of ϕ , the expected loss would be $h(E_{\alpha|\phi}(\alpha|\phi))$. The method of Strong et al. (2014) only applies to the special case where α is a vector and $h(E(\alpha)) = \min_d\{E(\alpha_d)\}$. To estimate $E_{\alpha|\phi}(\alpha|\phi)$ in more general problems, we use a similar principle to (7–8), by expressing

$$\alpha = E_{\alpha|\phi}(\alpha|\phi) + \epsilon = g(\phi) + \epsilon \quad (10)$$

then fitting a regression model $g()$ of α on ϕ allows us to estimate

$$E_\phi[E_{\theta|\phi}(L(d_\phi^*, \theta))] = E_\phi[h(E_{\alpha|\phi}(\alpha|\phi))] \approx \frac{1}{K} \sum_{k=1}^K h(\hat{g}(\phi^{(k)})).$$

Point estimation problems are also a special case of (9), for example, for estimation of a scalar α with quadratic loss, $h(E_\alpha(\alpha)) = E[(\alpha - E_\alpha(\alpha))^2] = \text{var}(\alpha)$. Therefore to calculate EVPPI in this case (equation 5), we estimate $\text{var}(\alpha|\phi^{(k)})$ by the squared residual $(\alpha - \hat{g}(\phi^{(k)}))^2$, substitute this for $h(\hat{g}(\phi^{(k)}))$ and estimate $E_\phi [\text{var}_{\alpha|\phi}(\alpha|\phi)]$ as the mean, over k , of the squared residuals. Equivalently we can estimate $\text{var}(\theta) - E_\phi [\text{var}_{\alpha|\phi}(\alpha|\phi)] = \text{var}_\phi(E_{\alpha|\phi}(\alpha|\phi))$ as the variance, over k , of the fitted values. Similarly, for vector $\boldsymbol{\alpha}$ and loss functions based on $\text{cov}(\boldsymbol{\alpha})$, we can fit regressions to get the marginal mean for each component α_d , and calculate the empirical covariance matrix of the residuals.

Several methods of nonparametric regression have been suggested. For small p , Strong et al. (2014) used generalized additive models, with tensor products of splines to represent interactions between components of $\boldsymbol{\phi}$. Where $\boldsymbol{\phi}$ included about $p = 5$ or more components, Gaussian process regression was recommended as a more efficient way of modelling interactions, though the resulting matrix computations rapidly become impractical as the MCMC sample size K increases. Heath et al. (2016) developed an integrated nested Laplace approximation for fitting Gaussian processes more efficiently where $p \geq 2$. For the application in § 4 (with $K = 150000$, $p \leq 3$), we have found multivariate adaptive regression splines (Friedman 1991) via the *earth* R package (Milborrow 2011) to be more efficient. Standard errors for the EVPPI estimates can be calculated in general by simulating from the asymptotic normal distribution of the regression coefficients (Mandel 2013).

3.2 Sample information

The regression method above can also be used to estimate the expected value of sample information $EVSI(\mathbf{y})$. This again requires a generalisation of the approach described by Strong et al. (2015) from finite decision problems to any problem satisfying condition (9), including point estimation. The method requires that the information provided by the

data \mathbf{y} can be expressed as a low-dimensional sufficient statistic $T(\mathbf{y})$, so that $E_{\alpha|\mathbf{y}}(\alpha|\mathbf{y}) = E_{\alpha|\mathbf{y}}(\alpha|T(\mathbf{y}))$. This could be a point estimator of the parameter μ (as in Figure 1) that \mathbf{y} gives direct information on. As in (10), we can write

$$\alpha = E_{\alpha|\mathbf{y}}(\alpha|T(\mathbf{y})) + \epsilon = g(T(\mathbf{y})) + \epsilon$$

and estimate $g(\cdot)$ using a regression fitted to a Monte Carlo sample of $(\alpha^{(k)}, T(\mathbf{y}^{(k)})) : k = 1, \dots, K$, where $\mathbf{y}^{(k)}$ are drawn from their posterior predictive distribution. Then the fitted values $\hat{g}(T(\mathbf{y}^{(k)}))$ enable the double expectation to be estimated as

$$E_{\mathbf{y}}[E_{\theta|\mathbf{y}}(L(d_{\mathbf{y}}^*, \boldsymbol{\theta}))] = E_{\mathbf{y}}[h(E_{\alpha|\mathbf{y}}(\alpha|\mathbf{y}))] \approx \frac{1}{K} \sum_{k=1}^K h(\hat{g}(T(\mathbf{y}^{(k)}))).$$

Then, for example, for point estimation with quadratic loss, this is the estimated residual variance from the regression, as in § 3.1.

4 The HIV prevalence MPES model

We consider the sub-model of the full HIV burden model (De Angelis et al. 2014, Kirwan et al. 2016) that estimates HIV prevalence in men who have sex with men (MSM), in London. We define three subgroups of MSM: those who have attended a genitourinary medicine (GUM) clinic in the past year (GMSM), those who have not (NGMSM), and previous MSM (PMSM), men who no longer have sex with men. We denote the proportion of all men who are in these subgroups by ρ_G , ρ_N and ρ_P respectively. For each group $g \in (G, N, P)$, we aim to estimate simultaneously these subgroup proportions ρ_g , prevalence of HIV in this group π_g and the proportion of infections that are diagnosed, δ_g . Given these parameters, further important quantities are easily derived: the prevalence of diagnosed ($\pi_g \delta_g = (\pi \delta)_g$) and undiagnosed ($\pi_g(1 - \delta_g) = \overline{(\pi \delta)}_g$) infection; and the numbers of MSM living with diagnosed ($\mu_{Dg} = \mu_{pop} \rho_g (\pi \delta)_g$) and undiagnosed ($\mu_{Ug} = \mu_{pop} \rho_g \overline{(\pi \delta)}_g$) infection, where μ_{pop} is

the number of men (MSM and non-MSM) living in London. Since the prevalence among PMSM is much lower, this subgroup is not examined in detail.

We construct a Bayesian model to link the unknown ρ_g, π_g, δ_g with the available evidence provided by various routinely-collected and survey datasets as well as expert belief. Figure 3 shows a directed acyclic graph representing this model, in the form of Figure 1, distinguishing founder nodes, observed data, and outputs of interest. The following sections explain in detail the quantities and relationships illustrated in Figure 3. All data and estimates refer to the year 2012 (unless indicated) and the Greater London area.

4.1 Subgroup membership

The total male population of London, μ_{pop} , is informed by published data y_{pop} (Office for National Statistics 2012), assumed to be a Poisson count: $y_{pop} \sim Po(\mu_{pop})$. A log-normal prior for μ_{pop} is assumed, $\log(\mu_{pop}) \sim N(0, 1000^2)$. The number of people in each group g is estimated as $r_g = \rho_g \mu_{pop}$. Estimates of the subgroup proportions ρ_g are informed by data from the National Survey of Sexual Attitudes and Lifestyles (Mercer et al. 2013): $y_G = 7$, $y_N = 38$, $y_P = 10$, out of $y_{NAT} = 824$ men, which we assume to come from a multinomial distribution with probabilities ρ_G, ρ_N, ρ_P given a uniform Dirichlet prior. Thus the expected number of people with HIV (diagnosed or undiagnosed) in group g is $\mu_g = \pi_g r_g$.

4.2 Registry of diagnosed infections and diagnosed prevalence

Individuals diagnosed with HIV and accessing care in the UK are reported to the HIV and AIDS Reporting System (Kirwan et al. 2016). From the 2012 version of this dataset, known as SOPHID (Surveillance of Prevalent HIV Infections Diagnosed), we obtain the reported number of HIV diagnoses for MSM, $y_M \sim Po(\mu_M)$, with $y_M = 8390$. A reporting bias of un-

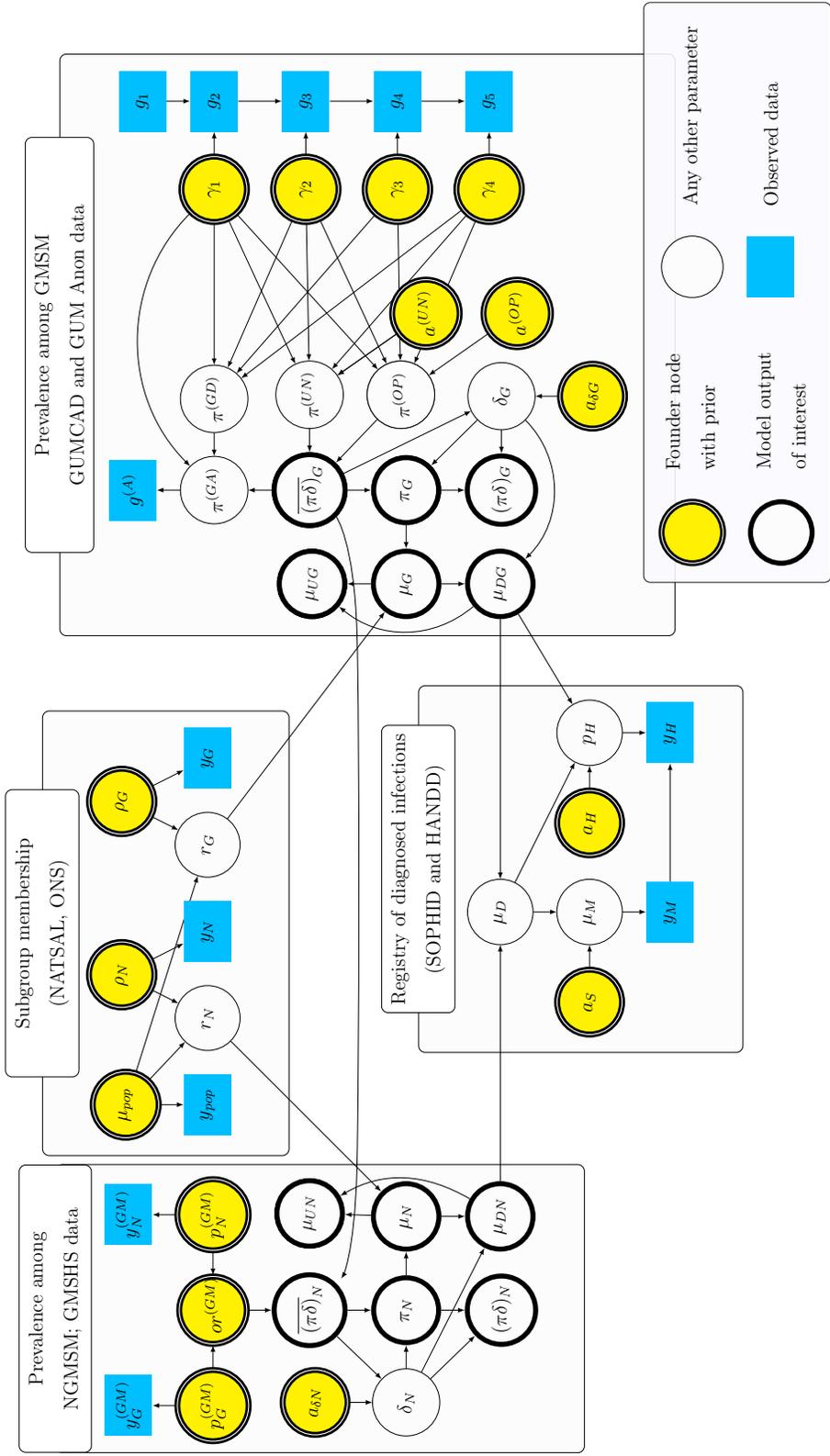


Figure 3: Directed acyclic graph for HIV prevalence estimation model.

known direction is assumed, through $\log(\mu_M) = a_S + \log(\mu_D)$ where $\exp(a_S) \sim N(1, 0.018^2)$, giving a prior 90% interval of about $(-3\%, 3\%)$ for the adjustment to the number of MSM HIV diagnoses μ_M . After adjustment, $\mu_D = \mu_{DG} + \mu_{DN} + \mu_{DP}$ is the expected number of diagnoses among MSM, summed from the expected numbers of diagnoses among GMSM, NGMSM and PMSM respectively. The following sections explain where μ_{DG}, μ_{DN} come from; μ_{DP} is modelled using similar techniques.

Since SOPHID did not record GUM clinic attendance, to strengthen the evidence on diagnosed prevalence in GMSM we include data from the HIV and AIDS New Diagnoses and Deaths Database (HANDD) (Kirwan et al. 2016), recording how many of the y_M prevalent diagnosed MSM were newly diagnosed in 2012 and reported to have been diagnosed initially in a GUM clinic. These new diagnoses, $y_H = 630$, are modelled as $y_H \sim \text{Bin}(y_M, p_H)$, where p_H is assumed to be a lower bound for the proportion of prevalent diagnosed MSM who have attended a GUM clinic in 2012. This bound is expressed through $p_H = a_H \mu_{DG} / \mu_D$, where $a_H \sim U(0, 1)$ is the unknown probability that a prevalent diagnosed MSM who has attended a GUM clinic in 2012 was newly diagnosed that year. y_H therefore gives us additional indirect information on μ_{DG} , the number of prevalent diagnosed GMSM.

The number of *diagnosed* infections is related to the total number of infections in each group g as $\mu_{Dg} = \delta_g \mu_g$. The proportion of infections that are diagnosed δ_g is not known, but given our inferences about the undiagnosed prevalence $\overline{(\pi\delta)}_g = \pi_g(1 - \delta_g)$ (explained in the subsequent sections), we can exploit the implicit constraint $1 - \delta_g > \overline{(\pi\delta)}_g$. Therefore we define $\delta_g = a_{\delta g}(1 - \overline{(\pi\delta)}_g)$, with $a_{\delta g} \sim U(0, 1)$, and the diagnosed prevalence $(\pi\delta)_g = \pi_g \delta_g$ in each group follows.

4.3 Undiagnosed prevalence among GSM

Information about undiagnosed infections in GSM is obtained from GUMCAD (Genitourinary Medicine Clinic Activity Dataset) (Kirwan et al. 2016) a registry of attendance episodes in GUM clinics. HIV tests are offered routinely to previously-undiagnosed patients. Thus we have a sequence of observations g_i , representing firstly the number of GUM clinic visitors ($g_1 = 35121$) and then the number of patients with no previous HIV diagnosis ($g_2 = 34187$), HIV tests offered ($g_3 = 30570$), HIV tests accepted ($g_4 = 29529$), and HIV diagnoses made ($g_5 = 855$). For $i = 2, \dots, 5$, $g_i \sim \text{Bin}(g_{i-1}, \gamma_{i-1})$, with priors $\gamma_1, \gamma_2, \gamma_3 \sim U(0, 1)$ and $\gamma_4 \sim U(0, 0.15)$ (see below). An HIV infection may therefore remain undiagnosed if either a test is not offered or the patient opts out of testing. We can then decompose the prevalence of undiagnosed infection $\overline{(\pi\delta)}_G$ into “unoffered” $\pi^{(UN)}$ and “opt-out” $\pi^{(OP)}$ components.

$$\overline{(\pi\delta)}_G = \pi^{(UN)} + \pi^{(OP)}. \quad (11)$$

Both of those require strong prior assumptions to estimate, which will later be relaxed in a sensitivity analysis (§4.5). Firstly, the prevalence of infection that remains undiagnosed due to an unoffered test is

$$\pi^{(UN)} = \gamma_1(1 - \gamma_2)p^{(UN)}$$

where $\gamma_1(1 - \gamma_2)$ is the proportion of clinic attenders that are undiagnosed but not offered a test, and $p^{(UN)}$ is the probability that a test would be positive for these people. We assume the prevalence in this group is between 0.5 and 1.5 times the prevalence in people actually tested, and $\text{logit}(p^{(UN)}) = \text{logit}(\gamma_4) + a^{(UN)}$, with $a^{(UN)} \sim U(\log(0.5), \log(1.5))$

Secondly, the prevalence of infection remaining undiagnosed due to refusing a test is

$$\pi^{(OP)} = \gamma_1\gamma_2(1 - \gamma_3)(\gamma_4 + a^{(EX)})$$

$\gamma_1\gamma_2(1 - \gamma_3)$ is the proportion of clinic attenders that are undiagnosed and offered a test but opt out. We assume this group has an underlying HIV prevalence higher than those given tests, but not more than 15%, so that the excess prevalence in this group is $a^{(EX)} = a^{(OP)}(0.15 - \gamma_4)$, where $a^{(OP)} \sim U(0, 1)$, and the prior on γ_4 is truncated above at 0.15.

A small amount of additional evidence on $\overline{(\pi\delta)}_G$ is available from another dataset, GUM Anon (Public Health England, London 2012), a convenience survey of men not previously diagnosed with HIV who had attended a GUM clinic in the previous year. This gives direct information about the prevalence of HIV among previously undiagnosed GSM,

$$\pi^{(GA)} = (\overline{(\pi\delta)}_G + \pi^{(GD)})/\gamma_1, \quad (12)$$

where $\pi^{(GD)} = \prod_1^4 \gamma_r$ is the prevalence of newly-diagnosed infection among clinic attenders. The data in GUM Anon are $g^{(A)} \sim Bin(g^{(AN)}, \pi^{(GA)})$, where $g^{(A)} = 4$ and $g^{(AN)} = 85$.

4.4 Undiagnosed prevalence among NGMSM

To inform undiagnosed HIV prevalence in NGMSM, we use data from the Gay Men's Sexual Health Survey (GMSHS) (Aghaizu et al. 2016), based on face-to-face interviews in selected venues where participants were offered anonymous HIV tests. While this group is likely to have a higher HIV prevalence than the general population of MSM, it is assumed that the *relative odds* of having HIV between NGMSM and GSM is the same as in the general population. The GMSHS data provide the numbers $y_g^{(GM)}$ out of $n_g^{(GM)}$ previously-undiagnosed people in group g who tested positive for HIV (20 out of 493 GSM and 20 out of 452 NGMSM) so that $y_g^{(GM)} \sim Bin(n_g^{(GM)}, p_g^{(GM)})$, with $p_g^{(GM)} \sim U(0, 1)$. Defining the odds $o(p) = p/(1 - p)$, we apply the resulting odds ratio $or^{(GM)} = o(p_N^{(GM)})/o(p_G^{(GM)})$ to the baseline estimated from GUMCAD (§ 4.3), giving $o(\overline{(\pi\delta)}_N) = o(\overline{(\pi\delta)}_G)or^{(GM)}$.

4.5 Alternative assumptions

The results presented in section 5 are for the above model assumptions, unless specified otherwise. Two alternative assumptions are also explored.

- (a) **Undiagnosed prevalence from GUM Anon only** To avoid the strong prior assumptions on prevalence among those not offered a test or refusing a test, which are necessary to use the GUMCAD data to infer $\overline{(\pi\delta)}_g$, we could infer $\overline{(\pi\delta)}_g$ from GUM Anon alone. To construct this model, we replace equation (11) by a $U(0, 1)$ prior on $\overline{(\pi\delta)}_g$, although the GUMCAD data are still used to estimate the parameters $\pi^{(GD)}$ and γ_1 relating the prevalence in GUM Anon to $\overline{(\pi\delta)}_g$.
- (b) **GUMCAD also informs diagnosed prevalence** Instead of being inferred indirectly through the graph, the diagnosed prevalence can be modelled directly as

$$(\pi\delta)_G = (1 - \gamma_1) + \gamma_1\gamma_2\gamma_3\gamma_4, \quad (13)$$

where $1 - \gamma_1$ is the probability of a previous diagnosis, and $\gamma_1\gamma_2\gamma_3\gamma_4$ is the probability of a new diagnosis, in GUMCAD. This is not done in the base case due to concerns about inconsistencies in reporting between GUMCAD and SOPHID/HANDD.

5 Value of Information results in the HIV model

The model outputs of interest (as in Figures 1,3) are $\boldsymbol{\alpha} = ((\pi\delta)_G, (\pi\delta)_N, \overline{(\pi\delta)}_G, \overline{(\pi\delta)}_N, \mu_{DG}, \mu_{DN}, \mu_{UG}, \mu_{UN}, \mu)$, the diagnosed and undiagnosed prevalences among both GMSM and NGMSM, and the corresponding absolute numbers of people living with HIV (or “case-counts”), and the total number of MSM with HIV $\mu = \mu_{DG} + \mu_{DN} + \mu_{UG} + \mu_{UN}$. Samples from the posterior are generated using Hamiltonian Monte Carlo methods in the Stan

software (Stan Development Team 2016). These are illustrated in Figure 4 along with the overall prevalence $\pi_g = (\pi\delta)_g + \overline{(\pi\delta)}_g$ in each group g , and each of these quantities summed over the two groups g . The estimates of diagnosed prevalence in all MSM (top panel) are reasonably precise, while the corresponding estimates for NGMSM and GMSM are more uncertain. Estimates of undiagnosed prevalence are lower and more precise. Full results under the two alternative assumptions are presented in the supplementary figures.

5.1 Partial perfect information (EVPPI) for single outputs

Defining the decision problem as point estimation of α with quadratic loss, we use EVPPI formula (5) to determine which parameters ϕ contribute most to the uncertainty about each component of α , thus which ϕ may be worth learning more precisely. We will take ϕ to include the founder nodes of the graph illustrated in Figure 3. Since they are related to the α through a network of deterministic functions, perfect knowledge of these implies perfect knowledge of α . Each of the ϕ are either directly informed by data or given a substantive prior distribution based on belief. In the former case, EVPPI measures the maximum potential value of collecting more data from the same source. In the latter case, it will not necessarily be feasible to collect data to improve the precision of the belief, but EVPPI is still useful as a measure of how much of the uncertainty in α is explained by the uncertainty in the parameter.

The results are presented in Figure 5 as a grid whose r, s entry is colored according to $EVPPI_{\alpha_s}(\phi_r)/\text{var}(\alpha_s)$, the proportion of variance in α_s which would be reduced if we learnt ϕ_r . The lighter cells correspond to ϕ_r with greater EVPPI. Standard errors in these and all following EVPPI and EVSI estimates, arising from uncertainty in the coefficients of the regression (10), were negligible, at less than 1% of the EVPPI or EVSI estimates.

The parameters $a_{\delta G}$ and $a_{\delta N}$, governing the proportions of HIV infections that are diag-

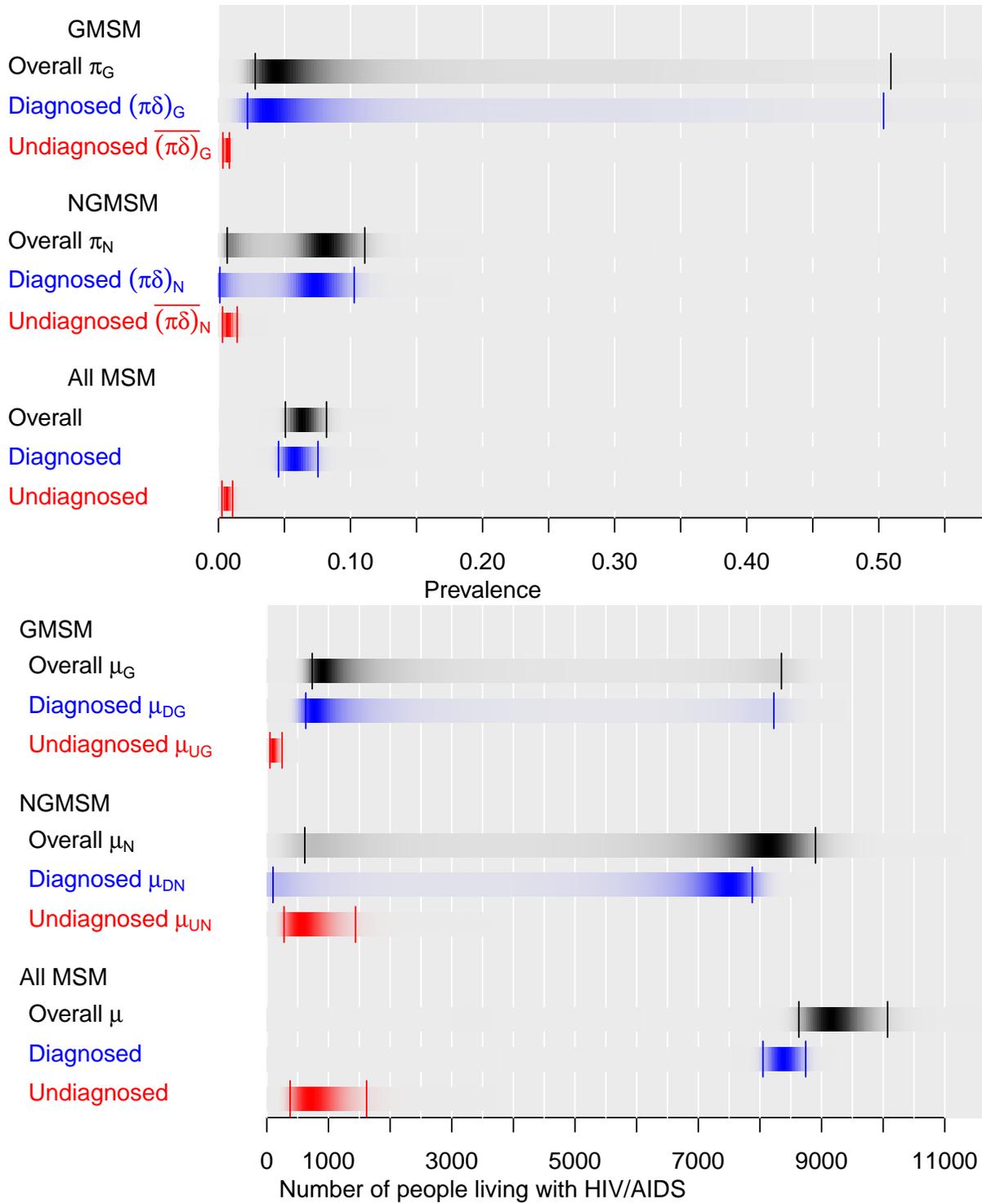


Figure 4: Posterior distributions of HIV prevalence (top) and numbers of MSM living with HIV/AIDS (bottom), London 2012. Darkness within each strip proportional to posterior density, with 95% credible intervals indicated.

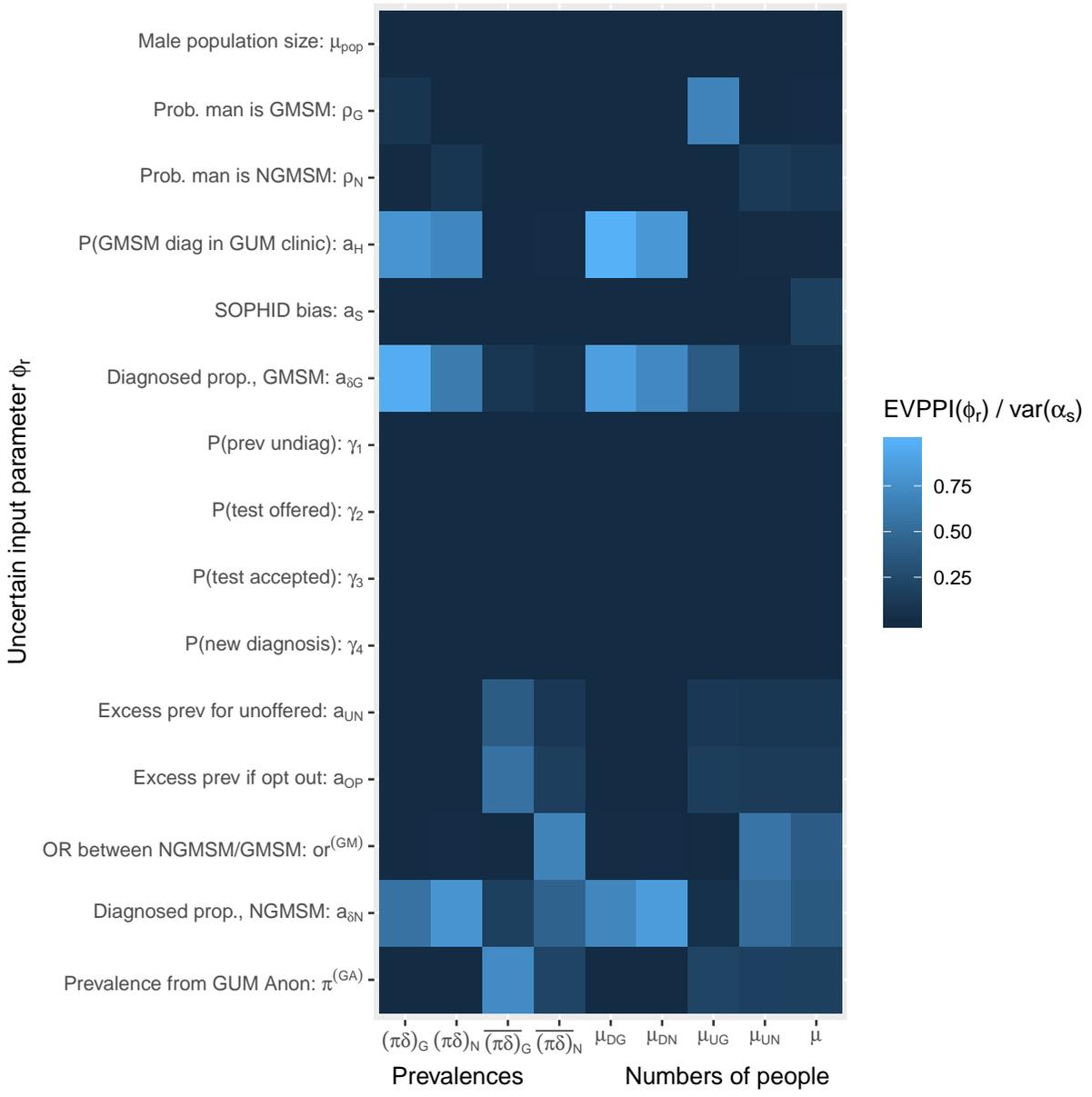


Figure 5: Expected value of partial perfect information in the HIV prevalence model.

nosed in each of the two groups, and the probability a_H that a GSM is newly diagnosed in a GUM clinic, explain most of the uncertainty in the *diagnosed* prevalences $(\pi\delta)_G, (\pi\delta)_N$ and the corresponding numbers of people diagnosed μ_{DG}, μ_{DN} . Direct data on any of these parameters would be difficult to obtain. However, if we were willing to make the assumption in (13), the estimates of diagnosed prevalence would become more precise, for example the posterior median (SD) of $(\pi\delta)_G$ would change from 0.06 (0.13) to 0.051 (0.001), though the extent of uncertainty around $(\pi\delta)_N, \mu_{DN}$ would not change substantively.

For the *undiagnosed* prevalences $\overline{(\pi\delta)}_G, \overline{(\pi\delta)}_N$ and undiagnosed case count μ_{UG} , Figure 5 shows that more GUM-Anon data (via $\pi^{(GA)}$), more GMSHS data (via $or^{(GM)}$) and more NATSAL data (via ρ_{UG}) respectively would give the greatest uncertainty reductions. These outcomes, however, are already precisely estimated in absolute terms (Figure 4). The number of NGMSM μ_{UN} with undiagnosed HIV is more uncertain, with 95% CrI (279,1442), and more GMSHS data would be potentially valuable to reduce this uncertainty.

If $\overline{(\pi\delta)}_G$ were informed only from the 4 infections out of 85 people observed in GUM Anon (alternative assumption (a)), the estimates of undiagnosed prevalence or case counts become extremely uncertain, for example, $\text{var}(\mu_{UN})$ increases from 302^2 to 2872^2 . We could reduce this uncertainty by collecting more GUM Anon data — since $EV PPI_{\mu_{UN}}(\pi^{(GA)})$ is $p = 62\%$ of $\text{var}(\mu_{UN})$, more GUM Anon data could reduce $\text{var}(\mu_{UN})$ to a minimum of $2872^2(1 - p) = 1777^2$ (note that the square root of the expected variance after learning data is not the same as the expected standard deviation).

5.2 Partial perfect information for multiple outputs

Staying with alternative assumption (a), suppose we wish to calculate the maximum potential value of extra GUM Anon data for *jointly* reducing the uncertainty about the number of GSM, NGMSM and PMSM with undiagnosed HIV, so that α is the vector

$(\mu_{UG}, \mu_{UN}, \mu_{UP})$. As described in § 2.3, we could simply calculate the standard EVPPI based on a scalar output α redefined as their sum, $\mu_U = \mu_{UG} + \mu_{UN} + \mu_{UP}$, the total number of MSM with undiagnosed HIV, whose posterior median is 5149 (SD 3280). This would ensure that any data expected to reduce the variance of any of these three outputs by the same (additive) amount would be valued equally. From this, we find that extra GUM Anon data would be expected to reduce $\text{var}(\mu_U)$ from 3280^2 to a minimum of 1801^2 . Since μ_U is dominated by NGMSM (posterior median of μ_{UN} is 4185), this is mostly explained by an expected reduction in $\text{var}(\mu_{UN})$ from 2864^2 to a minimum of 1770^2 .

Alternatively, suppose both the prevalences and the case counts are of interest, for example in NGMSM, so that $\boldsymbol{\alpha} = ((\overline{\pi\delta})_N, \mu_{UN})$. Since these two components are on very different scales, the Bayesian “D-optimality” criterion $v(\boldsymbol{\alpha}) = \det(\text{cov}(\boldsymbol{\alpha}))$ would be a preferable measure of overall expected loss due to uncertainty. We use this criterion to compare the maximum expected value of extra GUM Anon data and extra GMSHS data, which combine to estimate the outcomes for NGMSM as described in § 4.4. The EVPPI is interpreted as the expected reduction in the product of $\text{var}((\overline{\pi\delta})_N)$ and $\text{var}(\mu_{UN})$ given by extra GUM Anon or GMSHS data, adjusted for their covariance. This is 425 and 135 respectively, favouring extra data from GUM Anon. Though in this example, examining expected reductions in $\text{var}((\overline{\pi\delta})_N)$ or $\text{var}(\mu_{UN})$ separately would lead to the same conclusion, since $(\overline{\pi\delta})_N$ is defined as the proportion μ_{UN}/r_N of NGMSM with HIV, and GUM Anon and GMSHS are not informative about the number r_N of NGMSM, thus extra data informs μ_{UN} entirely through information on $(\overline{\pi\delta})_N$ (or vice versa).

5.3 Sample information (EVSI)

We now estimate the expected value of data with specific sample sizes for improving the precision of the estimated number of people μ_U with undiagnosed HIV. Using the GUMCAD

data and associated strong prior assumptions, the posterior median of μ_U is 804 (SD 320), compared to 5149 (SD 3280) with this information excluded (a). We compare the value of additional data from GUM Anon and additional data from GMSHS (on top of their original sample sizes of 85 and 945 respectively) for reducing these posterior standard deviations.

The EVSI is computed for a series of sample sizes n using the method in § 3.2. For GUM Anon (§ 4.3), the sufficient statistic $T(\mathbf{y})$ consists of the empirical HIV prevalence \mathbf{y}/n from an additional survey $\mathbf{y} \sim Bin(n, \pi^{(GA)})$. For GMSHS (§ 4.4), given a sample size n , $\mathbf{y} = (N_G^{(GM)}, Y_G^{(GM)}, Y_N^{(GM)})$, where $N_G^{(GM)}$ is the number of previously-undiagnosed MSM in the future sample of n who attend GUM clinics (the equivalent of the observed $n_G^{(GM)} = 493$). Then $Y_G^{(GM)}$ and $Y_N^{(GM)}$ are the numbers of men out of denominators $N_G^{(GM)}$ and $N_N^{(GM)} = n - N_G^{(GM)}$ (GMSM and NGMSM respectively) who test positive for HIV, the equivalents of the observed $y_G^{(GM)} = 20, y_N^{(GM)} = 492$. We take $T(\mathbf{y}) = o(\hat{p}_N^{(GM)}(\mathbf{y}))/o(\hat{p}_G^{(GM)}(\mathbf{y}))$, a point estimator of the odds ratio, where $\hat{p}_G^{(GM)}(\mathbf{y})$ is an estimator of the proportion of MSM in group g who have HIV. To avoid zeros in the denominator $o(\hat{p}_G^{(GM)}(\mathbf{y}))$, we use a Bayesian estimator $\hat{p}_G^{(GM)}(\mathbf{y}) = (Y_G^{(GM)} + 0.5)/(N_G^{(GM)} + 1)$, the posterior mean of a binomial proportion under a Jeffreys Beta(0.5,0.5) prior, rather than the empirical proportion $Y_G^{(GM)}/N_G^{(GM)}$.

Figure 6 shows $\text{var}(\mu_U) - EVSI(\mathbf{y})$, the expected variance remaining after data collection, under the two alternative assumptions. With the strong priors, μ_U is relatively well informed, and extra data from GUM Anon at realistic sample sizes (1000 or less) would not noticeably reduce $\text{var}(\mu_U)$. GMSHS data would be more valuable, through improving the estimate of μ_{UN} , the more uncertain contributor to $\mu_U = \mu_{UG} + \mu_{UN}$. 1000 extra observations from GMSHS would be expected to reduce $\text{var}(\mu_U)$ from 320^2 to 279^2 .

Without the strong prior information, $\text{var}(\mu_U) = 3280^2$ is substantially greater, and μ_U is only directly informed by the 85 observations from GUM Anon. Extra data from

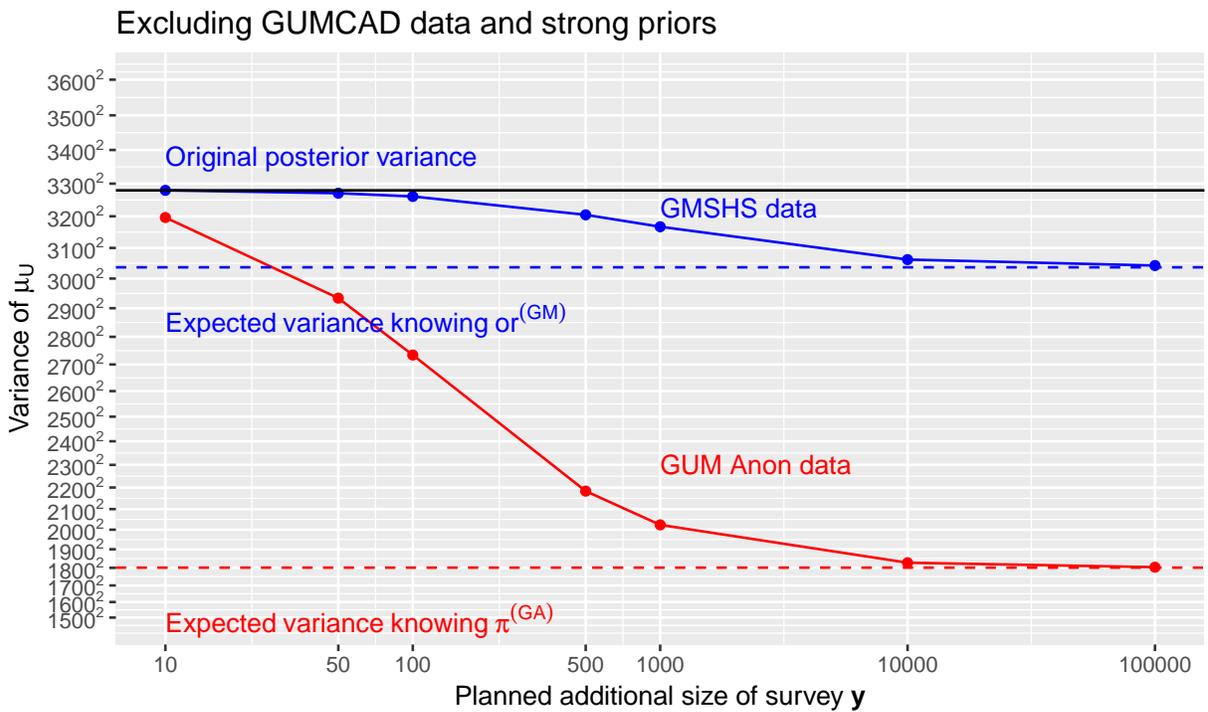
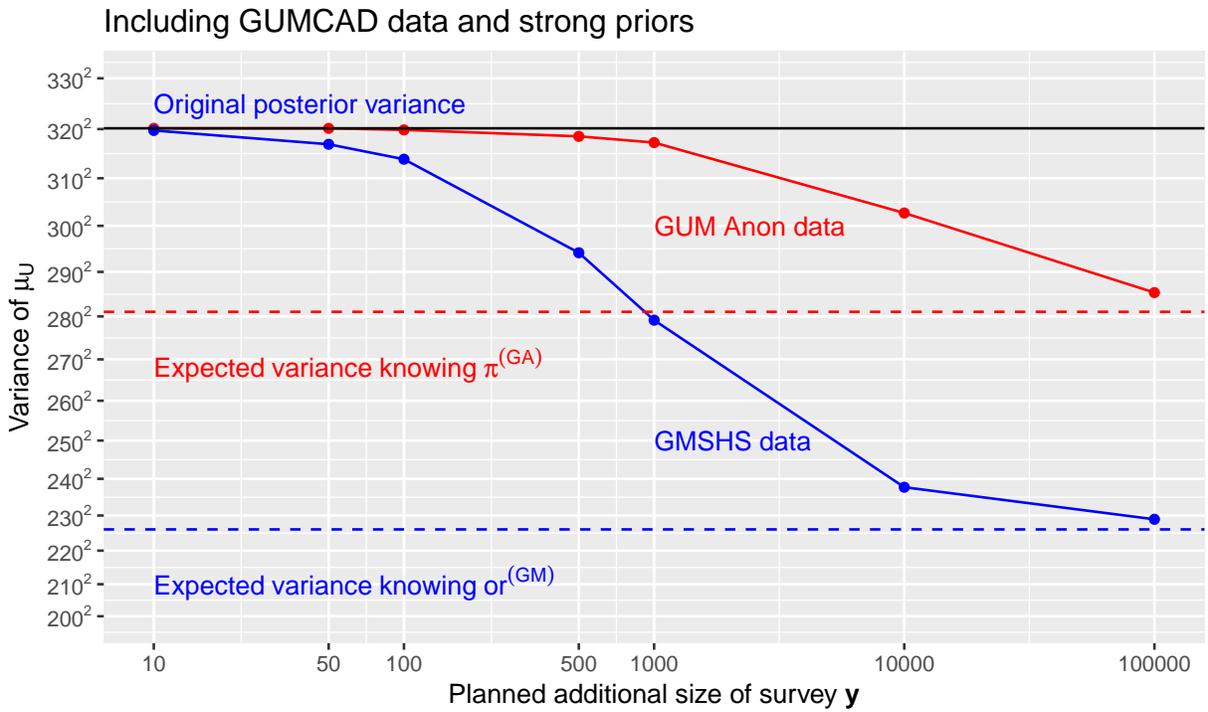


Figure 6: Expected value of sample information: value of additional data from GUM Anon or GMSHS for reducing the variance of the total number of MSM with undiagnosed HIV, $\mu_U = \mu_{UG} + \mu_{UN}$. The x axis is on the log scale. The y axis is the variance, with the labels as SD^2 .

this source would be valuable, for example, another 500 observations would be expected to reduce this variance to 2184². Relative to these improvements, GMSHS data of the same size would be much less valuable. GMSHS data however would be expected to give around the same *absolute* reductions in $\text{var}(\mu_U)$, whether or not the strong priors are included.

5.4 Net benefit of sampling

The benefits from improved precision of estimates of μ_U must be traded off with the costs of data collection, to determine an optimal sample size for extra survey data. Consider the scenario which excluded the GUMCAD data and associated strong priors. In the GUM Anon survey there was a cost of around £17 per participant, which is assumed to be the same for collecting further data from this source. The cost $c(\mathbf{y})$ is illustrated against sample sizes of \mathbf{y} from 1 to 400 by the straight line in Figure 7. Suppose also that the decision-maker is willing to pay £5,000 to reduce the variance of μ_U by $dv = 3271^2 - 2771^2$, which in this case would reduce the standard deviation by 500, from 3271 to 2771. The willingness to pay per unit variance is then $\lambda = 5000/dv$.

Collecting extra data \mathbf{y} will give an expected reduction in $\text{var}(\mu_U)$ of $EVSI(\mathbf{y})$, as illustrated in Figure 7. The resulting expected (monetary) benefit $E_{\mathbf{y}}(b(\mathbf{y}))$ (§2.4) is shown to be a nonlinear function of the sample size of \mathbf{y} , with an asymptote representing the expected value of partial perfect information on $\pi^{(GA)}$. Hence the expected *net benefit* of sampling $E_{\mathbf{y}}(b(\mathbf{y}) - c(\mathbf{y}))$ is illustrated in the bottom panel of Figure 7. The expected benefits of sampling always exceed the costs, and the net benefit is maximised at a sample size of 166. Also illustrated are the benefit and net benefit that would result if the decision maker was willing to pay twice or half the original amount, 2λ or $\lambda/2$. The corresponding optimal sample sizes would be 315 or 81 respectively.

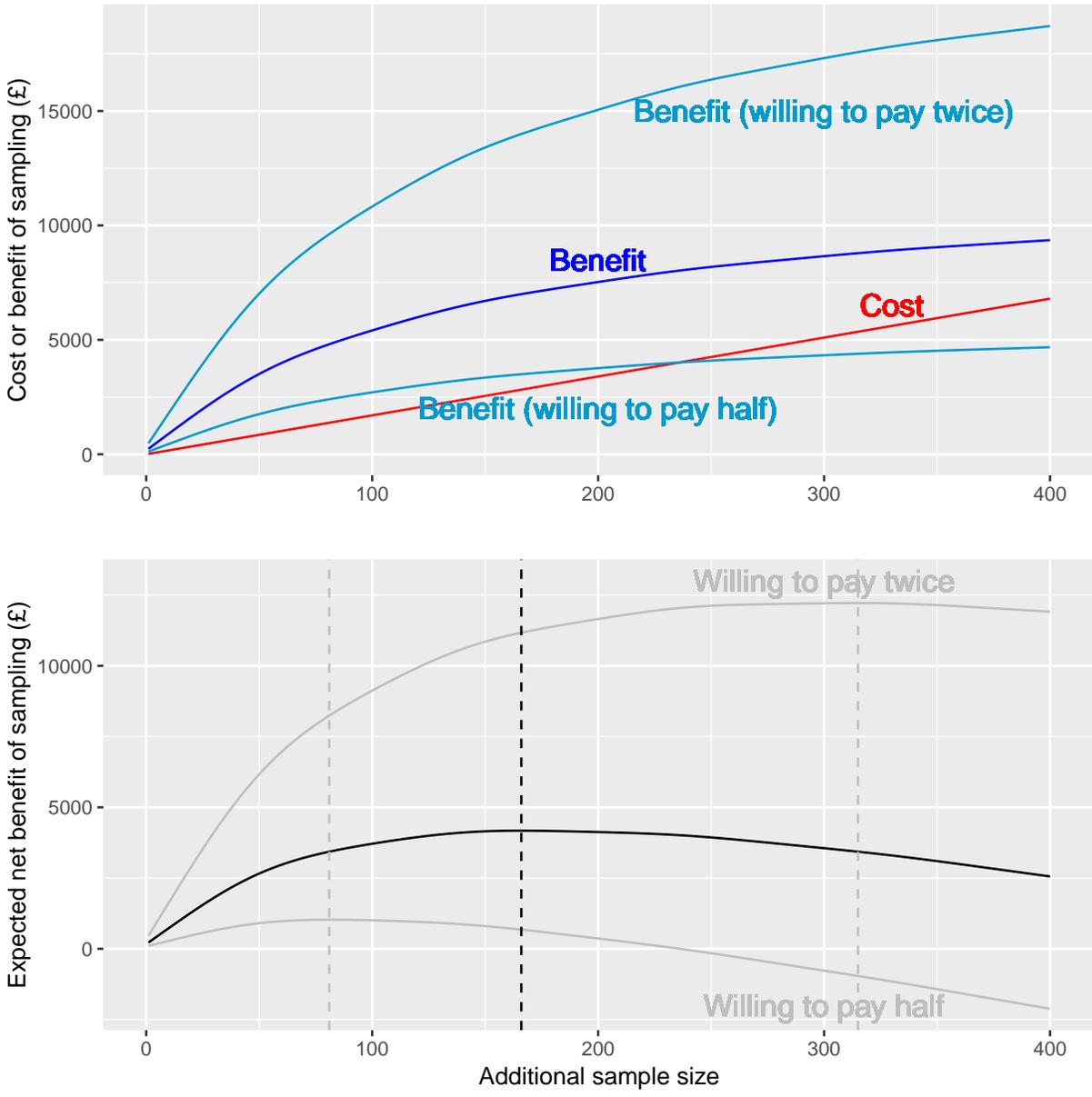


Figure 7: Expected cost, benefit and net benefit of sampling up to 400 extra participants from GUM Anon, if we wish to reduce the variance of μ_U , the number of MSM with undiagnosed HIV. The optimal sample size is illustrated as a dotted line.

6 Summary and potential further work

We have presented tools to find the most influential sources of uncertainty in MPES models and determine the expected value of extra data. We generalized methods, previously applied only in deterministic models, to complex graphical models, a class which also includes hierarchical models. We have shown how VoI methods developed for formal finite-choice decision problems can be extended to deal with estimation of single or multiple quantities.

While the purpose of our model was to estimate a quantity of interest to policy-makers, the same methods could be used for models to compare specific health policies. Sections §2.4 and §5.4 illustrated how the benefits from more precise estimates of HIV prevalence could be converted directly into a monetary value. An alternative approach would be to value the indirect health gains that would result from better data, through better-informed health policies. This would allow standard health economic principles to be used (see, e.g. Briggs et al. 2006). For example, the National Institute for Health and Care Excellence in the UK recommends that a new healthcare intervention is funded by the National Health Service if the cost per quality-adjusted life year (QALY) gained, compared to current practice, is less than around $\lambda = \pounds 20,000$, implying a willingness to pay of λ per QALY. This is a choice between two actions $d \in \{\text{accept, reject}\}$ intervention” (as in §2.3) with loss $L(d, \theta) = c_d(\theta) - \lambda q_d(\theta)$, where $q_d(\theta)$ is the expected QALY and $c_d(\theta)$ is the expected cost for a person under action d , from a health economic model with parameters θ . See, for example, Carmona et al. (2016), Baggaley et al. (2017), for how such models might be built for HIV testing interventions to increase the proportion δ of people who are diagnosed. Briefly, any QALY gains strongly depend on the underlying prevalence of HIV among the population receiving the intervention. Thus, improved estimates of prevalence will lead to more precise estimates of the QALY gains, and a better-informed decision about whether

to implement the intervention, which may result in a better use of health service resources. VoI methods may then be used to decide whether further information should be collected to support the decision.

In the HIV application, we found that structural assumptions, such as whether to include a particular piece of information, were influential to both the parameter estimates and the value of information. Such uncertainties might be parameterised (see, e.g. Strong et al. 2012), for example a particular prior or dataset of uncertain relevance could be discounted using an unknown weight (e.g. Neuenschwander et al. 2009). The EVPPI of the extra parameter would then quantify this uncertainty in the context of all other uncertainties, referred to as the “expected value of model improvement” by Strong & Oakley (2014).

Note that VoI refers to the expected value of *potential future* information, which differs from the *observed* value of a dataset x_i *currently* included in the model. The latter could be computed as the observed reduction in loss when the model is refitted without x_i . This could demonstrate the value of past data to the policymaker responsible for funding the collection of future data of the same type. For surveys or longitudinal studies conducted at regular intervals, VoI might be used to determine the expected value of future surveys or follow-up, although a full analysis would require modelling the expected changes through time in the quantities, such as disease prevalence or incidence, informed by the data.

While our method is broadly applicable, the details of computation for different decision problems and loss functions may be different. We discussed finite-action decisions and point estimation. A more general decision problem is to estimate the entire uncertainty distribution of $\boldsymbol{\theta}$. The standard posterior $p(\boldsymbol{\theta}|\mathbf{y})$ is then optimal under a log scoring rule (Bernardo & Smith 1994), and (following Lindley 1956) standard Bayesian design theory aims to maximise the information gain from new data \mathbf{y} , which we can write as $EVS I(\mathbf{y}) = E_{\boldsymbol{\theta}}(-\log(p(\boldsymbol{\theta}))) + E_{\mathbf{y}}E_{\boldsymbol{\theta}|\mathbf{y}}\{\log(p(\boldsymbol{\theta}|\mathbf{y}))\}$. Under linear models (Chaloner &

Verdinelli 1995), this is equivalent to minimising $\det(\text{cov}(\boldsymbol{\theta}))$, but more generally this is challenging to compute (Ryan et al. 2016).

Note that the VoI approach to sensitivity analysis is an example of the “global” approach, which examines the changes in model outputs given by varying parameters within the ranges of their belief distributions. The “local” approach is based on examining the posterior geometry resulting from small parameter perturbations around a base case, e.g. Roos et al. (2015) assess the robustness of hierarchical models to prior assumptions in this way. While the global approach is easier to interpret, as discussed by Oakley & O’Hagan (2004) and Roos et al. (2015), it conditions on one particular prior specification, and parameterising all potential prior beliefs or structural assumptions would be impractical.

The regression method for VoI computation that we described requires only a MCMC sample from the joint distribution of parameters of interest $\boldsymbol{\phi}$ and outputs $\boldsymbol{\alpha}$. Additionally for EVSI it requires that the information in the new data \mathbf{y} can be condensed into an analytic sufficient statistic $T(\mathbf{y})$. Alternative methods which exploit particular analytic structures of $g(\cdot)$, where α is a known function $g(\boldsymbol{\phi})$, thus avoiding a regression approximation, were discussed by Madan et al. (2014) for EVPPI and and Ades et al. (2004) for EVSI. Menzies (2016) also presented an importance resampling method for EVSI computation which needs only a single MCMC sample and not a sufficient statistic.

In conclusion, the consideration of future evidence requirements is an often-neglected part of statistical analysis. The Value of Information methods we have presented provide a practicable set of tools for achieving this aim in the context of Bayesian evidence synthesis.

References

- Ades, A. E. & Sutton, A. J. (2006), ‘Multiparameter evidence synthesis in epidemiology and medical decision-making: current approaches’, *Journal of the Royal Statistical Society, Series A* **169**(1), 5–35.
- Ades, A., Lu, G. & Claxton, K. (2004), ‘Expected value of sample information calculations in medical decision modeling’, *Medical Decision Making* **24**(2), 207–227.
- Aghaizu, A., Wayal, S., Nardone, A., Parsons, V., Copas, A., Mercey, D., Hart, G., Gilson, R. & Johnson, A. (2016), ‘Sexual behaviours, HIV testing, and the proportion of men at risk of transmitting and acquiring HIV in London, UK, 2000-13: a serial cross-sectional study’, *The Lancet HIV* **3**(9), e431–e440.
- Baggaley, R. F., Irvine, M. A., Leber, W., Cambiano, V., Figueroa, J., McMullen, H., Anderson, J., Santos, A. C., Terris-Prestholt, F., Miners, A., Hollingsworth, D. & Griffiths, C. J. (2017), ‘Cost-effectiveness of screening for HIV in primary care: a health economics modelling analysis’, *The Lancet HIV* **4**(10), e465–e474.
- Berger, J. O. (2013), *Statistical decision theory and Bayesian analysis*, Springer.
- Bernardo, J. M. & Smith, A. F. M. (1994), *Bayesian Theory*, Wiley, Chichester.
- Briggs, A., Sculpher, M. & Claxton, K. (2006), *Decision modelling for health economic evaluation*, Handbooks in Health Economic Evaluation, Oxford University Press, Oxford.
- Carmona, C., O’Rourke, D. & Robinson, S. (2016), ‘HIV testing: increasing uptake among people who may have undiagnosed HIV. Evidence review on the most cost effective ways to increase the uptake of HIV testing to reduce undiagnosed HIV among people

who may have been exposed to it’.

URL: <https://www.nice.org.uk/guidance/ng60/documents/evidence-review-5>

Chaloner, K. & Verdinelli, I. (1995), ‘Bayesian experimental design: A review’, *Statistical Science* pp. 273–304.

Claxton, K. P. & Sculpher, M. J. (2006), ‘Using value of information analysis to prioritise health research’, *Pharmacoeconomics* **24**(11), 1055–1068.

De Angelis, D., Presanis, A. M., Conti, S. & Ades, A. E. (2014), ‘Estimation of HIV burden through Bayesian evidence synthesis’, *Statistical Science* **29**(1), 9–17.

Felli, J. C. & Hazen, G. B. (1998), ‘Sensitivity analysis and the expected value of perfect information’, *Medical Decision Making* **18**(1), 95–109.

Friedman, J. H. (1991), ‘Multivariate adaptive regression splines’, *The Annals of Statistics* pp. 1–67.

Goubar, A., Ades, A. E., DeAngelis, D., McGarrigle, C. A., Mercer, C. H., Tookey, P. A., Fenton, K. & Gill, O. N. (2008), ‘Estimates of human immunodeficiency virus prevalence and proportion diagnosed based on Bayesian multiparameter synthesis of surveillance data’, *Journal of the Royal Statistical Society: Series A (Statistics in Society)* **171**(3), 541–580.

Han, C. & Chaloner, K. (2004), ‘Bayesian experimental design for nonlinear mixed-effects models with application to HIV dynamics’, *Biometrics* **60**(1), 25–33.

- Heath, A., Manolopoulou, I. & Baio, G. (2016), ‘Estimating the expected value of partial perfect information in health economic evaluations using integrated nested Laplace approximation’, *Statistics in Medicine* **35**(23), 4264–4280.
- Kirwan, P., Chau, C., Brown, A., Gill, O., Delpech, V. & contributors (2016), HIV in the UK — 2016 report, Technical report, Public Health England, London.
- Lamboni, M., Monod, H. & Makowski, D. (2011), ‘Multivariate sensitivity analysis to measure global contribution of input factors in dynamic models’, *Reliability Engineering & System Safety* **96**(4), 450–459.
- Lauritzen, S. L. (1996), *Graphical models*, Vol. 17, Clarendon Press.
- Lindley, D. V. (1956), ‘On a measure of the information provided by an experiment’, *The Annals of Mathematical Statistics* pp. 986–1005.
- Madan, J., Ades, A. E., Price, M., Maitland, K., Jemutai, J., Revill, P. & Welton, N. J. (2014), ‘Strategies for efficient computation of the expected value of partial perfect information’, *Medical Decision Making* **34**(3), 327–342.
- Mandel, M. (2013), ‘Simulation-based confidence intervals for functions with complicated derivatives’, *The American Statistician* **67**(2), 76–81.
- Menzies, N. A. (2016), ‘An efficient estimator for the expected value of sample information’, *Medical Decision Making* **36**(3), 308–320.
- Mercer, C., Tanton, C., Prah, P., Erens, B., Sonnenberg, P., Clifton, S., Macdowall, W., Lewis, R., Field, N., Datta, J., Copas, A., Phelps, A., Wellings, K. & Johnson, A. (2013), ‘Changes in sexual attitudes and lifestyles in Britain through the life course and over time:

findings from the National Surveys of Sexual Attitudes and Lifestyles (Natsal)', *Lancet* **382**(9907), 1781–1794.

Milborrow, S. (2011), *earth: Multivariate Adaptive Regression Splines*. R package. Derived from mda:mars by T. Hastie and R. Tibshirani.

URL: *<http://CRAN.R-project.org/package=earth>*

Neuenschwander, B., Branson, M. & Spiegelhalter, D. J. (2009), 'A note on the power prior', *Statistics in Medicine* **28**(28), 3562–3566.

Oakley, J. E. & O'Hagan, A. (2004), 'Probabilistic sensitivity analysis of complex models: a Bayesian approach', *Journal of the Royal Statistical Society, Series B* **66**(3), 751–769.

Office for National Statistics (2012), 'Mid-year population estimates'.

URL: *<https://www.ons.gov.uk/peoplepopulationandcommunity/populationandmigration/populationestimates>*

Parmigiani, G. & Inoue, L. (2009), *Decision theory: principles and approaches*, Wiley.

Presanis, A. M., Gill, O. N., Chadborn, T. R., Hill, C., Hope, V., Logan, L., Rice, B. D., Delpech, V. C., Ades, A. E. & De Angelis, D. (2010), 'Insights into the rise in HIV infections, 2001 to 2008: a Bayesian synthesis of prevalence evidence.', *AIDS (London, England)* **24**(18), 2849–58.

Public Health England, London (2012), 'UA survey of genitourinary medicine (GUM) clinic attendees (GUM anon survey)'.

URL: *<https://www.gov.uk/guidance/hiv-overall-prevalence\#ua-survey-of-genitourinary-medicine-gum-clinic-attendees-gum-anon-survey>*

Raiffa, H. & Schlaifer, R. (1961), *Applied statistical decision theory*, Harvard University.

- Roos, M., Martins, T. G., Held, L. & Rue, H. (2015), ‘Sensitivity analysis for Bayesian hierarchical models’, *Bayesian Analysis* **10**(2), 321–349.
- Ryan, E. G., Drovandi, C. C., McGree, J. M. & Pettitt, A. N. (2016), ‘A review of modern computational algorithms for Bayesian optimal design’, *International Statistical Review* **84**(1), 128–154.
- Saltelli, A., Tarantola, S., Campolongo, F. & Ratto, M. (2004), *Sensitivity analysis in practice: a guide to assessing scientific models*, John Wiley & Sons.
- Sobol’, I. M. (2001), ‘Global sensitivity indices for nonlinear mathematical models and their Monte Carlo estimates’, *Mathematics and Computers in Simulation* **55**(1), 271–280.
- Stan Development Team (2016), *Stan Modeling Language Users Guide and Reference Manual, Version 2.14.0*.
URL: <http://mc-stan.org>
- Strong, M., Oakley, J. & Chilcott, J. (2012), ‘Managing structural uncertainty in health economic decision models: a discrepancy approach’, *Journal of the Royal Statistical Society: Series C (Applied Statistics)* **61**(1), 25–45.
- Strong, M. & Oakley, J. E. (2014), ‘When is a model good enough? deriving the expected value of model improvement via specifying internal model discrepancies’, *SIAM/ASA Journal on Uncertainty Quantification* **2**(1), 106–125.
- Strong, M., Oakley, J. E. & Brennan, A. (2014), ‘Estimating multiparameter partial expected value of perfect information from a probabilistic sensitivity analysis sample: A nonparametric regression approach’, *Medical Decision Making* **34**(3), 311–326.

- Strong, M., Oakley, J. E., Brennan, A. & Breeze, P. (2015), 'Estimating the expected value of sample information using the probabilistic sensitivity analysis sample: A fast, nonparametric regression-based method', *Medical Decision Making* **35**(5), 570–83.
- Welton, N., Ades, A., Caldwell, D. & Peters, T. (2008), 'Research prioritization based on expected value of partial perfect information: a case-study on interventions to increase uptake of breast cancer screening', *Journal of the Royal Statistical Society: Series A (Statistics in Society)* **171**(4), 807–841.
- Willan, A. R. & Pinto, E. M. (2005), 'The value of information and optimal clinical trial design', *Statistics in Medicine* **24**(12), 1791–1806.